

ГОДИШНИК

НА

СОФИЙСКИЯ УНИВЕРСИТЕТ
„СВ. КЛИМЕНТ ОХРИДСКИ“

ФАКУЛТЕТ
ПО МАТЕМАТИКА И ИНФОРМАТИКА

Том 107
2020

ANNUAL

OF

SOFIA UNIVERSITY
“ST. KLIMENT OHRIDSKI”

FACULTY OF MATHEMATICS AND INFORMATICS

Volume 107
2020

СОФИЯ • 2020 • SOFIA

УНИВЕРСИТЕТСКО ИЗДАТЕЛСТВО „СВ. КЛИМЕНТ ОХРИДСКИ“
“ST. KLIMENT OHRIDSKI” UNIVERSITY PRESS

Annual of Sofia University “St. Kliment Ohridski”
Faculty of Mathematics and Informatics

Годишник на Софийския университет „Св. Климент Охридски”

Факултет по математика и информатика

Managing Editors: Geno Nikolov (Mathematics)
Krassen Stefanov (Informatics)

Editorial Board

Danyo Danev Darina Dicheva Ognyan Hristov Stefan Ivanov
Azniv Kasparian Vladimir Kostov Mikhail Krastanov Alexandra Soskova

Address for correspondence:

Faculty of Mathematics and Informatics
“St. Kliment Ohridski” University of Sofia Fax xx(359 2) 8687 180
5, J. Bourchier Blvd., P.O. Box 48 Electronic mail:
BG-1164 Sofia, Bulgaria annuaire@fmi.uni-sofia.bg

Aims and Scope. The *Annual* is the oldest Bulgarian journal, founded in 1904, devoted to pure and applied mathematics, mechanics and computer science. It is reviewed by *Zentralblatt für Mathematik*, *Mathematical Reviews* and the Russian *Referativnii Jurnal*. The *Annual* publishes significant and original research papers of authors both from Bulgaria and abroad in some selected areas that comply with the traditional scientific interests of the Faculty of Mathematics and Informatics at the “St. Kliment Ohridski” University of Sofia, i.e., algebra, geometry and topology, analysis, probability and statistics, mathematical logic, theory of approximations, numerical methods, computer science, classical, fluid and solid mechanics, and their fundamental applications.

© “St. Kliment Ohridski” University of Sofia
Faculty of Mathematics and Informatics
2020
ISSN 1313–9215 (Print)
ISSN 2603–5529 (Online)

CONTENTS

In Memoriam: Rumen Maleev (1943–2019)	3
STANIMIR TROYANSKI. Разнопосочни спомени за Румен Малеев – приятеля и колегата	7
ZHIVKO H. PETROV AND TATYANA L. TODOROVA. Representation of natural numbers by sum of four squares of almost-prime having a special form	13
BOGDANA A. GEORGIEVA. Review of continuum mechanics and its history. Part I: Deformation and stress. Conservation laws. Constitutive equa- tions	29
BOGDANA A. GEORGIEVA. Review of continuum mechanics and its history. Part II: The mechanics of thermoelastic media. Perfect fluids. Linearly viscous fluids	55
VLADIMIR PETROV KOSTOV. Univariate polynomials and the contractability of certain sets	79
NIKOLAY A. IVANOV. Examples of HNN–extensions with nontrivial quasi- kernels	105
GENO NIKOLOV, BORISLAVA PETROVA. On the regularity of certain three- row almost Hermitian incidence matrices	129

ГОДИШНИК НА СОФИЙСКИЯ УНИВЕРСИТЕТ СВ. КЛИМЕНТ ОХРИДСКИ“

ФАКУЛТЕТ ПО МАТЕМАТИКА И ИНФОРМАТИКА

Том 107

ANNUAL OF SOFIA UNIVERSITY ST. KLIMENT OHRIDSKI“

FACULTY OF MATHEMATICS AND INFORMATICS

Volume 107

In Memoriam

Rumen Maleev (1943–2019)



Professor Rumen Maleev passed away on December 16, 2019. With more than forty years service in Sofia University, he will be remembered by his colleagues, friends and students as one of the most distinguished Professors in the Faculty of Mathematics and Informatics. We present here a short CV of Rumen Maleev, followed by a paper (in Bulgarian) of Academician Stanimir Troyanski, who shares some memories about his long-term friendship and collaboration with Rumen Maleev.

CURRICULUM VITAE

of

Roumen Maleev

Date and place of birth:

August 17, 1943, Samokov, Bulgaria

Spoken foreign languages:

English, Russian, Romanian, German, French

Education, Degrees:

- Ms.C., University of Bucharest, Bucharest, Romania, 1967
- Ph.D., Sofia University, Sofia, Bulgaria, 1975
- D.Sc., Sofia University, Sofia, Bulgaria, 1996

Specializations:

- Moscow State University, Moscow, Russia, Academic year 1971/72
- Warsaw University, Warsaw, Poland, February – April, 1982

Professional experience:

- Researcher, Institute of Mathematics and Informatics, Bulgarian Academy of Sciences, Sofia, Bulgaria, 1967 - 1969
- Assistant Professor, Faculty of Mathematics and Informatics, Sofia University, Sofia, Bulgaria, 1970 - 1983
- Associate Professor, Faculty of Mathematics and Informatics, Sofia University, Sofia, Bulgaria, 1983 - 2006
- Full Professor, Faculty of Mathematics and Informatics, Sofia University, Sofia, Bulgaria, 2006 – 2011

Visiting positions:

- South Florida University, Tampa, Florida, USA, Spring semester 1991
- Athens University, Athens, Greece May – June, 1997

Administrative positions:

- Deputy Dean of the Faculty of Mathematics and Informatics, Sofia University, Sofia, Bulgaria, 1989 – 1995

- Head of the Department of Mathematical Analysis of the Faculty of Mathematics and Informatics, Sofia University, Sofia, Bulgaria, 1998 – 2000
- Member of the Specialized Scientific Council on Mathematics and Mechanics, 1995 – 2004
- Vice President of the Specialized Scientific Council on Mathematics and Mechanics, 1998 – 2004
- Secretary of the Scientific Commission for degrees and positions in Mathematics and Mechanics, 2003 – 2006
- Vice President of the Scientific Commission for degrees and positions in Mathematics and Mechanics, 2006 – 2009
Member of the Scientific Commission on Mathematics and Mechanics of the National Scientific Fund, 2010 –2012
- Elections expert in missions of Organization for Security and Cooperation in Europe(OSCE), 1997 – 2012

Research interests:

Geometry of Banach spaces, Approximation theory, Numerical analysis

List of selected publications

1. R. Maleev, An iterative method for equations with monotonic operators, *USSR Comput. Math. Math. Phys.*, **13** (1973), 280–286.
2. R. Maleev, S. Troyanski, The moduli of convexity and smoothness of the spaces $L_{p,q}$, *Annuare Univ. Sofia Math.*, **66** (1974), 331–339, (Russian).
3. R. Maleev, S. Troyanski, Unconditionally convergent and absolutely divergent series in Orlicz spaces, *C. R. Acad. Bulg. Sci.*, **27** (1974), 1029–1032, (Russian).
4. R. Maleev, S. Troyanski, On the moduli of convexity and smoothness in Orlicz spaces, *Studia Math.*, **54** (1975), 131–141.
5. R. Maleev, On conditionally convergent series in Orlicz spaces L_M , *Serdica*, **1** (1975), 178–182, (Russian).
6. R. Maleev, S. Markov, D. Vandev, Least square approximations using Hausdorff metric, in: *Mathematics and Education in Mathematics, 5-th Spring Conf. Bulg. Math. Union, Gabrovo, April 1975*, Publ. House of BAS, Sofia, 1990.
7. R. Maleev, On conditionally convergent series in Banach lattice, *C. R. Acad. Bulg. Sci.*, **32** (1979), 1015–1018.
8. A. Andreev, R. Maleev, Error estimation of the finite element method for one dimensional finite problems, *Serdica*, **6** (1980), 278–283.

9. R. Maleev, S. Troyanski, Order moduli of convexity and smoothness, *Funct. Anal. Appl.*, **17** (1983), 231–233.
10. R. Maleev, S. Troyanski, On cotypes of Banach lattices, in: *Constructive Function Theory'81*, BAS, Sofia 1983, 429–441.
11. R. Maleev, S. Troyanski, Smooth functions in Orlicz spaces, *Contemporary Math.*, **85** (1989), 355–370.
12. R. Maleev, Korovkin Theorem in rearrangement invariant function spaces, in: *Constructive Function Theory'84*, BAS, Sofia 1984, 578–582.
13. R. Maleev, S. Troyanski, Smooth norms in Orlicz spaces, *Canad. Math. Bull.*, **34** (1991), 74–82.
14. R. Maleev, Norms of best smoothness in Orlicz spaces, *Zeitschrift Anal. Anwendungen*, **12** (1993), 123–135.
15. R. Maleev, Higher order uniformly Gâteaux differentiable norms in Orlicz spaces, *Rocky Mount. Math. J.*, **25** (1995), 1117–1136.
16. R. Gonzalo, R. Maleev, Smooth functions in Orlicz function spaces, *Arch. Math.*, **69** (1997), 136–145.
17. R. Maleev, G. Nedev, B. Zlatanov, Gâteaux differentiability of bump functions in Banach spaces, *J. Math. Anal. Appl.*, **240** (1999), 311–323.
18. R. Maleev, B. Zlatanov, Cotype of weighted Orlicz sequence spaces, *C. R. Bulg. Acad. Sci.*, **53**, no. 3 (2000), 9–12.
19. R. Maleev, B. Zlatanov, Smoothness in Musielak–Orlicz sequence spaces, *C. R. Bulg. Acad. Sci.*, **55**, no. 6 (2002), 11–16.
20. R. Maleev, B. Zlatanov, Gâteaux differentiable norms in Musielak–Orlicz spaces, *Math. Balk.*, **20** (2006), 299–313.

Textbooks:

1. Bl. Sendov, R. Maleev, S. Markov, *Mathematics for Biologists*, Sofia, 1981, Nauka i Izkustvo, (Bulgarian).
2. Bl. Sendov, R. Maleev, S. Markov, S. Tashev, *Mathematics for Biologists*, Sofia, 1991, Publ. House of Sofia University, (Bulgarian).
3. P. Djakov, R. Levy, R. Maleev, S. Troyanski, *Differential and Integral Calculus. Functions of One Variable*, Demetra, Sofia, 2004, (Bulgarian).

ГОДИШНИК НА СОФИЙСКИЯ УНИВЕРСИТЕТ СВ. КЛИМЕНТ ОХРИДСКИ “

ФАКУЛТЕТ ПО МАТЕМАТИКА И ИНФОРМАТИКА

Том 107

ANNUAL OF SOFIA UNIVERSITY ST. KLIMENT OHRIDSKI“

FACULTY OF MATHEMATICS AND INFORMATICS

Volume 107

РАЗНОПОСОЧНИ СПОМЕНИ ЗА РУМЕН МАЛЕЕВ – ПРИЯТЕЛЯ
И КОЛЕГАТА

СТАНИМИР ТРОЯНСКИ

Румен е роден на 17 август 1943 г. в гр. Самоков в семейството на околийския инженер Ханс-Петър Малеев. Баща му е от смесен брак. Неговата майка (бабата на Румен) е немкиня, дъщеря на професор от Берлин. Дядото на Румен, Иван Малеев, завършва медицина във Франция и специализира педиатрия в Берлин. Запален турист, той е първият българин, който през 1903 г. изкачва Монблан.

След завръщането си в България основава Музикалното дружество в София, което е предшественик на Българското Музикално дружество. Целогодишно има запазена ложка в Софийската Опера.

Полага големи усилия за разпространяване на здравна култура. Издава първия ”Домашен лекар” през 1927 г. и дълги години се бори за откриване на лекарски кабинети в училищата. Иван Малеев и Анна (бабата немкиня) идват в България, раждат им се три момчета, които получават немско възпитание. Най-големият син е бащата на Румен, завършил Technische Hochschule в Берлин, специалност строително инженерство.

Усетих немското възпитание още със запознаването ми с бащата на Румен. Той се усмихна, подаде ми ръка и попита ”Ще пиеш ли бира”? Нещо измрънках. Той ми каза: Има два отговора: ”Да, моля!” и ”Не, благодаря!”.

Макар и ”де юре” бащата на Румен да не е засегнат от така наречените ”мероприятия на народната власт”, то на практика събитията след 1944 г. променят живота на семейството. Бащата на Румен съкращава името си от Ханс-Петър на Петър, загубва правото да работи като инженер и е пратен бригадир в строящия се Димитровград. След връщането си в София свири

в оркестъра на Музикалния театър. След възстановяването на инженерните права, работи и се пенсионира като обикновен инженер, въпреки че е от водещите специалисти. И всичко това, защото произхожда от немско семейство. Все пак редица нестандартни строителни обекти са минали през ръцете на бащата на Румен. Спомням си как инж. Малеев обсъждаше с нас (Румен и мен) построяването на първата ски-писта за скокове в България, как да оптимизира дължината и с какви криви да направи ски шанцата, така че да се увеличи скоростта, а от там дължината на скока.

С Румен се запознах през есента на далечната 1967 година. И двамата токущо бяхме постъпили на работа в Математическия Институт на БАН. Част от Института се намирал в една двуетажна сграда в Борисовата градина в началото на ул. Латинка, кв. Изток, срещу блока, в който живееше тогавашният Директор на Института проф. Л. Илиев. На първият етаж бяха кабинетите на по-старите сътрудници, А. Обретенов, А. Анчев, П. Русев, В. Чакалов, Е. Димитров, В. Спиридонов и други. На втория етаж бяха кабинетите на младите. В една голяма стая имаше 6-7 бюра. Там работеха М. Узунов, Ц. Игнатов, Д. Димитров, ние с Румен и други, на които вече не помня имената. Основната група, известни днес български математици, бяха аспиранти в Москва, Санкт Петербург и др. Ние с Румен ходехме на работа следобед. Със сериозна математика не се занимавахме, опознавахме се, разказвахме си кой къде е завършил, какво е специализирал. Румен беше завършил Букурещкия Университет, специалност – механика. Аз – Харковския университет, специалност – математика. Един следобед, като отидох на Латинка, видях Румен ядосан. Някой беше “изгравирал” отзад на стола му: МАЛЕЕВА- ЖИВКОВА с главни букви. Същата вечер на чаша вино той ми разказа по-подробно за семейството си.

През декември 1967 г. влязохме в казармата. За наш късмет се оказахме в едно отделение, в рота Инженерни войски, по-точно строителство на пътища и мостове. Както става в армията, всичко е бъркотия, но това не ни пречеше. Запознахме се с интересни личности. С нас служеха: Любомир Филипов, бъдещ Управител на Българска Народна Банка, Дончо Конакчиев, бъдещ вицепремиер в правителството на Жан Виденов, Тодор Гичев, по-късно професор в Университета по Архитектура, Строителство и Геодезия, Както можете да си представите, ние математиците нищо не разбирахме от строителство, но и двамата нито внимавахме на “лекциите”, нито четяхме учебниците по време на самоподготовка. Изобщо бяхме еталон за несретници. Веднъж Тодор Гичев се оплака от нас по време на лекции: “Другарю Пешлеевски, Малеев и Троянски непрекъснато си приказват и ми пречат да записвам!”. Все пак Румен, който беше получил сериозно образование по механика, и имаше допълнителни познания от баща си и брат си, също строителен инженер, нещо знайваше и с негова помощ и двамата взехме изпитите без да слушаме “лекциите” и ни произведоха младши лейтенанти. Случайно Румен прочете характеристиката, която му бяха дали в казармата: Умен, образован, но склонен да умува!

В казармата станахме приятели, на обед се увивахме заедно плътно с одеяла, за да се стоплим. В огромното спално помещение беше много студено.

След казармата нашето приятелство продължи, някак си станахме по-близки, въпреки че научните ни интереси бяха твърде различни. Румен беше задочен аспирант по механика в Букурещ. По онова време специалистите по механика бяха твърде много, благодарение на неизчерпаемия ентузиазъм на проф. Б. Долапчиев. Някои от младите механици се преориентираха към други области на математиката (напр. Е. Христов, К. Кирчев, Н. Никифоров и др.). Румен стана асистент в катедрата по Математическо Моделиране и започна да се интересува от Числени Методи и Теория на апроксимациите. Специализира Числени методи една година в Московския Университет, под ръководството на проф. А. Г. Дьяконов, публикува статия в Журнал вычислительной математики и математической физики. Постепенно, Теория на функциите стана пресечната точка на нашите научни интереси. Интензивно започнахме да работим заедно по задачи от Геометрия на Банаховите пространства, в които важен апарат е Конструктивната теория на Функциите. Намерихме асимптотически оценки за различни характеристики за гладкост и изпъкналост на единичното кълбо на Банахови пространства. Ще илюстрирам казаното с един пример. Получихме асимптотически оценки за модула на изпъкналост във функционалните пространства на Орлич $L_M[0, 1]$. По точно доказахме, че съществува функция на Орлич N , еквивалентна в безкрайността на M , такава че

$$\delta_X(\varepsilon) \geq C_M \varepsilon^2 \inf \left\{ \frac{M(u\nu)}{u^2 M(\nu)} : \varepsilon \leq u \leq 1 \leq \nu \right\},$$

където X е пространството $L_M[0, 1]$, снабдено с норма $\|\cdot\|_N$, породена от функцията N , а константата C_M , зависи само от M , е положителна тогава и само тогава, когато

$$\liminf_{t \rightarrow \infty} \left(\frac{tM'(t)}{M(t)} \right) > 1.$$

Впоследствие се оказа, че от работите на Т. Figiel и G. Pisier следва, че получената оценка е точна по порядък в класа от еквивалентни норми в $L_M[0, 1]$. Румен продължи изследванията за диференцируемост по Фреше и Гато от повисок ред на нормите в пространствата на Орлич и техните обобщения. Макар и не многобройни, работите на Румен са публикувани в известните специализирани списания по Анализ и по-специално по Функционален Анализ: *Studia Math.*, *Functional Analysis and Appl.*, *J. Math. Analysis and Appl.*, *Arch. Math.*, *Zeitschrift Anal. Anwend.*, както и в издания на Американското и Канадското Математически общества. Неговите резултати не остават незабелязани, споменати са в повечето монографии по Геометрия на Банаховите пространства: “Classical Banach spaces” на J. Lindenstrass, L. Tzafriri, “Series and sequences in Banach spaces” на J. Diestel, “Smoothnes and renormings in Banach spaces” на R. Deville, G. Godefroy, V. Zizler и много други. Оценките за остатъчния член

във формулата за развитие по Тейлор на нормата в $L_M[0, 1]$ са изложени подробно, с доказателства в монографията на P.Hájek, M.Johanis "Smooth Analysis in Banach spaces".

Румен беше внимателен, готов да изслуша всеки, да вникне в това, от което се интересуваше другият, и да работи по поставения проблем. Негови съавтори са: А. Андреев, Д. Върндев, Р. Гонзало (R.Gonzalo), П. Джаков, Б. Златанов, Р. Леви, С. Марков, Г. Недев, Бл. Сендов, С. Ташев и моя милост. Румен е съавтор на първия учебник по математика, ориентиран към студентите по биология, претърпял две издания, съавтор е и на учебника по Диференциално и Интегрално Смятане, по който се стараехме да четем лекции. Обясняваше точно и ясно на студентите.

Научната му кариера напредваше традиционно за учен: кандидат, доктор на математическите науки, доцент, професор, заместник декан на Факултета по Математика и Информатика, функционален декан в Ректората на Софийския Университет, ръководител на катедра Математически Анализ, Член на Комисията по математика при ВАК. След демократичните промени научните ни пътища малко се раздалечиха. Аз започнах да се занимавам с приложенията на Топологията и Дескриптивната теория на множествата в Геометрията на Банаховите пространства. На Румен тези направления не му се понравиха, все пак той беше механик по образование. Ние продължихме да дискутираме различни въпроси от Теорията на Функционалните пространства, по специално диференцируемост в пространствата на Орлич, Лоренц и техните обобщения. Заедно имахме дипломанти, ръководехме проекти по програмата ТЕМПУС. Посещавахме международни работни семинари по Банахови пространства в Пасеки (Чехия), Спецос (Гърция), Монс (Белгия) и др. Заедно ръководехме докторанта Б. Златанов, сега професор в Пловдивския Университет.

Паралелно със заниманията по математика, Румен започна да се занимава с отчитането на резултатите след провеждане на избори в нашата страна. Той бързо навлезе в тази област, напълно непозната за българското общество, по понятни причини. Неговите способности бяха забелязани от ръководството на Организацията за Сигурност и Сътрудничество в Европа (OSCE) и той беше поканен за експерт. Освен че имаше организаторски способности и аналитичен начин на мислене, Румен беше дипломат по рождение и владееше пет чужди езика. Румен като експерт оценяваше как се прилага законът, доколко изборите в дадена държава като цяло са демократични. В това си качество той посети Украйна, Таджикистан, Туркменистан, Канада, Италия, Черна гора, Словакия и др.

Румен боледува десет години. Бори се стоически! Стандартният отговор на въпроса ми как се чувства беше: "Бива!". Никога не говореше нищо за коварните болести, чийто изход бе предварително известен. Щом се намирах в София, отивахме на ресторант, понякога независимо че току-що е излязал от болницата. Шегуваше се, интересуваше се от последните новини от математическата колегия. През септември 2019 ми съобщи, че повече няма нужда да

ходи на процедури в болницата. На въпроса ми, кога е удобно да го посетя, той ми отговори: “Аз съм вкъщи, когато искаш, ела.” Посещавах го, говорихме дълго по телефона. Последно ми каза: “Щастлив бях, че имах такъв приятел!” Всъщност, късметлията бях аз. Общуването ми с Румен е червената линия в моя живот. И досега сънувам, че говорим по телефона. Най-съкровениите му мисли от тези последни дни вероятно знае единствено съпругата му, която се грижеше за него до кончината му на 16 декември 2019.

Мир на праха му !

Summary. This is Stanimir Troyanski’s personal account of the late Professor Rumen Maleev. The author shares his reminiscences about his dear friend and colleague.

Received on June 22, 2021

STANIMIR TROYANSKI
Bulgarian Academy of Sciences
Sofia
BULGARIA
E-mails: troyanski@math.bas.bg
stroya@um.es

ГОДИШНИК НА СОФИЙСКИЯ УНИВЕРСИТЕТ „СВ. КЛИМЕНТ ОХРИДСКИ“

ФАКУЛТЕТ ПО МАТЕМАТИКА И ИНФОРМАТИКА

Том 107

ANNUAL OF SOFIA UNIVERSITY „ST. KLIMENT OHRIDSKI“

FACULTY OF MATHEMATICS AND INFORMATICS

Volume 107

REPRESENTATION OF NATURAL NUMBERS BY SUMS OF FOUR SQUARES OF ALMOST-PRIME HAVING A SPECIAL FORM

ZHIVKO H. PETROV AND TATYANA. L. TODOROVA

In this paper we consider the equation $x_1^2 + x_2^2 + x_3^2 + x_4^2 = N$, where N is a sufficiently large integer and prove that if η is quadratic irrational number and $0 < \lambda < \frac{1}{10}$, then it has a solution in almost-prime numbers x_1, \dots, x_4 , such that $\{\eta x_i\} < N^{-\lambda}$ for $i = 1, \dots, 4$.

Keywords: Lagrange's equation, almost-primes, quadratic irrational numbers.

2020 Math. Subject Classification: Primary 11P05; Secondary 11N36.

1. INTRODUCTION AND STATEMENT OF THE RESULT

In 1770 Lagrange proved that for any positive integer N the equation

$$x_1^2 + x_2^2 + x_3^2 + x_4^2 = N \tag{1.1}$$

has a solution in integer numbers x_1, \dots, x_4 . Later Jacobi found an exact formula for the number of the solutions (see [8, Ch. 20]). A lot of researchers studied the equation (1.1) for solvability in integers satisfying additional conditions. There is a hypothesis stating that if N is sufficiently large and $N \equiv 4 \pmod{24}$ then (1.1) has a solution in primes. This hypothesis has not been proved so far, but several approximations to it have been established.

In 1994 J. Brüdern and E. Fouvry [1] proved that for any large $N \equiv 4 \pmod{24}$, the equation (1.1) has a solution in $x_1, \dots, x_4 \in \mathcal{P}_{34}$. (We say that integer n is an almost-prime of order r if n has at most r prime factors, counted with their multiplicities. We denote by \mathcal{P}_r the set of all almost-primes of order r .) This result was improved by D. R. Heath-Brown and D. I. Tolev [9]. They showed that, under the same restrictions for N , the equation (1.1) has a solution in prime x_1 and almost-prime $x_2, x_3, x_4 \in \mathcal{P}_{101}$. In their paper they also proved that the equation has a solution in $x_1, \dots, x_4 \in \mathcal{P}_{25}$. In 2020 Tak Wing Ching [2] improved this result with three of them being in \mathcal{P}_3 and the other in \mathcal{P}_4 .

On the other hand, let us consider a subset of the set of integers having the form

$$\mathcal{A} = \{n \mid a < \{\eta n\} < b\},$$

where η is a fixed quadratic irrational number, and $a, b \in [0, 1]$.

Denote by $I(N)$ the number of solutions of (1.1) in arbitrary integers and by $J(N)$ the number of solutions of (1.1) in integers from the set \mathcal{A} .

In 2011 S. A. Gritsenko and N. N. Motkina [6] proved that for any positive small ε , the following formula holds

$$J(N) = (b - a)^4 I(N) + O(N^{0.9+3\varepsilon}).$$

S. A. Gritsenko and N. N. Motkina consider many others additive problem in witch variables are in special set of numbers similar to \mathcal{A} . (See [4] – [5] and [7].) In 2013 A. V. Shutov [12] considered solvability of diophantine equation in integer numbers from \mathcal{A} . Further research in this area was made by A. V. Shutov and A. A. Zhukova [13].

We consider the equation (1.1), where x_i are almost-prime numbers and belong to a set similar to \mathcal{A} . Our result is

Theorem 1.1. *Let η be a quadratic irrational number, $0 < \lambda < \frac{1}{10}$ and $k = \left\lceil \frac{54}{1-10\lambda} \right\rceil$. Then for every sufficiently large integer N , the equation (1.1) has a solution in almost-prime numbers $x_1, \dots, x_4 \in \mathcal{P}_k$, such that $\{\eta x_i\} < N^{-\lambda}$, $i = 1, 2, 3, 4$.*

In the present paper we use the following notations.

We denote by N a sufficiently large odd integer and $P = N^{\frac{1}{2}}$. Letters a, b, k, l, m, n, q, p always stand for integers. By (n_1, \dots, n_k) we denote the greatest common divisor of n_1, \dots, n_k . Let $\|t\|$ denote the distance from t to the nearest integer. We denote by \vec{n} four dimensional vectors and let

$$|\vec{n}| = \max(|n_1|, \dots, |n_4|). \tag{1.2}$$

As usual, $\mu(q)$ is the Möbius function and $\tau(q)$ is the number of positive divisors of q . Sometimes we write $a \equiv b(q)$ as an abbreviation of $a \equiv b \pmod{q}$.

We write $\sum_{x \pmod{q}}$ for a sum over a complete system of residues modulo q and respectively $\sum_{x \pmod{q}}^*$ is a sum over a reduced system of residues modulo q . We also denote $e(t) = e^{2\pi it}$.

We use Vinogradov's notation $A \ll B$, which is equivalent to $A = O(B)$. By ε we denote an arbitrarily small positive number, which is not the same in different occurrences. The constants in the O -terms and \ll -symbols are absolute or depend on ε .

2. AUXILIARY RESULTS

Now we introduce some lemmas, which shall be used later.

Lemma 2.1. *Suppose that $D \in \mathbb{R}, D > 4$. There exist arithmetical functions $\lambda^\pm(d)$ (called Rosser's functions of level D) with the following properties:*

1. *For any positive integer d we have*

$$|\lambda^\pm(d)| \leq 1, \quad \lambda^\pm(d) = 0 \quad \text{if } d > D \quad \text{or} \quad \mu(d) = 0.$$

2. *If $n \in \mathbb{N}$ then*

$$\sum_{d|n} \lambda^-(d) \leq \sum_{d|n} \mu(d) \leq \sum_{d|n} \lambda^+(d).$$

3. *If $z \in \mathbb{R}$ is such that $z^2 \leq D$ and if*

$$P(z) = \prod_{2 < p < z} p, \quad \mathcal{B} = \prod_{2 < p < z} \left(1 - \frac{1}{p-1}\right), \quad \mathcal{N}^\pm = \sum_{d|P(z)} \frac{\lambda^\pm(d)}{\varphi(d)}, \quad s_0 = \frac{\log D}{\log z}, \quad (2.1)$$

then we have

$$\mathcal{B} \leq \mathcal{N}^+ \leq \mathcal{B} \left(F(s_0) + O\left((\log D)^{-\frac{1}{3}}\right) \right), \quad (2.2)$$

$$\mathcal{B} \geq \mathcal{N}^- \geq \mathcal{B} \left(f(s_0) + O\left((\log D)^{-\frac{1}{3}}\right) \right), \quad (2.3)$$

where $F(s)$ and $f(s)$ satisfy

$$\begin{aligned} F(s) &= 2e^\gamma s^{-1}, & \text{if } 2 \leq s \leq 3, \\ f(s) &= 2e^\gamma s^{-1} \log(s-1), & \text{if } 2 \leq s \leq 3, \\ (sF(s))' &= f(s-1), & \text{if } s > 3, \\ (sf(s))' &= F(s-1), & \text{if } s > 2. \end{aligned}$$

Here γ is Euler's constant.

Proof. See Greaves [3, Chapter 4]. □

Lemma 2.2. *Suppose that Λ_i, Λ_i^\pm are real numbers satisfying $\Lambda_i = 0$ or 1 , $\Lambda_i^- \leq \Lambda_i \leq \Lambda_i^+$, $i = 1, 2, 3, 4$. Then*

$$\begin{aligned} \Lambda_1 \Lambda_2 \Lambda_3 \Lambda_4 \geq & \Lambda_1^- \Lambda_2^+ \Lambda_3^+ \Lambda_4^+ + \Lambda_1^+ \Lambda_2^- \Lambda_3^+ \Lambda_4^+ + \Lambda_1^+ \Lambda_2^+ \Lambda_3^- \Lambda_4^+ + \\ & + \Lambda_1^+ \Lambda_2^+ \Lambda_3^+ \Lambda_4^- - 3\Lambda_1^+ \Lambda_2^+ \Lambda_3^+ \Lambda_4^+. \end{aligned} \quad (2.4)$$

Proof. The proof is similar to the proof of [1, Lemma 13]. □

Let

$$w_0(t) = \begin{cases} e^{t^2 - \frac{1}{25}} & \text{if } t \in \left(-\frac{4}{5}, \frac{4}{5}\right), \\ 0 & \text{if } t \notin \left(-\frac{4}{5}, \frac{4}{5}\right) \end{cases}$$

and

$$w(x) = w_0\left(\frac{x}{P} - \frac{1}{2}\right). \quad (2.5)$$

Lemma 2.3. *Let $u, \beta \in \mathbb{R}$ and*

$$J(\beta, u) = \int_{-\infty}^{+\infty} w_0\left(x - \frac{1}{2}\right) e(\beta x^2 + ux) dx. \quad (2.6)$$

Then:

1. *For every $k \in \mathbb{N}$ and $u \neq 0$ we have*

$$J(\beta, u) \ll_k \frac{1 + |\beta|^k}{|u|^k}.$$

2. *The following inequality hold*

$$J(\beta, u) \ll \min\left(1, |\beta|^{-\frac{1}{2}}\right).$$

Proof. See [9, Lemma 9]. □

Lemma 2.4. *Suppose that $\vec{u} \in \mathbb{Z}^4$ and*

$$J(\beta, \vec{u}) = \prod_{i=1}^4 J(\beta, u_i).$$

Then we have

$$\int_{-\infty}^{+\infty} |J(\beta, \vec{u})| d\gamma \ll |\vec{u}|^{-1+\varepsilon}.$$

Proof. Proof can be find in [9, Lemma 10]. □

Lemma 2.5. *There exists a function $\sigma(v, q, \gamma)$ defined for $-\frac{q}{2} < v \leq \frac{q}{2}$, $q \leq P$, $|\gamma| \leq \frac{P}{q}$, integrable with respect to γ , satisfying*

$$|\sigma(v, q, \gamma)| \leq \frac{1}{1 + |v|}$$

and also for every $a \in \mathbb{Z}$, $(a, q) = 1$ we have

$$\sum_{-\frac{q}{2} < v \leq \frac{q}{2}} e\left(\frac{\bar{a}v}{q}\right) \sigma(v, q, \gamma) = \begin{cases} 1 & \text{if } \gamma \in \mathcal{N}(a, q), \\ 0 & \text{otherwise,} \end{cases}$$

where

$$\mathcal{N}(a, q) = \left(-\frac{P^2}{q(q+q')}, \frac{P^2}{q(q+q'')} \right]$$

and

$$P < q + q', q + q'' \leq P + q, \quad aq' \equiv 1 \pmod{q}, \quad aq'' \equiv -1 \pmod{q}. \quad (2.7)$$

Proof. See [15, Lemma 45]. □

For $q \in \mathbb{N}$ and $m, n \in \mathbb{Z}$, the Gauss sum is defined by

$$G(q, m, n) = \sum_{x(q)} e\left(\frac{mx^2 + nx}{q}\right). \quad (2.8)$$

For $\vec{d} = \langle d_1, \dots, d_4 \rangle \in \mathbb{Z}^4$ and $\vec{n} = \langle n_1, \dots, n_4 \rangle \in \mathbb{Z}^4$ we denote

$$G(q, a\vec{d}^2, \vec{n}) = \prod_{i=1}^4 G(q, ad_i^2, n_i).$$

We need to estimate an exponential sum of the form

$$V_q = V_q(N, \vec{d}, v, \vec{n}) = \sum_{a(q)}^* e\left(\frac{\bar{a}v - Na}{q}\right) G(q, a\vec{d}^2, \vec{n}). \quad (2.9)$$

To estimate V_q we use the properties of the Gauss sum and the Kloosterman sum.

Lemma 2.6. *Suppose that $N, q \in \mathbb{N}$, $v \in \mathbb{Z}$ and $\vec{d}, \vec{n} \in \mathbb{Z}^4$. Then we have*

$$V_q(N, \vec{d}, v, \vec{n}) \ll q^{\frac{5}{2}} \tau(q)(q, N)^{\frac{1}{2}}(q, d_1)(q, d_2)(q, d_3)(q, d_4).$$

Moreover, if some of the conditions

$$(q, d_i) | n_i, \quad i = 1, \dots, 4$$

do not hold, then $V_q(N, \vec{d}, v, \vec{n}) = 0$.

Proof. This result is analogous to this one in [1, Lemma 1]. □

Lemma 2.7. (Liouville) *If η is an irrational number which is the root of a polynomial f of degree 2 with integer coefficients, then there exists a real number $A > 0$ such that, for all integers p, q , with $q > 0$,*

$$\left| \eta - \frac{p}{q} \right| \geq \frac{A}{q^2}.$$

Proof. See [11, Theorem 1A]. □

3. PROOF OF THE THEOREM

3.1. BEGINNING OF THE PROOF

Let N be a sufficiently large integer. We denote

$$z = N^\alpha, \quad P(z) = \prod_{p < z} p, \quad \delta = N^{-\lambda}.$$

We apply the well-known Vinogradov's "little cups" lemma (see [10, Chapter 1, Lemma A]) with parameters

$$\alpha_1 = \frac{\delta}{4}, \quad \beta_1 = \frac{3\delta}{4}, \quad \Delta = \frac{\delta}{2}, \quad r = [\log N]$$

and construct a function $\theta(t)$ which is periodic with period 1 and has the following properties:

$$\theta\left(\frac{\delta}{2}\right) = 1; \quad 0 < \theta(t) < 1 \quad \text{for} \quad 0 < t < \frac{\delta}{2} \quad \text{or} \quad \frac{\delta}{2} < t < \delta;$$

$$\theta(t) = 0 \quad \text{for} \quad \delta \leq t \leq 1.$$

Furthermore, from the Fourier series of $\theta(t)$ we find

$$\theta(t) = \frac{\delta}{2} + \sum_{\substack{0 < |m| \leq H \\ m \neq 0}} c(m) e(mt) + O(P^{-A}), \tag{3.1}$$

with

$$|c(m)| \leq \min \left(\frac{\delta}{2}, \frac{1}{|m|} \left(\frac{[\log N]}{\delta \pi |m|} \right)^{[\log N]} \right),$$

where A is arbitrary large constant and

$$H = \frac{[\log N]^2}{\delta}. \quad (3.2)$$

Let us denote

$$\theta(\eta\vec{x}) = \theta(\eta x_1)\theta(\eta x_2)\theta(\eta x_3)\theta(\eta x_4)$$

and

$$w(\vec{x}) = w(x_1)w(x_2)w(x_3)w(x_4).$$

We consider the sum

$$\Gamma = \sum_{\substack{x_1^2+x_2^2+x_3^2+x_4^2=N \\ (x_i, P(z))=1, i=1,2,3,4}} \theta(\eta\vec{x})w(\vec{x}).$$

From the condition $(x_i, P(z)) = 1$ it follows that any prime factor of x_i is greater than or equal to z . Suppose that x_i has l prime factors, counted with their multiplicities. Then we have

$$N^{\frac{1}{2}} \geq x_i \geq z^l = N^{\alpha l}$$

and hence $l \leq \frac{1}{2\alpha}$. This implies that if $\Gamma > 0$ then equation (1.1) has a solution in almost-prime numbers x_1, \dots, x_4 with at most $\lceil \frac{1}{2\alpha} \rceil$ prime factors, such that $\{\eta x_i\} < N^{-\lambda}$, $i = 1, \dots, 4$.

For $i = 1, 2, 3, 4$ we define

$$\Lambda_i = \sum_{d|(x_i, P(z))} \mu(d) = \begin{cases} 1 & \text{if } (x_i, P(z)) = 1, \\ 0 & \text{otherwise.} \end{cases} \quad (3.3)$$

Then we find that

$$\Gamma = \sum_{x_1^2+x_2^2+x_3^2+x_4^2=N} \Lambda_1\Lambda_2\Lambda_3\Lambda_4\theta(\eta\vec{x})w(\vec{x}).$$

We can write Γ as

$$\Gamma = \sum_{x_i \in \mathbb{Z}} \Lambda_1\Lambda_2\Lambda_3\Lambda_4\theta(\eta\vec{x})w(\vec{x}) \int_0^1 e(\alpha(x_1^2 + x_2^2 + x_3^2 + x_4^2 - N)) d\alpha.$$

Suppose that $\lambda^\pm(d)$ are the Rosser functions of level D (see Lemma 2.1). Let also denote

$$\Lambda_i^\pm = \sum_{d|(x_i, P(z))} \lambda^\pm(d), \quad i = 1, 2, 3, 4. \quad (3.4)$$

Then from Lemma 2.1, (3.3) and (3.4) we find that

$$\Lambda_i^- \leq \Lambda_i \leq \Lambda_i^+.$$

We use Lemma 2.2 and find that

$$\Gamma \geq \Gamma_1 + \Gamma_2 + \Gamma_3 + \Gamma_4 - 3\Gamma_5,$$

where $\Gamma_1, \dots, \Gamma_5$ are the contributions coming from the consecutive terms of the right side of (2.4). We have $\Gamma_1 = \Gamma_2 = \Gamma_3 = \Gamma_4$ and

$$\begin{aligned} \Gamma_1 &= \sum_{x_i \in \mathbb{Z}} \Lambda_1^- \Lambda_2^+ \Lambda_3^+ \Lambda_4^+ \theta(\eta \vec{x}) w(\vec{x}) \int_0^1 e(\alpha(x_1^2 + x_2^2 + x_3^2 + x_4^2 - N)) d\alpha, \\ \Gamma_5 &= \sum_{x_i \in \mathbb{Z}} \Lambda_1^+ \Lambda_2^+ \Lambda_3^+ \Lambda_4^+ \theta(\eta \vec{x}) w(\vec{x}) \int_0^1 e(\alpha(x_1^2 + x_2^2 + x_3^2 + x_4^2 - N)) d\alpha. \end{aligned}$$

Hence, we get

$$\Gamma \geq 4\Gamma_1 - 3\Gamma_5. \tag{3.5}$$

3.2. ASYMPTOTIC FORMULA FOR Γ_1

We shall find an asymptotic formula for the integral Γ_1 . We have

$$\begin{aligned} \Gamma_1 &= \sum_{d_i | P(z)} \lambda^-(d_1) \lambda^+(d_2) \lambda^+(d_3) \lambda^+(d_4) \sum_{x_i \equiv 0(d_i)} \theta(\eta \vec{x}) w(\vec{x}) \times \\ &\quad \times \int_0^1 e(\alpha(x_1^2 + \dots + x_4^2 - N)) d\alpha \\ &= \sum_{d_i | P(z)} \lambda^-(d_1) \lambda^+(d_2) \lambda^+(d_3) \lambda^+(d_4) \times \\ &\quad \times \int_0^1 \prod_{1 \leq i \leq 4} \left(\sum_{x \equiv 0(d_i)} \theta(\eta x) w(x) e(\alpha x^2) \right) e(-N\alpha) d\alpha. \end{aligned}$$

Let

$$S(\alpha, d, m) = \sum_{\substack{x \in \mathbb{Z} \\ x \equiv 0(d)}} w(x) e(\alpha x^2 + m\eta x). \tag{3.6}$$

Then using the Fourier series of $\theta(t)$ (see (3.1)), we find

$$\sum_{x \equiv 0(d)} \theta(\eta x) w(x) e(\alpha x^2) = \sum_{|m| \leq H} c(m) \sum_{x \equiv 0(d)} w(x) e(\alpha x^2 + m\eta x) + O(P^{-A}).$$

Denoting

$$S(\alpha, \vec{d}, \vec{m}) = S(\alpha, d_1, m_1) S(\alpha, d_2, m_2) S(\alpha, d_3, m_3) S(\alpha, d_4, m_4) \tag{3.7}$$

and

$$\lambda(\vec{d}) = \lambda^-(d_1) \lambda^+(d_2) \lambda^+(d_3) \lambda^+(d_4), \tag{3.8}$$

we find that

$$\Gamma_1 = \sum_{d_i|P(z)} \lambda(\vec{d}) \sum_{\substack{|m_i| \leq H \\ i=1,2,3,4}} c(m_i) \int_0^1 S(\alpha, \vec{d}, \vec{m}) e(-N\alpha) d\alpha + O(1).$$

We divide Γ_1 into two parts:

$$\Gamma_1 = \Gamma_1^0 + \Gamma_1^* + O(1),$$

where

$$\Gamma_1^0 = c^4(0) \sum_{d_i|P(z)} \lambda(\vec{d}) \sum_{\substack{x_i \equiv 0(d_i) \\ x_1^2 + x_2^2 + x_3^2 + x_4^2 = N}} w(\vec{x})$$

and

$$\Gamma_1^* = \sum_{d_i|P(z)} \lambda(\vec{d}) \sum_{\substack{0 < |m_i| \leq H \\ i=1,2,3,4}} c(m_i) \int_0^1 S(\alpha, \vec{d}, \vec{m}) e(-N\alpha) d\alpha. \quad (3.9)$$

Hence

$$\Gamma \geq 4\Gamma_1^0 - 3\Gamma_5^0 + O(\Gamma_1^*) + O(\Gamma_5^*) + O(1). \quad (3.10)$$

According to [1] and [9], for $D \leq P^{1/8-\varepsilon}$, $s = \frac{\log D}{\log z} = 3.13$ the estimate

$$4\Gamma_1^0 - 3\Gamma_5^0 \gg \frac{C\delta N}{(\log N)^4} + O(\delta P^{3/2+\varepsilon} D^4) \quad (3.11)$$

with some constant C is obtained. Thus it suffices to evaluate Γ_1^* and Γ_5^* .

3.3. ESTIMATION OF Γ_1^*

In this subsection we find the upper bound for Γ_1^* defined in (3.9). The function in the integral in Γ_1^* is periodic with period 1, so we can integrate over the interval \mathcal{I} defined as

$$\mathcal{I} = \left(\frac{1}{1 + [P]}, 1 + \frac{1}{1 + [P]} \right).$$

We apply the Kloosterman form of the Hardy-Littlewood circle method. We divide the interval only into large arcs. Using the properties of the Farey fractions, we represent \mathcal{I} as an union of disjoint intervals in the following way:

$$\mathcal{I} = \bigcup_{q \leq P} \bigcup_{\substack{a=1 \\ (a,q)=1}}^q \mathcal{L}(a, q),$$

where

$$\mathcal{L}(a, q) = \left[\frac{a}{q} - \frac{1}{q(q+q')}, \frac{a}{q} + \frac{1}{q(q+q'')} \right]$$

and where the integers q', q'' are specified in (2.7). Then

$$\Gamma_1^* = \sum_{d_i|P(z)} \lambda(\vec{d}) \sum_{\substack{0 < |m_i| \leq H \\ i=1,2,3,4}} c(m_i) \sum_{q \leq P} \sum_{\substack{a=1 \\ (a,q)=1}}^q \int_{\mathcal{L}(a,q)} S(\alpha, \vec{d}, \vec{m}) e(-N\alpha) d\alpha.$$

We change variable of integration $\alpha = \frac{a}{q} + \beta$ to get

$$\begin{aligned} \Gamma_1^* &= \sum_{d_i|P(z)} \lambda(\vec{d}) \sum_{\substack{0 < |m_i| \leq H \\ i=1,2,3,4}} c(m_i) \sum_{q \leq P} \sum_{\substack{a=1 \\ (a,q)=1}}^q \times \\ &\quad \times \int_{\mathcal{M}(a,q)} S\left(\frac{a}{q} + \beta, \vec{d}, \vec{m}\right) e\left(-N\left(\frac{a}{q} + \beta\right)\right) d\beta, \end{aligned}$$

where

$$\mathcal{M}(a, q) = \left[-\frac{1}{q(q+q')}, \frac{1}{q(q+q'')} \right].$$

From (2.7) we find that

$$\left[-\frac{1}{2qP}, \frac{1}{2qP} \right] \subset \mathcal{M}(a, q) \subset \left[-\frac{1}{qP}, \frac{1}{qP} \right]$$

and hence

$$|\beta| \leq \frac{1}{qP} \quad \text{for } \beta \in \mathcal{M}(a, q). \quad (3.12)$$

Now we consider the sum $S(\alpha, d_i, m_i)$ defined in (3.6). As η is irrational number, $\|s\eta\| \neq 0$ for all $s \in \mathbb{Z}$. Using that fact and working as in the proof of [9, Lemma 12], we find that for $\beta \in \mathcal{M}(a, q)$ we have

$$\begin{aligned} S\left(\frac{a}{q} + \beta, d_i, m_i\right) &= \frac{P}{d_i q} \sum_{|n - m_i d_i q \eta| < M_i} J\left(\beta P^2, \left(m_i \eta - \frac{n}{d_i q}\right) P\right) G(q, a d_i^2, n) + \\ &\quad + O(P^{-B}), \end{aligned} \quad (3.13)$$

where $G(q, m, n)$ and $J(\gamma, u)$ are defined respectively by (2.8) and (2.6), B is an arbitrarily large constant, $M_i = d_i P^\varepsilon$, $\varepsilon > 0$ is arbitrarily small and the constant in the O -term depends only on B and ε . We leave the verification of the last formula to the reader.

Let

$$F(P, \vec{d}) = \sum_{\substack{0 < |m_i| \leq H \\ i=1,2,3,4}} c(m_i) \sum_{q \leq P} \sum_{a(q)}^* e\left(-\frac{aN}{q}\right) \int_{\mathcal{M}(a,q)} S\left(\frac{a}{q} + \beta, \vec{d}, \vec{m}\right) e(-\beta N) d\beta.$$

It is obvious that

$$\Gamma_1^* = \sum_{d_i|P(z)} \lambda(\vec{d}) F(P, \vec{d}). \quad (3.14)$$

Using (3.13) and Lemma 2.3 we get

$$F(P, \vec{d}) = F^*(P, \vec{d}) + O(1), \quad (3.15)$$

where

$$F^*(P, \vec{d}) = \frac{P^4}{d_1 d_2 d_3 d_4} \sum_{\substack{0 < |m_i| \leq H \\ i=1,2,3,4}} c(m_i) \sum_{q \leq P} \frac{1}{q^4} \sum_{a(q)}^* e\left(-\frac{aN}{q}\right) \times \\ \times \sum_{|n_i - m_i d_i q \eta| < M_i} G(q, ad_i^2, \vec{n}) \int_{\mathcal{N}(a,q)} J\left(\beta P^2, \left(\vec{m}\eta - \frac{\vec{n}}{dq}\right)P\right) e(-\gamma) d\gamma.$$

Using Lemma 2.5 and working as in the proof of [14, Lemma 2] we find that

$$F^*(P, \vec{d}) = F'(P, \vec{d}) + O(P^{3/2+\varepsilon}), \quad (3.16)$$

where

$$F'(P, \vec{d}) = \frac{P^2}{d_1 d_2 d_3 d_4} \sum_{\substack{0 < |m_i| \leq H \\ i=1,2,3,4}} c(m_i) \sum_{q \leq P} \frac{1}{q^4} \sum_{\substack{|n_i - m_i d_i q \eta| < M_i \\ (q, d_i) | n_i, i=1, \dots, 4}} V_q(N, \vec{d}, 0, \vec{n}) \times \\ \times \int_{|\gamma| \leq \frac{P}{2q}} J\left(\gamma, \left(\vec{m}\eta - \frac{\vec{n}}{dq}\right)P\right) e(-\gamma) d\gamma,$$

and $V_q(N, \vec{d}, 0, \vec{n})$ is defined by (2.9). We represent the sum $F'(P, \vec{d})$ as

$$F'(P, \vec{d}) = F_1 + F_2, \quad (3.17)$$

where F_1 is the contribution of these addends with $q \leq Q$ and F_2 for addends with $Q < q \leq P$. Here Q is parameter, which we choose later. Using Lemma 2.3 (2), Lemma 2.6 and (3.1), we get

$$F_2 \ll \frac{P^2 \delta^4}{d_1 d_2 d_3 d_4} \sum_{\substack{0 < |m_i| \leq H \\ i=1,2,3,4}} \sum_{Q < q \leq P} \frac{q^{5/2} \tau(q) (q, N)^{1/2} (q, d_1) \dots (q, d_4)}{q^4} \times \\ \times \sum_{\substack{|n_i - m_i d_i q \eta| < M_i \\ (q, d_i) | n_i, i=1, \dots, 4}} 1. \quad (3.18)$$

It is clear that the sum over \vec{n} in the expression above is

$$\ll \prod_{1 \leq i \leq 4} \sum_{\substack{-\frac{M_i + m_i d_i q \eta}{(q, d_i)} < t_i < \frac{M_i + m_i d_i q \eta}{(q, d_i)}} 1 \ll \frac{M_1 M_2 M_3 M_4}{(q, d_1)(q, d_2)(q, d_3)(q, d_4)} \\ \ll \frac{P^\varepsilon d_1 d_2 d_3 d_4}{(q, d_1)(q, d_2)(q, d_3)(q, d_4)},$$

which, together with (3.18) and (3.2), gives

$$F_2 \ll P^{2+\varepsilon} \sum_{Q < q \leq P} \frac{\tau(q)(q, N)^{1/2}}{q^{3/2}}.$$

Now we apply Cauchy's inequality to get

$$\begin{aligned} F_2 &\ll P^{2+\varepsilon} \left(\sum_{Q < q \leq P} \frac{\tau^2(q)}{q} \right)^{\frac{1}{2}} \left(\sum_{Q < q \leq P} \frac{(q, N)}{q^2} \right)^{\frac{1}{2}} \\ &\ll P^{2+\varepsilon} \left(\sum_{\substack{t|N \\ t \leq P}} t \sum_{\substack{Q/t < q_1 \leq P/t}} \frac{1}{t^2 q_1^2} \right)^{\frac{1}{2}} \ll \frac{P^{2+\varepsilon}}{Q^{1/2}}. \end{aligned} \tag{3.19}$$

To evaluate F_1 we firstly apply Lemma 2.4 to get

$$\int_{|\gamma| \leq \frac{P}{2q}} \left| J \left(\gamma, \left(m\vec{\eta} - \frac{\vec{n}}{dq} \right) P \right) \right| d\gamma \ll \left(\left| \left(m\vec{\eta} - \frac{\vec{n}}{dq} \right) P \right| \right)^{-1+\varepsilon}.$$

Then using Lemma 2.6 and (3.2) we obtain

$$\begin{aligned} F_1 &\ll \frac{P^2}{d_1 d_2 d_3 d_4} \sum_{q \leq Q} \frac{q^{5/2} \tau(q)(q, N)^{1/2} (q, d_1) \dots (q, d_4)}{q^4} \times \\ &\quad \times \sum_{\substack{|n_i - m_i d_i q \eta| < M_i \\ (q, d_i) | n_i, i=1, \dots, 4}} \frac{1}{\left| \left(m\vec{\eta} - \frac{\vec{n}}{dq} \right) P \right|}. \end{aligned} \tag{3.20}$$

It is clear that if $n_i = (q, d_i)t_i$, $d_i = (q, d_i)d'_i$ and

$$\left| \left(m_i \eta - \frac{n_i}{d_i q} \right) P \right| = \frac{P(q, d_i)}{q d'_i} |t_i - m_i d'_i \eta q|,$$

then the sum over $\left(m\vec{\eta} - \frac{\vec{n}}{dq} \right) P$ in the expression above is

$$\ll \frac{q}{P} \sum_{\substack{|t_i - m_i d'_i \eta q| < \frac{M_i}{(q, d_i)}}} \frac{1}{\max_{1 \leq i \leq 4} (q, d_i) |t_i - m_i d'_i \eta q| / d_i}. \tag{3.21}$$

Let t_1^o is such that

$$|t_1^o - m_1 d'_1 \eta q| = | - m_1 d'_1 \eta q | = | m_1 d'_1 \eta q |.$$

As η is quadratic irrational number, then $|m_1 d'_1 \eta q| \neq 0$ and for $t_1 \neq t_1^o$ we have $|t_1 - m_1 d'_1 \eta q| \geq 1/2$. Hence

$$\max_{1 \leq i \leq 4} \frac{(q, d_i) |t_i - m_i d'_i \eta q|}{d_i} \gg \frac{(q, d_1)}{d_1},$$

which, together with (3.21), gives

$$\begin{aligned} & \frac{q}{P} \sum_{\substack{|t_i - m_i d'_i q \eta| < \frac{M_i}{(q, d_i)} \\ 1 \leq i \leq 4}} \frac{1}{\max_{1 \leq i \leq 4} (q, d_i) |t_i - m_i d'_i q \eta| / d_i} \\ & \ll \frac{q}{P} \left(\frac{d_1 M_1 M_2 M_3 M_4}{(q, d_1)^2 (q, d_2) (q, d_3) (q, d_4)} + \frac{d_1 M_2 M_3 M_4}{(q, d_1) (q, d_2) (q, d_3) (q, d_4) \|m_1 d'_1 q \eta\|} \right) \\ & \ll \frac{q P^{\varepsilon-1} D d_1 d_2 d_3 d_4}{(q, d_1)^2 (q, d_2) (q, d_3) (q, d_4)} + \frac{q P^{\varepsilon-1} d_1 d_2 d_3 d_4}{(q, d_1) (q, d_2) (q, d_3) (q, d_4) \|m_1 d'_1 q \eta\|}. \end{aligned} \quad (3.22)$$

As η is quadratic irrationality, it has periodic continued fraction and if $\frac{a_n}{b_n}$, $n \in \mathbb{N}$ is the n -th convergent, then $b_n \leq c^n$ for some constant $c > 0$. Using that $\|m_1 d'_1 q\| \leq \frac{HDQ}{(d_1, q)}$ and Liouville's inequality for quadratic numbers (see Lemma 2.7), we can find convergent $\frac{a}{b}$ to η with denominator such that

$$\frac{3HDQ}{(d_1, q)} < b \ll_c \frac{HDQ}{(d_1, q)}. \quad (3.23)$$

Since $(a, b) = 1$ we have that $m_1 d'_1 q \frac{a}{b} \notin \mathbb{Z}$. As $\left| \eta - \frac{a}{b} \right| < \frac{1}{b^2}$ and (3.23) we get

$$\begin{aligned} \|m_1 d'_1 q \eta\| & \geq \left\| m_1 d'_1 q \frac{a}{b} \right\| - \left\| m_1 d'_1 q \left(\eta - \frac{a}{b} \right) \right\| \geq \left\| m_1 d'_1 q \frac{a}{b} \right\| - \frac{|m_1| d'_1 q}{b^2} \\ & > \frac{1}{b} - \frac{|m_1| d'_1 q (d_1, q)}{3bHDQ} \geq \frac{1}{b} - \frac{|m_1| d_1 q}{3bHDQ} \\ & > \frac{1}{b} - \frac{|m_1|}{3bH} \geq \frac{1}{b} - \frac{1}{3b} = \frac{2}{3b} \\ & \gg \frac{(d_1, q)}{HDQ}. \end{aligned}$$

From (3.21) and (3.22) it follows that

$$\sum_{\substack{|n_i - m_i d_i q \eta| < M_i \\ (q, d_i) | n_i, i=1, \dots, 4}} \frac{1}{|(\vec{m}\eta - \frac{\vec{n}}{d})P|} \ll \frac{q P^{\varepsilon-1} d_1 d_2 d_3 d_4 HDQ}{(q, d_1)^2 (q, d_2) (q, d_3) (q, d_4)}.$$

Then for F_1 (see (3.20)) we obtain

$$F_1 \ll \frac{P^{1+\varepsilon} DQ}{\delta} \sum_{q \leq Q} \frac{\tau(q) (q, N)^{1/2}}{q^{1/2}}. \quad (3.24)$$

Applying Cauchy's inequality we get

$$\begin{aligned}
 F_1 &\ll \frac{P^{1+\varepsilon}DQ}{\delta} \left(\sum_{q \leq Q} \tau^2(q) \right)^{\frac{1}{2}} \left(\sum_{q \leq Q} \frac{(q, N)}{q} \right)^{\frac{1}{2}} \\
 &\ll \frac{P^{1+\varepsilon}DQ}{\delta} \cdot Q^{1/2} (\log Q)^{3/2} \left(\sum_{\substack{t|N \\ t \leq Q}} \sum_{q_1 \leq \frac{Q}{t}} \frac{1}{q_1} \right)^{\frac{1}{2}} \\
 &\ll \frac{P^{1+\varepsilon}DQ^{3/2}}{\delta}.
 \end{aligned} \tag{3.25}$$

We choose $Q = \delta^{1/2}P^{1/2}D^{-1/2}$. Then

$$F_1, F_2 \ll P^{7/4+\varepsilon} \delta^{-1/4} D^{1/4}.$$

From (3.14), (3.15), (3.16), (3.17) it follows that

$$\Gamma_1^* \ll D^{17/4} P^{7/4+\varepsilon} \delta^{-1/4}.$$

The estimate of Γ_5^* goes along the same lines.

3.4. END OF THE PROOF OF THEOREM 1.1

From (3.10) and (3.11) we get

$$\Gamma \gg \frac{\delta N}{(\log N)^4} + D^{17/4} P^{7/4+\varepsilon} \delta^{-1/4}.$$

Then for a fixed small $\varepsilon > 0$, $\lambda < \frac{1-8\varepsilon}{10}$, $D < N^{\frac{1-10\lambda-8\varepsilon}{34}}$ and $z = D^{1/3,13}$ we get $\Gamma \gg \frac{\delta N}{(\log N)^4}$. So the equation (1.1) have solutions in almost-prime numbers $x_1, \dots, x_4 \in \mathcal{P}_k$, $k = \left\lceil \frac{53,21}{1-10\lambda-8\varepsilon} \right\rceil$ such that $\{\eta x_i\} < N^{-\lambda}$, $i = 1, 2, 3, 4$.

ACKNOWLEDGEMENTS. The authors thank Professor Doychin Tolev for his helpful comments and suggestions. Zhivko Petrov was partially supported by the Sofia University Research Fund through Grant 80-10-43/2020. Tatiana Todorova was partially supported by the Sofia University Research Fund through Grant 80-10-151/2020.

4. REFERENCES

- [1] Brüdern, J. and Fouvry, E.: Lagrange's Four Squares Theorem with almost prime variables. *J. Reine Angew. Math.* **454** (1994), 59–96.
- [2] Ching, Tak Wing: Lagrange's equation with almost-prime variables. *J. Number Theory* **212** (2020), 233–264.
- [3] Greaves, G.: *Sieves in number theory*, Springer, 2001.
- [4] Gritsenko, S. A. and Mot'kina, N. N.: Hua Loo Keng's problem involving primes of a special type. *IAP Tajikistan* **52**(7) (2009), 497–500.
- [5] Gritsenko, S. A. and Mot'kina, N. N.: On the solvability of Waring's equation involving natural numbers of a special type. *Chebyshevskii Sb.* **17**(1), (2016), 37–51.
- [6] Gritsenko, S. A. and Mot'kina, N. N.: Representation of natural numbers by sums of four squares of integers having a special form. *J. Math. Sci.* **173**(2) (2011), 194–200.
- [7] Gritsenko, S. A. and Mot'kina, N. N.: Waring's problem involving natural numbers of a special type. *Chebyshevskii Sb.* **15** (3), (2014), 31–47.
- [8] Hardy, G. H. and Wright, E. M.: *An introduction to the theory of numbers*, 5-th ed., Oxford Univ. Press, 1979.
- [9] Heath-Brown, D. R. and Tolev, D. I.: Lagrange's four squares theorem with one prime and three almost-prime variables. *J. Reine Angew. Math.* **558** (2003), 159–224.
- [10] Karatsuba, A. A.: *Basic analytic number theory*, Springer, 1993.
- [11] Schidt, M. W.: *Diophantine approximation*, Springer, 1980.
- [12] Shutov, A. V.: On the additive problem with fractional numbers. *BSUSB Series "Mathematics. Physics"* **30** **5** (148) (2013), 11–120.
- [13] Shutov, A. V. and Zhukova, A. A.: Binary additive problem with numbers of a special type. *Chebyshevskii Sb.* **16** (3) (2015), 246–275.
- [14] Todorova, T. L. and Tolev, D. I.: On the equation $x_1^2 + x_2^2 + x_3^2 + x_4^2 = N$ with variables such that $x_1 x_2 x_3 x_4 + 1$ is an almost-prime. *Tatra Mountains Math. Publications* **59**. **I** (2014), 1–26.
- [15] Tolev, D. I.: *Additive problems in number theory*. Doctoral Dissertation, Moscow University "M. Lomonosov", Moscow, 2001, (in Russian).

Received on December 17, 2020

ZHIVKO H. PETROV AND TATYANA L. TODOROVA

Faculty of Mathematics and Informatics

Sofia University St Kliment Ohridski

5 James Bourchier Blvd.

1164 Sofia

BULGARIA

E-mails: zhpetrov@fmi.uni-sofia.bg

tl@fmi.uni-sofia.bg

ГОДИШНИК НА СОФИЙСКИЯ УНИВЕРСИТЕТ СВ. КЛИМЕНТ ОХРИДСКИ “

ФАКУЛТЕТ ПО МАТЕМАТИКА И ИНФОРМАТИКА

Том 107

ANNUAL OF SOFIA UNIVERSITY ST. KLIMENT OHRIDSKI“

FACULTY OF MATHEMATICS AND INFORMATICS

Volume 107

REVIEW OF CONTINUUM MECHANICS AND ITS HISTORY
PART I. DEFORMATION AND STRESS. CONSERVATION LAWS.
CONSTITUTIVE EQUATIONS

BOGDANA A. GEORGIEVA

This is a review of continuum mechanics and its history, citing its original sources. It “bridges” the contributions of Bernoulli, Euler, Lagrange, Cauchy, Helmholtz, St Venant, Stokes, Fresnel, Cesaro, and others, written in a period of two centuries in 5 languages, in a coherent and historically accurate presentation in the contemporary notation. The only prerequisite knowledge to understand the paper is advanced calculus and elementary differential equations. Some valuable, but little known, results are reviewed in detail, like the exact solution of Cesaro to the system of differential equations which every continuous medium obeys, as well as his derivation of the conditions of St Venant for compatibility of the deformations. The last section presents the contemporary applications of continuum mechanics. The review continues with Part II. The Mechanics of Thermoelastic Media. Perfect Fluids, reference [45]. It discusses the consequences of Navier’s system of linear elasticity and approaches for its solution. It also gives a perspective of how waves propagate in continuous media. Reviewed are perfect fluids and linearly viscous fluids. At the end, Part II discusses the conditions for compatibility of the stresses.

Keywords: Mechanics of continuous media, continuum mechanics, history of continuum mechanics, elasticity, theory of elasticity.

2020 Math. Subject Classification: 74-02, 74B-05.

1. INTRODUCTION

Mechanics of continuous media is one of the classical branches of applied mathematics, which was built by several of the most prominent mathematicians of the 18th, 19th and the early 20th centuries. In addition to being a discipline

of its own, it is the heart of several modern branches of applied mathematics: fluid mechanics, gas dynamics, theory of elasticity, theory of deformable solids and others. Its applications penetrate almost every aspect of contemporary applied mathematics and mathematical physics. Over the centuries so much material accumulated in this subject, that at present only a few mathematicians know what is a fundamental notion in it and what is an application or a consequence of its core results. It is important that the mathematicians of today do know continuum mechanics not only for this knowledge itself, but also for the correct vision and proper sight of Mathematics and Science that it gives. It will help them size their own gauge to the contemporary needs of their profession. In addition to its powerful applications, continuum mechanics is precious for its esthetics - it is a part of the most elegant and sophisticated classical mathematics and reading it gives a pleasure and a professional growth.

The first attempt to discuss local features of the motion of a continuous medium in more than one dimension occurs in an isolated passage by D. Bernoulli from 1738 ([1], §11, paragraph 4). We are surrounded by matter in the form of continuous media – deformable solids, liquids and gasses. Let us begin at the moment of time $t = 0$ with a continuous medium, like a gallon of water, which we can easily imagine fills the volume V , with a shape specified by our imagination. Atomic structure is not considered. If the water is not held in a vessel, when we “unfreeze” time, it will move under the law of gravity and the laws of conservation of mass, momentum and energy, in a perfectly deterministic manner, continuously changing its shape, and eventually splash on the floor. This is a simple example of a motion of a continuous medium and is suitable to demonstrate what is meant by “material coordinates” and by “spatial coordinates”. **Material coordinates**, also called **Lagrangian coordinates**, are denoted by (X_1, X_2, X_3) and are the coordinates of the material points of the continuous medium at time $t = 0$. Lagrange introduced them in 1788 in [54], part II, section II. **Spatial coordinates**, also known as **Eulerian coordinates**, are denoted by (x_1, x_2, x_3) and are the coordinates of the points of 3-dimensional space (in which we observe the medium) occupied by the medium at time $t > 0$. Since the material coordinates are the coordinates of the material points at an arbitrary initial time $t = 0$, they can serve for all time as names for the *particles* of the material. The spatial coordinates, on the other hand, we think of as assigned once and for all to a point in the Euclidean space. They are the names of *places*. The motion $\mathbf{x} = \mathbf{x}(\mathbf{X}, t)$ chronicles the places \mathbf{x} occupied by the particle \mathbf{X} in the course of time. Under external influences - forces and heating - the continuous body deforms. *The goal of Mechanics of Continua is to find the family of transformations*

$$x_i = x_i(X_1, X_2, X_3, t), \quad i = 1, 2, 3, \quad (1)$$

giving the Eulerian coordinates as functions of the Lagrangian coordinates for $t \geq 0$.

This motion is perfectly deterministic, obeying only the natural laws, that we will present. We will arrive at a system of 20 partial differential equations for 20

unknown functions. This system is one of the finest triumphs of the symbiosis between mathematics and physics. We sketch the solution to this system and present the conditions for its existence and uniqueness. We give credit to the mathematicians and physicists who built this discipline by citing the date, name and the historical reference where the result was published for the first time.

The general theory of the motion of a continuous medium, which is understood of as a family of deformations continuously varying in time, is almost exclusively due to Euler, published in the period 1745 – 1766 in [25] – [41], and Cauchy, published in the period 1815 – 1841 in [3] – [18]. Important special results were added by D’Alembert in 1749 in [22], Green in 1839 in [46], Stokes in 1845 in [61], Helmholtz in 1858 in [48] and Cesaro in 1906 in [19].

2. STRAIN

The change in length and relative direction occasioned by the transformation is called **strain**. The term is due to Rankine [56] in 1851. Let us begin its study by defining the displacement vector \mathbf{u} , with components $u_i = x_i - X_i$, where $x_i = x_i(X_1, X_2, X_3, t)$ $i = 1, 2, 3$. The components u_i can be expressed in Lagrangian or in Eulerian coordinates, depending on need. Let P_0 be an arbitrary point of the continuous medium at time $t = 0$ and let Q_0 be a neighboring point, such that in a fixed Cartesian coordinate system $O_{e_1e_2e_3}$ P_0 has coordinates (X_1, X_2, X_3) , i.e. the radius vectors to P_0 is \mathbf{X} and to the point Q_0 is $\mathbf{X} + d\mathbf{X}$. At time $t > 0$ the material point P_0 occupies new geometric point P with coordinates (x_1, x_2, x_3) , i.e. P has radius vector \mathbf{x} and hence the new geometric location of the material point Q_0 is Q with a radius vector $\mathbf{x} + d\mathbf{x}$. To study the deformation that has occurred, we need to see how much has the distance between the two neighboring points P_0 and Q_0 changed. For that we calculate

$$(d\mathbf{x})^2 - (d\mathbf{X})^2 = \left(\frac{\partial x_k}{\partial X_i} \frac{\partial x_k}{\partial X_j} - \delta_{ij} \right) dX_i dX_j = \left(\delta_{ij} - \frac{\partial X_k}{\partial x_i} \frac{\partial X_k}{\partial x_j} \right) dx_i dx_j. \quad (2)$$

Here and throughout the paper each index takes the values 1, 2 and 3 and the summation convention on repeated indexes is assumed. We see that all the information about the deformation is contained in the coefficients of $dX_i dX_j$ and respectively of $dx_i dx_j$ in (2). These sets of coefficients

$$E_{ij} \equiv \frac{1}{2} \left(\frac{\partial x_k}{\partial X_i} \frac{\partial x_k}{\partial X_j} - \delta_{ij} \right) \quad \text{and} \quad e_{ij} \equiv \frac{1}{2} \left(\delta_{ij} - \frac{\partial X_k}{\partial x_i} \frac{\partial X_k}{\partial x_j} \right)$$

satisfy the transformation laws for tensors of rang 2 and are called the Lagrangian and the Eulerian **tensors of finite deformations** or **finite strain tensors**. The difference $(d\mathbf{x})^2 - (d\mathbf{X})^2$ is a measure for the size of the deformation in the vicinity of P_0 . Because dX_i and dx_i are arbitrary, the necessary and sufficient condition this difference to be 0 is $E_{ij} = 0$ or equivalently $e_{ij} = 0$. In that case the deformation

near that point is 0 and the motion is that of a rigid body. Written in terms of the gradients $\partial u_i/\partial X_j$ or $\partial u_i/\partial x_j$ of the displacement vector \mathbf{u} , E_{ij} and e_{ij} are

$$E_{ij} = \frac{1}{2} \left(\frac{\partial u_i}{\partial X_j} + \frac{\partial u_j}{\partial X_i} + \frac{\partial u_k}{\partial X_i} \frac{\partial u_k}{\partial X_j} \right) \quad \text{and} \quad e_{ij} = \frac{1}{2} \left(\frac{\partial u_i}{\partial x_j} + \frac{\partial u_j}{\partial x_i} - \frac{\partial u_k}{\partial x_i} \frac{\partial u_k}{\partial x_j} \right).$$

We will now make a crucial assumption - that the deformations which we will study are small. This means that the gradients $\partial u_i/\partial X_j$ and $\partial u_i/\partial x_j$ of the displacement \mathbf{u} are small in comparison to 1, and hence the products of these gradients may be ignored in the presence of the gradients themselves. In this manner we obtain the tensors \bar{E}_{ij} and \bar{e}_{ij} . A calculation based on the same assumption shows that they are equal and we give them the common name ε_{ij} . This is the **tensor of (infinitesimal) deformations** or the (infinitesimal) **strain tensor**

$$\varepsilon_{ij} \equiv \frac{1}{2} \left(\frac{\partial u_i}{\partial x_j} + \frac{\partial u_j}{\partial x_i} \right).$$

The strain tensor ε_{ij} was introduced by Green in 1841 in [47] and by St Venant in 1844 in [59]. It is the most popular strain measure even today. The vanishing of ε_{ij} is necessary and sufficient for a rigid displacement. The general deformation $dX \rightarrow dx$ as well as the displacement gradients $\partial u_i/\partial X_j$ and $\partial u_i/\partial x_j$ as measures of local changes of length and angle are due to Lagrange 1762, [53] §XLIV and 1788 [54] Part II, Sect.11. The fully general spatial description is due to Euler, dates 1752, and was first published in 1757 in [31] and then in 1761 in [33]. The theory of finite strain is the creation of Cauchy published in 1823 [4], in 1827 [7] and in 1841 [18]. The theory of infinitesimal strain was first developed by Euler. It was fully elaborated by Cauchy, who obtained it by specialization from his general theory of finite strain.

We will now explain the geometry of the process of deformation. The component ε_{11} of the strain tensor is the relative elongation of a linear element in the direction of the unit coordinate vector \mathbf{e}_1 , and similarly for ε_{22} and ε_{33} . The component ε_{23} is half of the change (as a result of the deformation) of the angle between two lines, that initially had the directions of the unit coordinate vectors \mathbf{e}_2 and \mathbf{e}_3 . Even more surprising is the fact that, at each point inside the deforming medium, the deformations can not take an arbitrary shape. Instead, they form quadratic surfaces only, called **surfaces of Cauchy**. This is not hard to see and is worth the effort. Let us denote by ε the relative elongation in direction of the vector $d\mathbf{X}$, with length dX

$$\varepsilon \equiv \frac{dx - dX}{dX}.$$

Consider the difference

$$(dx)^2 - (dX)^2 = (dx - dX)(dx + dX) = 2\varepsilon_{ij} dX_i dX_j, \quad (3)$$

and observe the smallness of the deformations, i.e. that $dx \approx dX$. Then by dividing both sides of (3) by $dX dX$ we see that

$$\varepsilon = \varepsilon_{ij} \frac{dX_i}{dX} \frac{dX_j}{dX}. \quad (4)$$

Hence for any vector with components (ξ_1, ξ_2, ξ_3) and magnitude ξ , the last formula (4) gives $\xi^2 \varepsilon = \varepsilon_{ij} \xi_i \xi_j$. For each direction we can select ξ in such a way that $\xi^2 \varepsilon = \pm k^2$, where k is a positive constant and the sign is chosen so that the square of the length of vector ξ to be positive. It follows that at any point of the deforming medium the strain takes the shape of the **quadratic surface**

$$\varepsilon_{ij} \xi_i \xi_j = \pm k^2$$

called **surface of deformations of Cauchy at the point P_0** . From this geometric picture it is clear that the elongation ε in the direction of the vector ξ is inversely proportional to the square of the distance from the center of the surface (the point P_0) to the intersection of the vector ξ with that surface.

Because the vector $(\varepsilon_{1j} \xi_j, \varepsilon_{2j} \xi_j, \varepsilon_{3j} \xi_j)$ is normal to the quadratic surface of Cauchy, we see that the relative displacement at P_0 due to the pure deformation is in the direction of the normal to that surface at the point of intersection of the surface with this vector.

After these observations, it is plausible to seek lines through P_0 with directions that do not change under pure deformation. Of course, these are the lines along the eigenvectors of the strain tensor ε_{ij} . It is symmetric and hence has 3 real eigenvalues, called **main deformations** or **Cauchy principal stretches**, ε_I , ε_{II} , and ε_{III} . To each of them corresponds an eigenvector, called **main direction** or **main axis of the strain tensor**. Cauchy published these results first in 1823 [4] and again in 1827 [7]. To different main deformations correspond main directions that are orthogonal. We can select the axes of the coordinate system to coincide with the main axes of the tensor of deformations and, as a result, obtain the simplest form of the quadratic surface of Cauchy

$$\varepsilon_I \xi_1^2 + \varepsilon_{II} \xi_2^2 + \varepsilon_{III} \xi_3^2 = \pm k^2.$$

The invariants of the tensors \mathbf{E} and \mathbf{e} were first published by Cauchy in 1827 in [7].

3. CONDITIONS FOR COMPATIBILITY OF THE DEFORMATIONS

Common sense tells us that the deformations that take place in a medium are not independent of each other. If we stretch an elastic membrane with a rectangular shape along one of its diagonals, the other diagonal will shrink. St. Venant proved in 1860 that in order for the six functions $\varepsilon_{ij}(x_1, x_2, x_3)$ to adequately define the components of the tensor of deformations ε_{ij} , so that the 6 partial differential equations

$$u_{i,j} + u_{j,i} = 2\varepsilon_{ij} \tag{5}$$

have a unique solution $\mathbf{u}(x_1, x_2, x_3)$, they must satisfy the system of 6 PDEs

$$\varepsilon_{ij,kl} + \varepsilon_{kl,ij} - \varepsilon_{ik,jl} - \varepsilon_{jl,ik} = 0. \tag{6}$$

The notation $, j$ denotes partial differentiation with respect to x_j . The 6 restrictions (6) on the components ε_{ij} of the tensor of deformations are called **conditions for compatability of the deformations** and their fulfillment is a necessary and sufficient condition for the existence of the solution vector \mathbf{u} to the system (5), which of course, has the physical meaning of the displacement vector $\mathbf{u} \equiv \mathbf{x} - \mathbf{X}$ in the process (1) of the deformation of the continuous medium. The derivation of the compatability conditions is exceptionally original. On the way of deriving the compatability conditions, an analytic formula for the displacement \mathbf{u} itself is derived, thus obtaining a result of even greater significance. Due to lack of space, this derivation is not presented here, but it is sketched. This method of obtaining the displacement \mathbf{u} is due to E. Cesaro [19], who published it in 1906. Volterra presents it in [62], citing Cesaro. Contemporary references on it are Ivanov [49] and Sokolnikoff [58]. The solution to (5) has components

$$u_j = u_j^0 + \omega_{jk}^0(x_k - x_k^0) + \int_{P_0}^P (\varepsilon_{jl} + (x_k - y_k)(\varepsilon_{jl,k} - \varepsilon_{kl,j})) dy_l, \quad j = 1, 2, 3. \quad (7)$$

Here

$$\omega_{ij} = \frac{1}{2} \left(\frac{\partial u_i}{\partial x_j} - \frac{\partial u_j}{\partial x_i} \right)$$

is the tensor of small rotations, introduced by Euler in 1761, §§46-47. In the components of the exact solution (7) of Cesaro u_j^0 are the components of the translation and ω_{jk}^0 are those of the tensor of rotation in an arbitrary point P_0 of the deforming body, and are assumed known. The first term in the solution (7) for u_j represents the translation and the second term represents the rotation of the continuous medium as a rigid body. The third term in u_j represents its deformation. Because the displacement \mathbf{u} is unique, its components u_j must not depend on the path of integration, so the integrands of the 3 integrals must be total differentials. Demanding this, yields the 6 equations (6) of St Venant for compatability of the deformations.

The compatability conditions were first published by Kirchhoff in 1859 in [52], but without a statement of their meaning, which was first explained by St Venant in his memoir [60]. St Venant obtained these conditions in a different way, than the one presented in this section. Submitted them to Scocietà Philomathique in 1860, who published them in 1864.

4. STRESS

The notion of stress arose in special case studies of theories of flexible, elastic and fluid bodies. Galileo (1638), Pardies (1673), James Bernoulli (1691-1704), Hermann (1716), Coulomb (1776), John Bernoulli (1743), and Euler (1749-1752) published studies on this notion. The general concept and mathematical theory are due to Cauchy, published in 1823 [4] and in 1827 [7]. Cauchy achieved the general

theory of stress by adopting the common features and discarding the special aspects of the foregoing theories. The term **stress** was introduced by Rankine in 1856 in [57].

The field of stress vectors is not an ordinary vector field. Rather, since the stress vectors across two different surfaces through the same point are generally different, at any given time, the stress vectors $\sigma(\mathbf{x}, t, \mathbf{n})$ depend both on the position vector \mathbf{x} and on the direction \mathbf{n} of the normal to the surface. We wish to extract all the information about the stress at a point of the body into a single mathematical object, and separate it from the information about the direction. This is accomplished by the **stress tensor** σ_{ij} .

To derive the components of that tensor we take a tetrahedron having 3 edges coming out of an arbitrarily fixed point P , parallel to the coordinate axes. The force acting on the medium occupying the volume V of the tetrahedron is $\int_V \rho \mathbf{f} dV$, where $\rho(\mathbf{x}, t)$ is the mass density and $\mathbf{f}(\mathbf{x}, t)$ is the mass force acting on ρdV . Examples of mass forces are gravity and the centrifugal force in a rotating body. Surface forces act on every surface inside the medium or on its surface. Those forces are modeled with the stress vector $\sigma(\mathbf{x}, t, \mathbf{n})$. The force acting on a portion S of a surface is $\int_S \sigma dS$. The orientation of S is given by the outward unit normal $\mathbf{n}(\mathbf{x}, t) = n_i(\mathbf{x}, t) \mathbf{e}_i$ to the surface at that point. (The dimension of the vector σ is pressure.)

We assume that all forces acting on the tetrahedron balance out

$$\sum_{j=1}^3 \int_{\Delta S_j} \sigma(\mathbf{x}, t, -\mathbf{e}_j) dS + \int_{\Delta S} \sigma(\mathbf{x}, t, \mathbf{n}) dS + \int_{\Delta V} \rho \mathbf{f} dV = \mathbf{0}, \quad (8)$$

where ΔS_j is the face perpendicular to \mathbf{e}_j , ΔS is the fourth face and ΔV is the part of 3-space occupied by the tetrahedron. We make use of the mean-value theorem in equation (8). Denote the radius-vector to the point P by \mathbf{x} , make use of $\Delta S_j = \Delta S \cos(\mathbf{n}, \mathbf{e}_j) = \Delta S n_j$, $\Delta V = h\Delta S/3$, and let the altitude h from P approach 0. We get

$$\sigma(\mathbf{x}, t, -\mathbf{e}_j) n_j + \sigma(\mathbf{x}, t, \mathbf{n}) = \mathbf{0}. \quad (9)$$

If we now denote by $\sigma_{ij}(\mathbf{x}, t)$ the components of the stress vector with a normal \mathbf{e}_j , $\sigma_{ij}(\mathbf{x}, t) = \sigma_i(\mathbf{x}, t, \mathbf{e}_j)$, from the last vector equation (9) we get

$$\sigma_i(\mathbf{x}, t, \mathbf{n}) = \sigma_{ij}(\mathbf{x}, t) n_j.$$

This important result is **Cauchy's fundamental theorem** and expresses the relationship between the components of the stress vector and the components of the stress tensor. All the information about the stress at a point is "extracted" in the stress tensor itself, and is "separated" from the orientation \mathbf{n} of the surface. Cauchy published this formula in 1823 [4] and in 1827 [6].

The geometry of the stress at a point of a deforming medium is also that of quadratic surfaces. Consider the stress vector σ , acting on a surface element with a

unit normal \mathbf{n} at a fixed point P of the body. Its components are $\sigma_i = \sigma_{ij} n_j$. Let us denote by σ_N the magnitude of its projection on \mathbf{n} . σ_N is called **normal stress** and can be expressed as

$$\sigma_N = \sigma_i n_i = \sigma_{ij} n_i n_j .$$

If ξ is a vector having the direction of the unit normal \mathbf{n} and size ξ , then from the last equation follows that $\xi^2 \sigma_N = \sigma_{ij} \xi_i \xi_j$, where ξ_i are the components of ξ . Select the size ξ of the vector ξ in such a way that $\xi^2 \sigma_N = \pm k^2$, where k is a fixed positive constant and the sign is chosen so that the length of ξ defined with this equation be positive. Then the “tip” of an *arbitrary* vector ξ with base at P , and magnitude ξ satisfying $\xi^2 \sigma_N = \pm k^2$, lies on the surface

$$\sigma_{ij} \xi_i \xi_j = \pm k^2$$

called **quadratic surface of the stress tensor** or **surface of Cauchy of the stress** at the point P . The stress tensor is symmetric and hence has 3 real eigenvalues, called **main stresses**. The corresponding eigenvectors are called **main directions** or **main axes**. If we choose a coordinate system with coordinate axes along the main axes of the stress tensor, the quadratic surface of the stress at the point acquires the form

$$\sigma_I \xi_1^2 + \sigma_{II} \xi_2^2 + \sigma_{III} \xi_3^2 = \pm k^2 ,$$

where $\sigma_I, \sigma_{II}, \sigma_{III}$ are the main stresses of σ_{ij} at that point. At a surface element with a normal \mathbf{n} along a main axes of the stress tensor, the stress vector σ has the direction of the normal.

5. CONSERVATION OF MASS, MOMENTUM AND MOMENT OF MOMENTUM

In contemporary mathematics and mathematical physics conservation laws are a main goal of study. Researchers obtain them from variational principles via the famous first theorem of Emmy Noether. In Mechanics of Continua, however, history went differently. All the laws of conservation, namely the conservation of mass, energy, momentum, and moment of momentum, were discovered by judicious guessing and verification with the physical experiment. They are all empirical laws. Much later they were derived from deliberately calculated for this purpose Lagrangians.

The law of conservation of mass is the statement that the mass, contained in any portion of the body with volume V , does not change during the deformation

$$\frac{d}{dt} \int_V \rho dV = 0 .$$

This can be rewritten as $\int_V \partial \rho / \partial t dV + \int_S v_n \rho dS = 0$, where $v_n = \mathbf{v} \cdot \mathbf{n}$ is the component of the velocity of the points on the surface S of V along the outward unit

normal \mathbf{n} to S . Thus, $\int_V (\frac{\partial \rho}{\partial t} + (\rho v_i)_{,i}) dV = 0$ where $v_i(\mathbf{x}, t)$ are the components of the velocity. If the integrand is continuous, we obtain the differential form of the **law of conservation of mass**

$$\frac{\partial \rho}{\partial t} + (\rho v_i)_{,i} = 0. \quad (10)$$

The law of conservation of mass was first discovered by Euler in 1757, reference [31], §§16-17.

In mechanics of continua the so-called **equations of motion** play the same role as do the equations of Newton in mechanics of rigid bodies. These equations of motion of a continuous medium follow from the law of **conservation of momentum**, which states that “The total time derivative of the momentum of an arbitrarily fixed portion of the deforming body is equal to the sum of all forces (mass forces $\mathbf{f}(\mathbf{x}, t)$ and surface forces $\sigma(\mathbf{x}, t)$) that act on it”

$$\frac{d}{dt} \int_V \rho v_i dV = \int_V \rho f_i dV + \int_S \sigma_i dS, \quad i = 1, 2, 3. \quad (11)$$

A simple calculation shows that, if mass is conserved, for any continuously differentiable function $g(x, t)$ it is true that

$$\frac{d}{dt} \int_V \rho g(x, t) dV = \int_V \rho \frac{dg}{dt} dV.$$

With $g = v_i$ this formula simplifies the law of conservation of momentum (11) to

$$\int_V \rho \frac{dv_i}{dt} dV = \int_V \rho f_i dV + \int_S \sigma_{ij} n_j dS. \quad (12)$$

Applying Gauss’ theorem to the surface integral in (11), combining the resulting 2 integrals, and assuming continuity, we obtain the **equations of motion of a continuous medium**

$$\sigma_{ij,j} + \rho f_i = \rho \frac{dv_i}{dt}, \quad i = 1, 2, 3. \quad (13)$$

These equations were first published by Cauchy in 1827 in [9], and also in 1827 in [11].

The **law of conservation of moment of momentum** asserts that “the time rate of change of the moment of momentum is equal to the sum of the moments of the mass forces and the surface forces that act on the body”, i.e.,

$$\frac{d}{dt} \int_V \rho e_{ijk} x_j v_k dV = \int_V \rho e_{ijk} x_j f_k dV + \int_S e_{ijk} x_j \sigma_k dS,$$

where the moments are written with respect to the origin of the coordinate system.

The laws of conservation of momentum and of moment of momentum are both due to Euler and were introduced by him in 1775, [43], §§26–28. While the memoire

is about rigid bodies, these two laws are expressly stated to hold for any continuous medium.

The law of conservation of moment of momentum is fully equivalent to the symmetry of the stress tensor

$$\sigma_{ij} = \sigma_{ji} .$$

This important result is known as **Cauchy's fundamental theorem**, and was published by him in 1827, [6]. It was discovered (but not published) by Fresnel in 1822, who published it in 1868, [44].

6. CONSERVATION OF ENERGY

In mechanics of rigid bodies thermal effects and thermal consequences of the motion are either considered separately from the equations of motion or completely ignored, if they do not affect the motion in consideration. For example, we ignore the heat generated during the friction between the surface of a cube sliding on a plane and that plane. In Mechanics of Continua heat generation and thermal effects can not be ignored or even considered separately from the equations of motion. The reason is that when a deformation takes place, heat is generated/lost throughout the entire volume where the deformation occurs. This thermal energy affects significantly the motion and the deformation. It becomes a cycle: the deformation generates heat and that heat in turn affects the distance between the particles of the continuous medium, thus causing deformation. The dynamics of a continuous medium and the thermal laws are intertwined and must be studied simultaneously.

That heat is a mode of motion was widely believed in the 18th century. Both Daniel Bernoulli [1] in 1738 and Euler [35] in 1765 constructed kinetic molecular models in which temperature may be identified with the kinetic energy of the molecules. The general and phenomenological principle, independent of molecular interpretation, was known to Carnot by 1824, as proved by his memoir [2]. The first clear statement of the interconvertibility of heat and mechanical work, that any equation of energy balance should contain terms that represent non-mechanical transfer of energy, are those of Joule [50], [51] from 1843 and 1845 and of Waterston [64] from 1843.

Let us now consider the **law of conservation of energy**. It states that "The total time derivative of the sum of the kinetic energy and the internal energy is equal to the sum of the power of the external forces and the in-flow of all other kinds of energies per unit of time"

$$\frac{dK}{dt} + \frac{dE}{dt} = W + Q, \tag{14}$$

where $K = \int_V \rho v_i v_i / 2 dV$ is the kinetic energy, $W = \int_V \rho f_i v_i dV + \int_S \sigma_i v_i dS$ is the power of the external forces, $Q = - \int_S q_i n_i dS + \int_V \rho r dV$ is the in-flow

of heat per unit of time. Here $\mathbf{q} = q_i(\mathbf{x}, t) e_i$ is the vector of heat flow and $r(\mathbf{x}, t)$ is the specific heat source. For simplicity, we assume that there is only in-flow of thermal energy. We also assume the existence of a function $\epsilon(\mathbf{x}, t)$ called **specific internal energy** such that

$$\int_{V(t)} \rho \epsilon dV = E,$$

where E is the total internal energy of the part of the body with volume V at time t . The general law of conservation of energy (when heat effects are included), i.e. equation (14), is called “**the first law of thermodynamics**”. The first one to formulate this important law was Duhem [24], Chapter III, §3, in 1892 .

In the special case $Q = 0$ the first law of thermodynamics reduces to the **law of conservation of mechanical energy**

$$\frac{dK}{dt} + \int_V \sigma_{ij} d_{ij} dV = W,$$

where

$$d_{ij} \equiv \frac{1}{2}(v_{i,j} + v_{j,i}) = d_{ji}$$

is the tensor of rate of deformations, introduced by Euler [41], §§ 9–12, in 1769. By a simple, but tedious calculation, substituting dK/dt , E and Q into the general law of conservation of energy (14), transforming the surface integral into a volume integral, and assuming continuity, we obtain the **differential form of the general law of conservation of energy**

$$\rho \frac{d\epsilon}{dt} = \sigma_{ij} d_{ij} - q_{i,i} + \rho r. \quad (15)$$

That use of a differential equation expressing balance of energy is necessary, except in specially simple circumstance, was first emphasized by Duhem [23], Vol. I, Livre II, Chapter III, in 1891. In 1769 Euler [41], §13, showed that the vanishing of all components of the tensor of rate of deformations is the criterion for a rigid motion.

7. ENTROPY

In the present section we define and explain the concept of entropy and the second law of thermodynamics.

Let us begin with some history. During the Industrial Revolution in Western Europe, it was observed that the steam engines of locomotives and other engines that transform thermal energy into mechanical energy can not achieve efficiency of 100%. In 1865 Rudolf Clausius [21], §14, introduced the concept of entropy for the lost thermal energy in steam engines, i.e., the heat which remained unconverted into mechanical energy. **Entropy** is defined by

$$d\eta = c \frac{d\theta}{\theta}$$

where $\eta(\mathbf{x}, t)$ is the entropy for unit mass, c is the specific heat and $\theta(\mathbf{x}, t)$ is the absolute temperature of the body.

The **inequality of Clausius - Duhem** is

$$\frac{d}{dt} \int_V \rho \eta \, dV \geq - \int_S \frac{q_i}{\theta} n_i \, dS + \int_V \rho \frac{r}{\theta} \, dV,$$

where $r(\mathbf{x}, t)$ is the specific heat source, and $\mathbf{q} = q_i \mathbf{e}_i$ is the vector of heat flow. It has the direction of motion of heat. The normal \mathbf{n} is outward to the surface S . The first integral in the right hand side is the flow of entropy per unit time through the surface S of the volume V and the second integral is the creation of entropy inside V by outside sources per unit time. This inequality is one of the fundamental empirical laws of thermodynamics – the second law of thermodynamics. It is due to Clausius [20] (1854). The meaning of the second law of thermodynamics is can be explained as follows. It is known from experience that a substance at uniform temperature and free fro sources of heat may consume mechanical work, but can not give it out. That is, whatever work is not recoverable is lost, not created. Also, in a body at rest and subject to no sources of heat, the flow of heat is from the hotter to the colder parts, not vice versa.

Using the well known formula

$$\frac{d}{dt} \int_V \rho f \, dV = \int_V \rho \frac{df}{dt} \, dV,$$

where ρ is the mass density, which holds for any continuously differentiable function $f(\mathbf{x}, t)$, we obtain the differential form of the inequality of Clausius - Duhem:

$$\rho \frac{d\eta}{dt} + \left(\frac{q_i}{\theta} \right)_{,i} - \rho \frac{r}{\theta} \geq 0. \quad (16)$$

8. CONSTITUTIVE EQUATIONS

We consider the differential forms of: the law of conservation of mass (10), the law of conservation of energy (15), the equations of motion of a continuous medium (13), and the inequality of Clausius-Duhem (16) as a system. These are 5 scalar differential equations and 1 inequality for the 16 unknown functions u_i , ρ , σ_{ij} , ϵ , η and θ . We take in consideration the symmetry of the stress tensor $\sigma_{ij} = \sigma_{ji}$, the definition of $d_{ij} = (v_{i,j} + v_{j,i})/2$, and assume that \mathbf{f} and r are given. It is remarkable, but not surprising, that physics provides the additional equations necessary to solve this system. These are the so called **constitutive equations** and contain information about the specific material of the medium. An elastic is very different from water, which is very different from an oil or a gas. The constitutive equations characterize the mechanical and thermal properties of the medium.

In experiments and observations, the motion of the material particles of the continuous medium and its temperature can be observed and measured, so from mathematical stand point the components u_i of the displacement vector, the temperature θ , as well as their derivatives, will be the independent variables in the constitutive equations, which we are trying to build. All of the rest of the variables will dependent on these ones and will be dependent variables. These are: σ_{ij} , ϵ , q_i and η , a total of 11 such variables. The mass density ρ is also a dependent variable. For it we already have a differential equation, relating it to the rest of the variables, namely the law of conservation of mass.

Because the constitutive equations characterize the properties of the materials, they must remain invariant under a rotation or a translation. This requirement is met if the variables (both independent and dependent), which those equations relate, are themselves independent of such transformations. It is easy to show that such variables are:

$$\Sigma_{kl} = \sigma_{ij} \frac{\partial x_i}{\partial X_k} \frac{\partial x_j}{\partial X_l}, \quad Q_j = q_i \frac{\partial x_i}{\partial X_j}$$

as well as the scalar functions ϵ and η . Thus, in the constitutive equations which we are trying to construct, it will be reasonable to regard as independent variables the temperature θ , the coordinates X_i , the gradient $\delta\theta/\delta X_i$ of the temperature and the tensor of deformations E_{ij} . Hence for a thermoelastic medium the constitutive equations are :

$$\Sigma_{ij} = \Sigma_{ij}(\mathbf{E}, \theta, \mathbf{G}, \mathbf{X}), \quad \mathbf{Q}_i = \mathbf{Q}_i(\mathbf{E}, \theta, \mathbf{G}, \mathbf{X}), \quad \epsilon = \epsilon(\mathbf{E}, \theta, \mathbf{G}, \mathbf{X}), \quad \eta = \eta(\mathbf{E}, \theta, \mathbf{G}, \mathbf{X}),$$

where \mathbf{G} denotes the gradient of the temperature with respect to the Lagrangian coordinates X_i , $i = 1, 2, 3$. Using

$$(\partial x_i / \partial X_k)(\partial X_k / \partial x_j) = \delta_{ij}, \quad (\partial X_i / \partial x_k)(\partial x_k / \partial X_j) = \delta_{ij},$$

we invert the equations for Σ_{kl} and Q_j to obtain

$$\begin{aligned} \sigma_{ij} &= \Sigma_{kl}(\mathbf{E}, \theta, \mathbf{G}, \mathbf{X}) X_{k,i} X_{l,j}, & q_i &= Q_j(\mathbf{E}, \theta, \mathbf{G}, \mathbf{X}) X_{j,i}, \\ \epsilon &= \epsilon(\mathbf{E}, \theta, \mathbf{G}, \mathbf{X}), & \eta &= \eta(\mathbf{E}, \theta, \mathbf{G}, \mathbf{X}). \end{aligned}$$

Using the inequality of Clausius-Duhem we will be able to see the form of the constitutive equations in more detail. For this, a new function, **free energy**, is introduced:

$$\psi \equiv \epsilon - \eta \theta.$$

Obviously $\psi = \psi(\mathbf{E}, \theta, \mathbf{G}, \mathbf{X})$ and we assume that it is symmetric with respect to E_{ij} and E_{ji} . (This is possible, because $E_{ij} = E_{ji}$ and so we can replace E_{ij} and E_{ji} in ψ with $(E_{ij} + E_{ji})/2$.) By elementary mathematical manipulations we eliminate r from the inequality of Clausius-Duhem (16) to obtain

$$-\rho \frac{d\psi}{dt} - \rho \eta \frac{d\theta}{dt} + \sigma_{ij} d_{ij} - \frac{q_i \theta_{,i}}{\theta} \geq 0. \quad (17)$$

Substituting ψ in inequality (17), we get

$$-\rho \frac{\partial \psi}{\partial E_{ij}} \frac{\partial E_{ij}}{\partial t} - \rho \frac{\partial \psi}{\partial \theta} \frac{d\theta}{dt} - \rho \frac{\partial \psi}{\partial G_i} \frac{dG_i}{dt} - \rho \eta \frac{d\theta}{dt} + \sigma_{ij} d_{ij} - \frac{q_i \theta_{,i}}{\theta} \geq 0. \quad (18)$$

Let us now do the calculation

$$\begin{aligned} \frac{\partial E_{ij}}{\partial t} &= \frac{1}{2} \frac{\partial}{\partial t} \left(\frac{\partial x_k}{\partial X_i} \frac{\partial x_k}{\partial X_j} - \delta_{ij} \right) = \frac{1}{2} \left(\frac{\partial v_l}{\partial X_i} \frac{\partial x_l}{\partial X_j} + \frac{\partial x_k}{\partial X_i} \frac{\partial v_k}{\partial X_j} \right) \\ &= \frac{1}{2} \left(\frac{\partial v_l}{\partial x_k} \frac{\partial x_k}{\partial X_i} \frac{\partial x_l}{\partial X_j} + \frac{\partial x_k}{\partial X_i} \frac{\partial v_k}{\partial x_l} \frac{\partial x_l}{\partial X_j} \right) = d_{kl} \frac{\partial x_k}{\partial X_i} \frac{\partial x_l}{\partial X_j} = \frac{d\varepsilon_{kl}}{dt} \frac{\partial x_k}{\partial X_i} \frac{\partial x_l}{\partial X_j}. \end{aligned}$$

We substitute this result in the last inequality (18) to obtain

$$\left(\sigma_{kl} - \rho \frac{\partial \psi}{\partial E_{ij}} \frac{\partial x_k}{\partial X_i} \frac{\partial x_l}{\partial X_j} \right) \frac{d\varepsilon_{kl}}{dt} - \rho \left(\eta + \frac{\partial \psi}{\partial \theta} \right) \frac{d\theta}{dt} - \rho \frac{\partial \psi}{\partial G_i} \frac{dG_i}{dt} - \frac{q_i \theta_{,i}}{\theta} \geq 0. \quad (19)$$

The inequality (19) is linear with respect to the three variables $d\varepsilon_{kl}/dt$, $d\theta/dt$ and dG_i/dt with coefficients which do not depend on them. Because $d\varepsilon_{kl}/dt$, $d\theta/dt$ and dG_i/dt are independent of each other (since u , θ and their gradients at an arbitrary point are independent variables), it follows that a necessary and sufficient condition for inequality (19) to hold is that the coefficients of these three variables are zeros. Thus,

$$\sigma_{kl} = \rho \frac{\partial \psi}{\partial E_{ij}} \frac{\partial x_k}{\partial X_i} \frac{\partial x_l}{\partial X_j}, \quad \eta = -\frac{\partial \psi}{\partial \theta}, \quad \frac{\partial \psi}{\partial G_i} = 0, \quad q_i \theta_{,i} \leq 0.$$

Hence ψ does not depend on G_i , i.e. $\psi = \psi(\mathbf{E}, \theta, \mathbf{X})$. Traditionally, the left hand side of the inequality $q_i \theta_{,i} \leq 0$ is written as

$$q_i \theta_{,i} = Q_j X_{j,i} \frac{\partial \theta}{\partial X_k} X_{k,i} = Q_j(\mathbf{E}, \theta, \mathbf{G}, \mathbf{X}) G_k X_{j,i} X_{k,i}.$$

Let us summarize what we have accomplished in this section. To the original system of 5 differential equations for the 16 unknown functions, stated in the beginning of the section, we added 7 new unknowns (E_{ij} and ψ) and their defining equations

$$E_{ij} \equiv \frac{1}{2} \left(\frac{\partial u_i}{\partial X_j} + \frac{\partial u_j}{\partial X_i} + \frac{\partial u_k}{\partial X_i} \frac{\partial u_k}{\partial X_j} \right), \quad \psi \equiv \epsilon - \eta \theta,$$

and also added 7 equations – for σ_{ij} and η . So we have a total of 19 equations for 23 unknowns and the inequality $q_i \theta_{,i} \leq 0$. Thus, we need 4 more equations. These are the equations that specify the nature of the free energy $\psi = \psi(\mathbf{E}, \theta, \mathbf{X})$ and that of the heat flow $\mathbf{q} = \mathbf{q}(\mathbf{E}, \theta, \mathbf{G}, \mathbf{X})$.

For historical references on the constitutive equations of continuous media we refer the reader to Truesdell and Toupin [63].

9. VISCOELASTIC MEDIUM

We are interested in deriving the equations of motion of a viscoelastic medium. We use a “dot” above a letter to denote the time derivative of the variable.

Let us assume that the continuous medium we consider has a constant density ρ , constant temperature θ and constant entropy η . Let us also assume that the stresses depend not only on the deformations, but also on the time derivatives of the deformations, namely that

$$\Sigma_{ij} = \Sigma_{ij}(\mathbf{E}, \dot{\mathbf{E}}, \mathbf{X}).$$

The specific internal energy ϵ depends on the same variables.

Assuming that the deformations are small, the formula

$$\sigma_{ij} = \rho \frac{\partial \psi}{\partial E_{kl}} \frac{\partial x_i}{\partial X_k} \frac{\partial x_j}{\partial X_l}$$

derived above, which is valid for any continuous medium even in the case of large deformations and with no restrictions on the form that the free energy ψ , acquires the form

$$\sigma_{ij} = \rho \frac{\partial \psi}{\partial \varepsilon_{ij}}.$$

Let us assume that the free energy ψ is a quadratic function of the deformations and their time derivatives, namely,

$$\rho\psi = a + \alpha_{ij}\varepsilon_{ij} + \frac{1}{2}c_{ijkl}\varepsilon_{ij}\varepsilon_{kl} + \beta_{ij}\dot{\varepsilon}_{ij} + \frac{1}{2}\beta_{ijkl}\varepsilon_{ij}\dot{\varepsilon}_{kl} + \frac{1}{2}\gamma_{ijkl}\dot{\varepsilon}_{ij}\dot{\varepsilon}_{kl}.$$

Thus we arrive at the system of equations which an elastic medium with viscosity, a constant density ρ , constant temperature θ and constant entropy η obeys:

$$\begin{aligned} 2\varepsilon_{ij} &= u_{i,j} + u_{j,i} \\ \sigma_{ij,j} + \rho f_i &= \rho \ddot{u}_i, \quad i = 1, 2, 3 \\ \sigma_{ij} &= \rho \frac{\partial \psi}{\partial \varepsilon_{ij}}. \end{aligned}$$

Let us now calculate σ_{ij} by differentiating ψ with respect to the deformations. We obtain

$$\sigma_{ij} = \alpha_{ij} + c_{ijkl}\varepsilon_{kl} + \beta_{ijkl}\dot{\varepsilon}_{kl}.$$

If there are no stresses in a nondeformed state, $\alpha_{ij} = 0$, so

$$\sigma_{ij} = c_{ijkl}\varepsilon_{kl} + \beta_{ijkl}\dot{\varepsilon}_{kl}.$$

Then,

$$\sigma_{ij,j} = c_{ijkl} \frac{\partial \varepsilon_{kl}}{\partial x_j} + \beta_{ijkl} \frac{\partial \dot{\varepsilon}_{kl}}{\partial x_j}$$

$$\begin{aligned}
&= c_{ijkl} \frac{1}{2} \frac{\partial}{\partial x_j} (u_{k,l} + u_{l,k}) + \beta_{ijkl} \frac{1}{2} \frac{\partial}{\partial x_j} (\dot{u}_{k,l} + \dot{u}_{l,k}) \\
&= c_{ijkl} \frac{1}{2} (u_{k,lj} + u_{l,kj}) + \beta_{ijkl} \frac{1}{2} (\dot{u}_{k,lj} + \dot{u}_{l,kj}).
\end{aligned}$$

Thus, the equations of motion of a viscoelastic medium with a constant density, constant temperature and constant entropy are:

$$c_{ijkl} \frac{1}{2} (u_{k,lj} + u_{l,kj}) + \beta_{ijkl} \frac{1}{2} (\dot{u}_{k,lj} + \dot{u}_{l,kj}) + \rho f_i = \rho \ddot{u}_i \quad i = 1, 2, 3.$$

In the one-dimensional case these equations become the single equations for the displacement $u = u(x, t)$

$$c u_{xx} + \beta \dot{u}_{xx} + \rho f = \rho \ddot{u}.$$

This equation can also be written as

$$u_{tt} - \frac{\beta}{\rho} u_{xxt} - \frac{c}{\rho} u_{xx} - f = 0,$$

where $f = f(x, t)$ is given and ρ , β and c are known constants.

10. LINEAR THERMOELASTIC MEDIUM

In this section we will reach our ultimate goal – to derive the system of 20 PDEs, for 20 unknown functions, that governs the motion of a continuous medium.

Let us get started by rewriting the general law of conservation of energy (15) in a simpler form. For this, substitute in it $\epsilon = \psi + \eta \theta$ and use $\psi = \psi(\mathbf{E}, \theta, \mathbf{X})$. Then the law acquires the form

$$\rho \left(\frac{\partial \psi}{\partial E_{kl}} \frac{\partial E_{kl}}{\partial t} + \frac{\partial \psi}{\partial \theta} \frac{d\theta}{dt} + \frac{d\eta}{dt} \theta + \eta \frac{d\theta}{dt} \right) = \sigma_{ij} d_{ij} - q_{i,i} + \rho r$$

and with the help of

$$\frac{\partial E_{ij}}{\partial t} = d_{kl} \frac{\partial x_k}{\partial X_i} \frac{\partial x_l}{\partial X_j}$$

it becomes

$$\rho \left(\frac{\partial \psi}{\partial E_{kl}} \frac{\partial x_i}{\partial X_k} \frac{\partial x_j}{\partial X_l} d_{ij} + \frac{\partial \psi}{\partial \theta} \frac{d\theta}{dt} + \frac{d\eta}{dt} \theta + \eta \frac{d\theta}{dt} \right) = \sigma_{ij} d_{ij} - q_{i,i} + \rho r. \quad (20)$$

Now substitute σ_{ij} and η with

$$\sigma_{ij} = \rho \frac{\partial \psi}{\partial E_{kl}} \frac{\partial x_i}{\partial X_k} \frac{\partial x_j}{\partial X_l}, \quad \eta = - \frac{\partial \psi}{\partial \theta}$$

and 4 terms in the above law (20) cancel out. The law of conservation of energy becomes

$$\rho \theta \frac{d\eta}{dt} + q_{i,i} = \rho r. \quad (21)$$

A linear thermoelastic homogeneous medium is one for which the following assumptions hold:

1. The deformations are small, so the product of the gradients of the displacements are ignored. Also we substitute E_{ij} with ε_{ij} ;
2. The mass density ρ does not change during the deformation process;
3. The free energy ψ is a quadratic function of the components ε_{ij} and of the temperature change $T = \theta - T_0$. Also $|T|/T_0$ is small with respect to 1, thus $\theta \approx T_0$.
4. The components of the heat flow \mathbf{q} are linear functions of ε_{ij} , T and $T_{,i}$.

With these assumptions the gradient of the temperature becomes

$$G_i = \frac{\partial \theta}{\partial X_i} = \frac{\partial T}{\partial X_i} = \frac{\partial T}{\partial x_j} \frac{\partial x_j}{\partial X_i} = \frac{\partial T}{\partial x_j} \left(\delta_{ij} + \frac{\partial u_j}{\partial X_i} \right) = \frac{\partial T}{\partial x_i} + \frac{\partial T}{\partial x_j} \frac{\partial u_j}{\partial X_i},$$

where $u_j = x_j - X_j$. We ignore the product of the gradients, and obtain $G_i = \partial T / \partial x_i$. In the calculations that follow we will substitute Q_i with q_i , because $q_i = Q_j X_{j,i} = (\delta_{ij} - u_{j,i}) Q_j = Q_i - Q_j u_{j,i}$ and we ignore $Q_j u_{j,i}$ in the presence of Q_i .

To find the form of the functions ψ and q_i we develop them in Taylor series around their undeformed values, which are 0's. In the series for ψ we will keep terms up to and including second order, and in the series for q_i we will keep only the linear terms:

$$\rho \psi = a - \rho \eta_0 T - \frac{c_\varepsilon}{2T_0} T^2 + \alpha_{ij} \varepsilon_{ij} - \chi_{ij} \varepsilon_{ij} T + \frac{1}{2} c_{ijkl} \varepsilon_{ij} \varepsilon_{kl},$$

$$q_i = a_i + b_i T - k_{ij} T_{,j} + d_{ijk} \varepsilon_{jk}.$$

In these Taylor expansions the constants will be determined by the calculations that follow. Because of the requirement that ψ is symmetric with respect to the components ε_{ij} and ε_{ji} of the strain tensor, we have the following relations among the constants in its Taylor polynomial: $\alpha_{ij} = \alpha_{ji}$, $\chi_{ij} = \chi_{ji}$, $c_{ijkl} = c_{jikl} = c_{ijlk} = c_{klij}$.

A short calculation shows that in the theory of small deformations

$$\sigma_{ij} = \rho \partial \psi / \partial \varepsilon_{ij}. \quad (22)$$

We also remember from the previous section that $\eta = -\partial \psi / \partial \theta$. So

$$\eta = -\partial \psi / \partial \theta = -(\partial \psi / \partial T)(\partial T / \partial \theta) = -\partial \psi / \partial T. \quad (23)$$

Substitute the Taylor expansion for ψ in the last formulae for σ_{ij} and η to get

$$\sigma_{ij} = \rho \frac{\partial \psi}{\partial \varepsilon_{ij}} = \alpha_{ij} - \chi_{ij} T + c_{ijkl} \varepsilon_{kl},$$

$$\rho \eta = -\frac{\partial(\rho \psi)}{\partial T} = \rho \eta_0 + \frac{c_\varepsilon}{T_0} T + \chi_{ij} \varepsilon_{ij}.$$

We assume that when there is no deformation, i.e. $\varepsilon_{ij} = 0$, $T = 0$, there are no stresses, so $\sigma_{ij} = \alpha_{ij} = 0$. Also, $Q_j(\mathbf{E}, \theta, \mathbf{0}, \mathbf{X}) = 0$. Substituting 0 for Q_j in $q_i = Q_i - Q_j u_{j,i}$, we get $q_i(\varepsilon_{kl}, T, T_{,k})|_{T_1=T_2=T_3=0} = 0$. Thus, when there are no deformations, $T = 0$ and $q_i = 0$, and we obtain the following equation which relates the constants in the Taylor expansion for q_i , namely $0 = a_i + b_i T + d_{ijk} \varepsilon_{jk}$. But 1, T and ε_{jk} are linearly independent functions, so from this equation we conclude that the coefficients of these three linearly independent functions are zeros, i.e. $a_i = b_i = d_{ijk} = 0$. Substituting these constants in the Taylor expansion for q_i , we get $q_i = -k_{ij} T_{,j}$. With this expression for q_i the inequality $q_i \theta_{,i} \leq 0$ becomes $k_{ij} T_{,j} T_{,i} \geq 0$.

Thus, we arrive at the system of partial differential equations that every (linear) **continuous medium** obeys:

$$\begin{array}{ll} \sigma_{ij,j} + \rho f_i = \rho \ddot{u}_i & \text{equations of motion} \\ \rho T_0 \frac{\partial \eta}{\partial t} + q_{i,i} = \rho r & \text{law of conservation of energy} \\ \sigma_{ij} = c_{ijkl} \varepsilon_{kl} - \chi_{ij} T & \text{constitutive equations for the stress tensor} \\ \rho \eta = \rho \eta_0 + \frac{c_\varepsilon}{T_0} T + \chi_{ij} \varepsilon_{ij} & \text{constitutive equation for the entropy} \\ q_i = -k_{ij} T_{,j} & \text{constitutive equation for the heat flow} \\ \varepsilon_{ij} = \frac{1}{2} (u_{i,j} + u_{j,i}) & \text{equations of strain.} \end{array}$$

These are 20 equations for the 20 unknown functions σ_{ij} , u_i , q_i , ε_{ij} , T , η . The mass density ρ does not change during the deformation process, so ρ coincides with the initial mass density which we consider known. If we substitute the expressions for σ_{ij} , q_i , ε_{ij} , η from the last 4 lines of this system in the first two lines - the equations of motion and the law of conservation of energy, we obtain the equations

$$\begin{array}{ll} c_{ijkl} u_{k,jl} - \chi_{ij} T_{,j} + \rho f_i = \rho \ddot{u}_i, \quad i = 1, 2, 3 & \text{equations of motion} \\ k_{ij} T_{,ij} - c_\varepsilon \frac{\partial T}{\partial t} - \chi_{ij} T_0 \frac{\partial u_{i,j}}{\partial t} + \rho r = 0 & \text{equation of thermoconductivity} \end{array}$$

for the unknown functions u_i , T . These 4 equations are valid for any **thermoelastic anisotropic medium**, that is a medium with different mechanical and thermal properties in different directions. Some crystals are examples of such media. For isotropic media the constants in the constitutive equations remain unchanged under rotation of the body. Hence for such a medium $\chi_{ij} = \chi \delta_{ij}$, $k_{ij} = k \delta_{ij}$, $c_{ijkl} = \lambda \delta_{ij} \delta_{kl} + \mu \delta_{ik} \delta_{jl} + \nu \delta_{il} \delta_{jk}$. From the symmetries $c_{ijkl} = c_{jikl} = c_{ijlk} = c_{klij}$ it is clear that $\mu = \nu$, and consequently $c_{ijkl} = \lambda \delta_{ij} \delta_{kl} + \mu(\delta_{ik} \delta_{jl} + \delta_{il} \delta_{jk})$.

The constants λ and μ are called constants of Lamè. Thus, the equations of a **thermoelastic isotropic medium** are

$$\begin{aligned}(\lambda + \mu) u_{j,j i} + \mu u_{i,j j} - \chi T_{,i} + \rho f_i &= \rho \ddot{u}_i, \\ k T_{,i i} - c_\varepsilon \frac{\partial T}{\partial t} - \chi T_0 \frac{\partial u_{i,i}}{\partial t} + \rho r &= 0.\end{aligned}$$

The system of the general equations of linear elasticity in the case of absence of thermal effects was first derived by Navier [55] in 1821.

The system of 20 differential equations above or equivalently the system of 4 equations for thermoelastic anisotropic medium can be solved with suitable initial and boundary conditions. If the system of PDEs in question has a solution (u_1, u_2, u_3) , it is given by the formula (7) of Cesaro. This solution is unique, provided that $c_\varepsilon > 0$ and the quadratic form $c_{ijkl} \varepsilon_{ij} \varepsilon_{kl}$ is positive definite. The fact that the solution (7) of Cesaro satisfies the whole system is demonstrated by a direct substitution in the equations. The proof of uniqueness uses an identity, relating the variables involved in the system of PDEs. It is delightfully elegant and surprisingly short, see Ivanov [49] or Sokolnikoff [58].

11. TWO PROBLEMS

In this section we consider a couple of concrete problems.

Problem 1. Let us first consider an elastic body undergoing spherically symmetric deformation. Then the displacement vector is of the form

$$\mathbf{u} = u(r)\mathbf{e}_r, \quad r \neq 0$$

where \mathbf{e}_r is the unit vector along the radial direction. For such a displacement, compute (i) the corresponding stress components, (ii) the normal stress on a spherical surface $r = \text{constant}$ and (iii) the normal stress on a radial plane. Then determine $u(r)$ so that Navier's equation of equilibrium with zero body force is satisfied.

Solution. (i) The given form of the displacement vector can be rewritten as

$$\mathbf{u} = u(r)\frac{1}{r}\mathbf{x} = \phi(r)\mathbf{x},$$

where

$$\phi(r) = \frac{1}{r}u(r).$$

From this we find that $u_i = \phi(r)x_i$, so that

$$u_{i,j} = \phi(r)\delta_{ij} + \phi'(r)\left(\frac{1}{r}x_j\right)x_i = u_{j,i}.$$

Hence

$$u_{k,k} = 3\phi(r) + r\phi'(r).$$

Let us now substitute these last two results in the stress–displacement relation

$$\sigma = \lambda(\operatorname{div}\mathbf{u})\mathbf{I} + \mu(\nabla\mathbf{u} + \nabla\mathbf{u}^T)$$

and make use of the fact that $\phi(r) = (1/r)u(r)$. We obtain the following expression for the stresses associated with the given displacement field:

$$\sigma_{ij} = 2\left((\lambda + \mu)\delta_{ij} - 2\mu\frac{1}{r}x_ix_j\right)\frac{1}{r}u(r) + \left(\lambda\delta_{ij} + 2\mu\frac{1}{r^2}x_ix_j\right)u'(r).$$

(ii) For a spherical surface $r = \text{constant}$, we have $\mathbf{n} = \mathbf{e}_r$, so that $n_i = x_i/r$. Hence, by the formula

$$\sigma_N = \sigma_{ik}n_in_k,$$

enabling us to determine the normal stress σ_N directly from the stress components σ_{ik} , the normal stress σ_r on this surface is given by $\sigma_r = \sigma_{ij}n_in_j = (\sigma_{ij}x_ix_j)/r^2$. Using the expression for σ_{ij} obtained in part (i) of this problem, we get

$$\sigma_r = 2\lambda\frac{1}{r}u(r) + (\lambda + 2\mu)u'(r).$$

This normal stress is the radial stress.

(iii) If \mathbf{n} is the unit normal to a radial plane, we have $\mathbf{n} \cdot \mathbf{e}_r = 0$, and the normal stress σ_N on the plane is given by $\sigma_N = \sigma_{ij}n_in_j$. Another use of the expression for σ_{ij} obtained in part (i), we arrive at the following expression for the normal stress:

$$\sigma_\varphi = 2(\lambda + \mu)\frac{1}{r}u(r) + \lambda u'(r).$$

This normal stress is the peripheral stress.

(iv) Finally, to determine $u(r)$, we return to the expressions for $u_{i,j}$ and $u_{k,k}$ obtained in part (i) of the problem and calculate that

$$u_{i,ij} = u_{k,ki} = \left(\phi''(r) + \frac{4}{r}\phi'(r)\right)x_i.$$

Substituting these into Navier's equation of equilibrium

$$\mu\nabla^2u_i + (\lambda + \mu)u_{k,ki} + f_i = 0,$$

with $f_i = 0$, we see that it is satisfied if $\phi(r)$ obeys the following differential equation:

$$\frac{d^2\phi}{dr^2} + \frac{4}{r}\frac{d\phi}{dr} = 0.$$

The general solution of this equation is

$$\phi(r) = \frac{A}{r^3} + B,$$

where A and B are arbitrary constants. Thus,

$$u(r) = \frac{A}{r^2} + Br,$$

which is the sought solution of Navier's equation of equilibrium with zero body force.

The interested reader is invited to apply the ideas demonstrated in the above problem to solve the following

Problem 2. An elastic body undergoes a deformation, which is symmetric about the x_3 axes. Then the displacement vector is of the form

$$\mathbf{u} = u(R)\mathbf{e}_R, \quad R \neq 0,$$

where $R^2 = x_1^2 + x_2^2$ and \mathbf{e}_R is the unit vector along the radial direction in the cylindrical polar coordinate system with x_3 axis as axis. For this displacement compute (i) the corresponding stress components; (ii) the normal stress on a cylindrical surface $R = \text{constant}$; and (iii) the normal stress on a plane containing the x_3 axis. Also, determine $u(R)$ such that the Navier's equation of equilibrium with zero body force is satisfied.

12. THE CONTEMPORARY APPLICATIONS

In many applications the analytic solution (7) of Cesaro, to the system which a continuous medium obeys, can be obtained. Examples of such applications are the elongation, the twisting and the bending of cylindrical elastic beams; the stretching of a beam by its own weight; the twisting of a rectangular beam by two pairs of forces applied at each end of the beam; the twisting of circular cylinder with one base fixed and the other subjected to a pair of forces creating a torque; the displacement of a bended beam; and many others. Some 2-dimensional problems, like the displacement of an elastic membrane, subjected to uniform pressure from one side, have analytic solutions that use harmonic functions. The solution for the twisting of hollow, tube-like, beams also uses harmonic functions. The solution for the twisting of a cylinder by forces applied to its surface, and that for the bending of a tube with a circular or an elliptical cross-section, uses conformal maps. Most of these problems, solved in all detail, can be found in Sokolnikoff [58].

During the mid-1950s and 1960s the computer started to become a major tool for solving problems in continuum mechanics. At first the finite difference methods and the Rayleigh-Ritz method (using the theorem of minimal potential theory), were employed. Both of these methods required the solution of large numbers of simultaneous equations and faced the danger of the system becoming ill-conditioned as the number of equations increased. Finite difference methods have a long history, including contributions by Newton, Laplace, Gauss, Bessel and others. The method of finite differences replaces the defining differential equation with

equivalent difference equations. The boundary conditions are satisfied at discrete points by specifying either the function or its derivatives. The result of this analysis are numerical values of the function at discrete points throughout the body.

Computer simulations of exploding stars, the expansion of the early Universe, and the evolution of nebulae are so unbelievably realistic, only because they obey the equations of continuum mechanics. May be less dramatic, but significant from an applied point of view, is the fact that the flow of water or the spilling of oil can be modeled with the system for the motion of that continuous medium, and be presented visually in real time.

Modern cosmological simulations following the evolution of large portions of the Universe use numerical methods from hydrodynamics, more specifically the numerical solutions of the equations of compressible fluids. Simulations of merging clusters of galaxies are made this way. More specifically, the equations of motion for a compressible fluid are solved using a Lagrangian formulation in which the fluid is partitioned into elements, a subset of which is represented by particles of known mass and specific energy. Continuous fields are represented by interpolating between particles using a smoothing kernel, which is normally defined in terms of a sphere containing a fixed number of neighbors, centered on the particle in question. This method uses an artificial viscosity.

Continuum mechanics has become a fundamental science in investigations in tissue biomechanics. Soft tissue constitutive equations have been developed and the stresses and strains are being calculated for skin, tendon, ligament and bone. As new materials are being developed, they are being modeled as a continuum. Continuum mechanics is also being used in nanotechnology even on that small of a scale.

The most prominent relevant texts in Russian are listed as references [65] – [68].

Making an exhaustive list of the contemporary applications of Continuum Mechanics is impossible, as the subject is vast, vibrant, and multidisciplinary and develops literary every day. New branches of the subject are the nonlinear theory of elasticity, relativistic continuum mechanics and computational fluid dynamics. In recent years it has found connections with biomechanics and nanomechanics. A few of the most recent applications of continuum mechanics are: memory effects, the qualitative studies of the equations of Navier-Stokes, cross-diffusion systems from biology and physics, the decay of acceleration waves, and the fluid animation implementing numerical solutions to the 3D Navier-Stokes equations.

ACKNOWLEDGEMENTS. I am indebted to Professor Tsolo Ivanov, Professor Emeritus at the Department of Mechatronics, Robotics and Mechanics, Sofia University “St Kliment Ohridski”, Bulgaria, for valuable conversations in continuum mechanics.

13. REFERENCES

- [1] Bernoulli, D.: *Hydrodynamica sive de Viribus. Motibus Fluidorum Commentarii. Argentorati.* (1738).
- [2] Carnot, S.: *Reflexions sur la Puissance Motrice du Feu et sur les Machines Propres a Developper cette Puissance.* Paris (1824) = *Ann. École Norm* (2) **1** (1872), 393–457.
- [3] Cauchy, A.-L.: *Theorie de la propagation des ondes a la surface d'un fluide pesant d'une profondeur indefinie* (1815). *Mem. divers savants* (2) **1** (1816), 3-312 = *Oeuvres* (1) **1**, 5–318.
- [4] Cauchy, A.-L.: *Recherches sur l'équilibre et le mouvement intérieur des corps solides ou fluides, élastiques ou non élastiques.* *Bull. Soc. Philomath*, 9–13 (1823) = *Oeuvres* (2) **2**, 300–304.
- [5] Cauchy, A.-L.: *Memoire sur une espece particuliere de mouvement des fluides.* *J. Ecole Polytech.* **12** (1823), cahier 19, 204–214 = *Oeuvres* (2) **1**, 264–274.
- [6] Cauchy, A.-L.: *De la pression ou tension dans un corps solide.* *Ex. de math.* **2** (1827), 42–56 = *Oeuvres* (2) **7**, 60–78.
- [7] Cauchy, A.-L.: *Sur la condensation et la dilatation des corps solides.* *Ex. de math.* **2** (1827), 60–69 = *Oeuvres* (2) **7**, 82–83.
- [8] Cauchy, A.-L.: *Sur les moments d'inertie.* *Ex. de math.* **2** (1827), 93–103 = *Oeuvres* (2) **7**, 124–136.
- [9] Cauchy, A.-L.: *Sur les relations qui existent dans l'état d'équilibre d'un corps solide ou fluide, entre les pressions ou tensions et les forces acceleratrices.* *Ex. de math.* **2** (1827), 108–111 = *Oeuvres* (2) **7**, 141–145.
- [10] Cauchy, A.-L.: *Sur les centers, les plans principaux et les axes principaux des surfaces du second degre.* *Ex. de math.* **3** (1828), 1–22 = *Oeuvres* (2) **8**, 9–35.
- [11] Cauchy, A.-L.: *Sur les equations qui expriment les conditions d'équilibre, ou les lois du mouvement interieur d'un corps solide, elastique, ou non elastique.* *Ex. de math.* **3** (1828), 160–187 = *Oeuvres* (2) **8**, 195–226.
- [12] Cauchy, A.-L.: *Sur quelques theoremes relatifs a la condensation ou a la dilatation des corps.* *Ex. de math.* **3** (1828), 237–244 = *Oeuvres* (2) **8** (1828), 278–287.
- [13] Cauchy, A.-L.: *Sur les pressions ou tensions supportees en un point donne d'un corps solide par trois plans perpendiculaires entre eux.* *Ex. de math.* **4** (1829), 30–40 = *Oeuvres* (2) **9**, 41–52.
- [14] Cauchy, A.-L.: *Sur la relation qui existe entre les pressions ou tensions supportees par deux plans quelconques en un point donne d'un corps solide.* *Ex. de math.* **4** (1829), 41–46 = *Oeuvres* (2) **9**, 53–55.
- [15] Cauchy, A.-L.: *Sur les corps solides ou fluides dans lesquels la condensation ou dilatation lineaire est la meme en tous sens autour de chaque point.* *Ex. de math.* **4** (1829), 214–216 = *Oeuvres* (2) **9**, 254–258.
- [16] Cauchy, A.-L.: *Sur l'équilibre et le mouvement interieur des corps consideres comme des masses continues.* *Ex. de math.* **4** (1829), 293–319 = *Oeuvres* (2) **9**, 243–369.
- [17] Cauchy, A.-L.: *Sur les diverses methodes a l'aide desquelles on peut etablir les equations qui representent les lois d'équilibre, ou le mouvement interieur des corps solides ou fluides.* *Bull. sci. math. soc. prop. conn.* **13** (1830), 169–176.

- [18] Cauchy, A.-L.: Mèmoire sur les dilatations, les condensations et les rotations produits par un changement de forme dans un système de points matériels. *Ex. d'an. phys. math.* **2**, 302–330 (1841) = *Oeuvres* (2) **12**, 343–377.
- [19] Cesaro, E.: Rendiconto dell' Accademia delle scienze fisiche e matematiche (Società reale di Napoli) (1906).
- [20] Clausius, R.: Ueber eine veränderte Form des zweiten Hauptsatzes der mechanischen Wärmetheorie. *Ann. Physik* **93** (1854), 481–506.
- [21] Clausius, R.: Uber verschiedene für die Anwendung bequeme Formen der Hauptgleichungen der mechanischen Wärmetheorie. *Vjschr. nat. Ges. Zurich* **10** (1865), 1–59.
- [22] D'Alembert, J.L.: *Traite de l'Equilibre et du Mouvement des Fluides pour servir de Suite au Traite de Dynamique*. Paris (1749); 2nd ed., 1770.
- [23] Duchem, P.: *Hydrodynamique, Électisité, Acoustique*. Paris. Hermann, Vol. I, II (1891).
- [24] Duchem, P.: Commentaire aux principes de la thermodynamique, Première partie. *J. math pures appl.* (4) **8** (1892), 269–330.
- [25] Euler, L.: De motu corporum in superficiebus mobilibus. *Ousc. var. arg.* **1** (1746), 1–136 = *Opera omnia* (2) **6**, 75–174.
- [26] Euler, L.: Decouverte d'un nouveau principe de mécanique. *Mem. Acad. Sci. Berlin* **6** (1750), 185–217 = *Opera omnia* (2) **5**, 81–108.
- [27] Euler, L.: Recherches sur l'effet d'une machine hydraulique proposée par Mr. Segner Professeur a Gottingue. *Mem. Acad. Sci. Berlin* **6** (1750), 311–354 = *Opera omnia* (2) **15**, 1–39.
- [28] Euler, L.: Recherche sur une nouvelle manière d'élever de l'eau proposée par Mr. De Mour. *Mem. Acad. Sci. Berlin* **7** (1751), 305–330 = *Opera omnia* (2) **15**, 134–156.
- [29] Euler, L.: Théorie plus complète des machines qui sont mises en mouvement par la réaction de l'eau. *Mem. Acad. Sci. Berlin* **10** (1754), 227–295 = *Opera omnia* (2) **15**, 157–218.
- [30] Euler, L.: 1757 Principes généraux de l'état d'équilibre des fluides. *Mèm. Acad. Sci. Berlin* **11**, 217–273 (1755) = *Opera omnia* (2) **12**, 2–53.
- [31] Euler, L.: 1757 Principes généraux du mouvement des fluides. *Mèm. Acad. Sci. Berlin* **11**, 274–315 (1755) = *Opera omnia* (2) **12**, 54–91.
- [32] Euler, L.: Continuation des recherches sur la théorie du mouvement des fluides. *Mèm. Acad. Sci. Berlin* **11** (1755), 316–361 = *Opera omnia* (2) **12**, 92–132.
- [33] Euler, L.: 1761 Principia motus fluidorum (1752-1755). *Novi Comm. Acad. Sci. Petrop.* **6** (1756-1757), 271–311 = *Opera omnia* (2) **12**, 133–168.
- [34] Euler, L.: Lettre de M. Euler a M. La Grange, Recherches sur la propagation des branlemens dans un milieu élastique. *Misc. Taur.* **2**² (1760-1761), 1–10 = *Opera omnia* (2) **10**, 255–263.
- [35] Euler, L.: De motu fluidorum a diverso caloris gradu oriundo. *Novi Comm. Acad. Sci. Petrop.* **11** (1765), 232–267 = *Opera omnia* (2) **12**, 244–271.
- [36] Euler, L.: Recherches sur la connaissance mécanique des corps. *Mem. Acad. Sci. Berlin* **14** (1758), 131–153.

- [37] Euler, L.: Du mouvement de rotation des corps solides autour d'un axe variable. *Mem. Acad. Sci Berlin* **14** (1758), 154–193.
- [38] Euler, L.: Supplement aux recherches sur la propagation du son. *Mem. Acad. Sci Berlin* **15** (1759), 210–240 = *Opera omnia* (3) **1**, 452–483.
- [39] Euler, L.: Recherche sur le mouvement des rivieres (1751). *Mem. Acad. Sci Berlin* **16**] (1760), 101–118 = *Opera omnia* (2) **12**, 212–288.
- [40] Euler, L.: Theoria Motus Corporum Solidorum seu Rigidorum ex Primis nostrae Cognitionis Principiis Stabilita et ad Omnis Motus, qui in hujusmodi Corpora Cadere Possunt, *Accomodata* (1765), Rostock = *Opera omnia* (2) **3, 4**, 3–293.
- [41] Euler, L.: Sectio secunda de principiis motus fluidorum. *Novi Comm. Acad. Sci. Petrop.* **14** (1769), 270–386 = *Opera omnia* (2) **13**, 73–153.
- [42] Euler, L.: De gemina methodo tam aequilibrium quam motus corporum flexibilium determinandi et utriusque egregio consensu. *Novi Comm. Acad. Sci. Petrop.* **20**, 286–303 (1775).
- [43] Euler, L.: Formulae generales pro translatione quacunque corporum rigidorum. *Novi Comm. Acad. Sci. Petrop.* **20** (1775), 189–207.
- [44] Fresnel, A.: Second supplément au mémoire sur la double réfraction. *Oeuvres* **2** (1822), 369–442.
- [45] Georgieva, B.: Review of continuum mechanics and its history. Part II. The mechanics of thermoelastic media. Perfect fluids. Linearly viscous fluids. *Ann. Sofia Univ., Fac Math and Inf.*, **107** (2020), 55–77.
- [46] Green, G.: On the laws of reflection and refraction of light at the common surface of two non-crystalized media. *Trans. Cambridge Phil. Soc.* **7** (1839), 1–24.
- [47] Green, G.: On the propagation of light in crystalized media. *Trans. Cambridge Phil. Soc.* **7** (1841), 121–140.
- [48] Helmholtz, H.: Uber Integrale der hydrodynamischen Gleichungen, welche den Wirbelbewegungen entsprechen. *J. Reine Angew. Math.* **55** (1858), 25–55.
- [49] Ivanov, Ts.: *Theory of Elasticity*. Sofia University Publishing House (1995).
- [50] Joule, J. P.: On the caloric effects of magneto-electricity, and on the mechanical value of heat. *Phil. Mag.* (3) **23** (1843) , 263–276, 347–355, 435–443.
- [51] Joule, J. P.: On the existence of an equivalent relation between heat and the ordinary forms of mechanical power. *Phil. Mag.* (3) **27** (1845), 205–207.
- [52] Kirchhoff, G.: Uber das Gleichgewicht und die Bewegung eines unendlich dunnen elastischen Stabes. *J. Reine Angew. Math.* **56** (1859), 285–313.
- [53] Lagrange, J. L.: Nouvelles recherches sur la nature et la propagation du son. *Misc. Taur.* **2**² (1760-1761), 11–172 = *Oeuvres* **1** (1762), 151–316.
- [54] Lagrange, J. L.: *Mechanique Analitique*, Paris, *Oeuvres* **11, 12** (1788).
- [55] Navier, C. L. M. H.: Sur les lois des mouvements des fluides, en ayant egard a l'adhesion des moleules. *Ann. Chimie* **19** (1821), 244–260.
- [56] Rankine, W. J. M.: Laws of the elasticity of solid bodies. *Cambr. Dubl. Math. J.* **6** (1851), 41–80, 178–181.
- [57] Rankine, W. J. M.: On axes of elasticity and crystalline forms. *Phil. Trans. Roy. Soc. Lond.* **46** (1856), 261–285.

- [58] Sokolnikoff, I. S.: *Mathematical Theory of Elasticity*. McGraw-Hill Book Company Inc., 2nd Ed. (1956).
- [59] St Venant, A.-J.-C. B.: Sur les pressions qui se développent a l'intérieur des corps solides lorsque les déplacements de leurs points, sans altérer l'élasticité, ne peuvent cependant pas être considérés comme très petits. *Bull. Soc. Philomath.* **5** (1844) , 26–28.
- [60] St Venant, A.-J.-C. B.: Theorie de l'elasticite des solides, ou cinematique de leurs deformations. *L'Institut* **32**¹ (1864), 389–390.
- [61] Stokes, G. G.: On the theories of the internal friction of fluids in motion, and of the equilibrium and motion of elastic solids. *Trans. Cambridge Phil. Soc.* **8** (1844-1849), 287–319.
- [62] Voltera, V.: L'Equilibre des corps élastiques. *Annales de l'ecole normal supérieure* **24** (1907), 416–417.
- [63] Truesdell, C., Toupin, R.: *The Classical Field Theories*. Encyclopedia Of Physics **245** (S. Flugge, Ed.), Vol. III. Principles of Classical Mechanics and Field Theory. Springer-Verlag, Berlin (1960).
- [64] Waterston, J. J.: Note on the physical constitution of gaseous fluids and a theory of heat. Appendix to *Thoughts on the Mental Functions*. Edinburgh (1843).
- [65] Ландау, Л. Д., Лифшиц, Е. М.: *Теория упругости*, Наука, Москва, 1975.
- [66] Лурье, А. И.: *Теория упругости*, Наука, Москва, 1970.
- [67] Работнов, Ю. Н., *Механика твердого деформируемого тела*, Наука, Москва, 1988.
- [68] Седов, Л. И.: *Механика сплошной среды*, том 1, Наука, Москва, 1970.

Received on November 20, 2020

BOGDANA A. GEORGIEVA
 Faculty of Mathematics and Informatics
 Dept. of Mechatronics, Robotics and Mechanics
 Sofia University St Kliment Ohridski
 5 James Bourchier Blvd.
 1164 Sofia
 BULGARIA
 E-mail: georgieva@fmi.uni-sofia.bg

ГОДИШНИК НА СОФИЙСКИЯ УНИВЕРСИТЕТ СВ. КЛИМЕНТ ОХРИДСКИ “

ФАКУЛТЕТ ПО МАТЕМАТИКА И ИНФОРМАТИКА

Том 107

ANNUAL OF SOFIA UNIVERSITY ST. KLIMENT OHRIDSKI“

FACULTY OF MATHEMATICS AND INFORMATICS

Volume 107

REVIEW OF CONTINUUM MECHANICS AND ITS HISTORY
PART II. THE MECHANICS OF THERMOELASTIC MEDIA.
PERFECT FLUIDS. LINEARLY VISCOUS FLUIDS

BOGDANA A. GEORGIEVA

This paper is the second of a series of two articles reviewing the contributions of continuum mechanics and its history. The review is written for the mathematician who is not a specialist in this field, and aims to give an in-depth overview of the mathematics as well as a historical perspective of this field. The first of the two papers [10], Part I. “Deformation and Stress. Conservation Laws. Constitutive Equations”, starts at the very origins of continuum mechanics and brings the reader up to the 1820's when Navier publishes the system of the general equations of linear elasticity in 1821. The present paper continues, discussing the consequences of this system, some of its simplifications and approaches for solution. It also gives a perspective of how waves propagate in continuous media. Reviewed are also perfect fluids and linearly viscous fluids. At the end, the paper discusses the conditions for compatibility of the stresses.

Keywords: Mechanics of continuous media, continuum mechanics, hydrodynamics, history of continuum mechanics, elasticity, theory of elasticity.

2020 Math. Subject Classification: XX76-02.

1. INTRODUCTION

The first attempt to discuss the motion of a continuous medium in more than one dimension occurs in an isolated passage by D. Bernoulli from 1738 [2], §11, paragraph 4.

We are surrounded by matter in the form of continuous media – deformable solids, liquids and gasses. To study how they move in response to forces, while

obeying the natural laws, we need two sets of coordinates. **Material coordinates**, also called **Lagrangian coordinates**, are denoted by (X_1, X_2, X_3) and are the coordinates of the material points of the continuous medium at time $t = 0$. Lagrange introduced them in 1788 in [22], part II, section II. **Spatial coordinates**, also known as **Eulerian coordinates**, are denoted by (x_1, x_2, x_3) and are the coordinates of the points of 3-dimensional space (in which we observe the medium) occupied by the medium at time $t > 0$. Since the material coordinates are the coordinates of the material points at an arbitrary initial time $t = 0$, they can serve for all time as names for the *particles* of the material. The spatial coordinates, on the other hand, we think of as assigned once and for all to a point in the Euclidean space. They are the names of *places*. The motion $\mathbf{x} = \mathbf{x}(\mathbf{X}, t)$ chronicles the places \mathbf{x} occupied by the particle \mathbf{X} in the course of time. Under external influences - forces and heating - the continuous body deforms. *The goal of Continuum Mechanics is to find the family of transformations*

$$x_i = x_i(X_1, X_2, X_3, t), \quad i = 1, 2, 3, \quad (1)$$

giving the Eulerian coordinates as functions of the Lagrangian coordinates for $t \geq 0$. This motion is deterministic, obeying only the natural laws.

The general theory of the motion of a continuous medium, understood as a family of deformations continuously varying in time, is almost exclusively due to Euler, published in the period 1745 – 1766, references [25] – [40] in the first paper of this review, and Cauchy, published in the period 1815 – 1841, references [3] – [18] in the first paper of this review. Important special results were added by D’Alembert in 1749, Green in 1839, Stokes in 1845, Helmholtz in 1858 and Cesaro in 1906, also cited in the first part of this review.

2. LINEAR THERMOELASTIC CONTINUOUS MEDIA

The systems of equations

$$c_{ijkl} u_{k,jl} - \chi_{ij} T_{,j} + \rho f_i = \rho \ddot{u}_i, \quad i = 1, 2, 3 \quad \text{equations of motion} \quad (2)$$

$$k_{ij} T_{,ij} - c_\varepsilon \frac{\partial T}{\partial t} - \chi_{ij} T_0 \frac{\partial u_{i,j}}{\partial t} + \rho r = 0 \quad \text{equation of thermal conductivity} \quad (3)$$

for the unknown functions u_i , T , which are valid for any linear **thermoelastic** anisotropic medium, were first published in 1821 by Navier [30]. Here and in the rest of the paper a dot above a variable denotes a differentiation with respect to time and two dots denote a double differentiation with respect to time. Here u_i are the three components of the vector of displacement \mathbf{u} , T is the temperature difference, and are functions of the space coordinates (x_1, x_2, x_3) and the time t . As usually in the literature, a first partial derivative with respect to a space coordinate is denoted by one lower index after a comma, a second partial derivative with respect to space

coordinates is denoted by two lower indexes after a comma; ρ is the mass density, f_i are the components of the assigned (mass) force, r is the heat source, and c_{ijkl} , χ_{ij} , k_{ij} , c_ε , T_0 are constants.

The system (2), (3) in the case of an isotropic body acquires the form:

$$(\lambda + \mu) u_{j,ji} + \mu u_{i,jj} - \chi T_{,i} + \rho f_i = \rho \ddot{u}_i \quad (4)$$

$$k T_{,ii} - c_\varepsilon \frac{\partial T}{\partial t} - \chi T_0 \frac{\partial u_{i,i}}{\partial t} + \rho r = 0. \quad (5)$$

Both these systems simplify significantly if the process is isothermal or adiabatic. A process is called **isothermal** if the changes that are taking place are “slow”, so that the change in temperature is small and can be ignored. In the notation we use, $T = 0$. In that case we do not consider at all the equation of thermal conductivity. So the remaining equations are

$$c_{ijkl} u_{k,jl} + \rho f_i = \rho \ddot{u}_i, \quad i = 1, 2, 3 \quad (6)$$

for an anisotropic body and

$$(\lambda + \mu) u_{j,ji} + \mu u_{i,jj} + \rho f_i = \rho \ddot{u}_i \quad (7)$$

for an isotropic medium. The equations (6) and (7) are known as **the isothermal equations of elasticity** for an anisotropic and isotropic medium, respectively. The constants c_{ijkl} , λ , μ are called isothermal constants.

The process is called adiabatic if the changes that take place in the medium are “fast”, so that the heat exchange that takes place between different parts of the body, being a “slower” process, can be ignored, that is, $q_i = 0$. Of course, in this case there are no sources of heat.

It is interesting that an adiabatic process is also isoentropic, that is, has a constant entropy $\eta = \eta_0 = \text{constant}$. This can be seen from the equations

$$\rho T_0 \frac{\partial \eta}{\partial t} + q_{i,i} = \rho r \quad \text{law of conservation of energy}$$

$$\rho \eta = \rho \eta_0 + \frac{c_\varepsilon}{T_0} T + \chi_{ij} \varepsilon_{ij} \quad \text{constitutive equation for the entropy}$$

which were derived in Part I of this review, as a part of the system of 20 equations which any linear thermoelastic continuous medium obeys. From the constitutive equation for the entropy, for an anisotropic body, we get

$$T = -\frac{\chi_{ij} T_0}{c_\varepsilon} \varepsilon_{ij}.$$

For an isotropic body the relationship between the entropy and the deformations is

$$T = -\frac{\chi T_0}{c_\varepsilon} \varepsilon_{ii}.$$

After substituting these last two equations for T in the equations of motion (2) for an elastic anisotropic medium and in equations (4) for a thermoelastic isotropic medium, these equations acquire the form:

$$c_{ijkl}^a u_{k,jl} + \rho f_i = \rho \ddot{u}_i, \quad i = 1, 2, 3 \quad (8)$$

and respectively

$$(\lambda^a + \mu^a) u_{j,ji} + \mu^a u_{i,jj} + \rho f_i = \rho \ddot{u}_i, \quad (9)$$

where the adiabatic constants c_{ijkl}^a , λ^a and μ^a are related to their corresponding isothermal constants via the equations

$$c_{ijkl}^a = c_{ijkl} + \frac{\chi_{ij} \chi_{kl} T_0}{c_\varepsilon}$$

$$\lambda^a = \lambda + \frac{\chi^2 T_0}{c_\varepsilon}, \quad \mu^a = \mu.$$

In the remaining of the paper we will drop the upper index of the constants in the adiabatic equations of elasticity, namely in equations (8) and equations (9), so they will not differ in form from their corresponding isothermal equations (6) and (7). We will call these equations **the equations of elasticity**.

For an isothermal process, the constitutive equations

$$\sigma_{ij} = c_{ijkl} \varepsilon_{kl} - \chi_{ij} T \quad (10)$$

for the components of the stress tensor, acquire the form

$$\sigma_{ij} = c_{ijkl} \varepsilon_{kl}. \quad (11)$$

In the case of adiabatic process the relationship among stresses and deformations is analogous, if the constants c_{ijkl} are the adiabatic constants.

For isotropic bodies from

$$c_{ijkl} = \lambda \delta_{ij} \delta_{kl} + \mu (\delta_{ik} \delta_{jl} + \delta_{il} \delta_{jk})$$

follows that

$$\sigma_{ij} = \lambda \varepsilon_{kk} \delta_{ij} + 2\mu \varepsilon_{ij}. \quad (12)$$

The equations (11) or respectively (12), giving the relationship between the stresses and the deformations, are known as **the generalized law of Hooke**. Equations (12) with $\lambda = \mu$ were derived from a molecular model by Navier published in 1821 [28], [29]; more generally by Poisson [32] in 1829.

The **elasticities** λ , μ and c_{ijkl} in equations (11) and (12) are material constants or functions of the temperature or entropy. Their physical dimensions are those of stress, and they bear no physical connection with the mathematically analogous viscosities appearing in the Navier-Poisson law, discussed in section 8 “Linearly Viscous Fluids” of this paper.

The equations (11) or respectively (12) together with the equations of motion

$$\sigma_{ij,j} + \rho f_i = \rho \ddot{u}_i$$

of a linear continuous medium and the equations of strain

$$\varepsilon_{ij} = \frac{1}{2}(u_{i,j} + u_{j,i})$$

form **the system of equations of elasticity** for an anisotropic and isotropic body respectively. This system consists of 15 equations for the 15 unknown functions u_i , ε_{ij} and σ_{ij} , where ε_{ij} are the components of the strain tensor and σ_{ij} are the components of the stress tensor. Unlike that, the systems (6) and (7) are systems with 3 equations each for the three unknown displacements u_i . These equations are typically called **the equations of elasticity in displacements or equations of Lamè**, and can be solved with appropriate initial and boundary conditions.

Because of the symmetries $c_{ijkl} = c_{jikl} = c_{ijlk} = c_{klij}$, the number of the independent components of the tensor c_{ijkl} is significantly smaller than that of a general tensor of rank 4. Thus, it is appropriate to replace couples of indexes with a single index via the following scheme: 11 \rightarrow 1, 22 \rightarrow 2, 33 \rightarrow 3, 23 and 32 \rightarrow 4, 31 and 13 \rightarrow 5, 12 and 21 \rightarrow 6. The following notation is also used to denote the components of the stress tensor and those of the tensor of deformations:

$$\sigma_1 = \sigma_{11}, \sigma_2 = \sigma_{22}, \sigma_3 = \sigma_{33}, \sigma_4 = \sigma_{23}, \sigma_5 = \sigma_{31}, \sigma_6 = \sigma_{12}$$

$$\varepsilon_1 = \varepsilon_{11}, \varepsilon_2 = \varepsilon_{22}, \varepsilon_3 = \varepsilon_{33}, \varepsilon_4 = 2\varepsilon_{23}, \varepsilon_5 = 2\varepsilon_{31}, \varepsilon_6 = 2\varepsilon_{12}.$$

Then the generalized law of Hooke (11) acquires the form

$$\sigma_\alpha = c_{\alpha\beta} \varepsilon_\beta, \tag{13}$$

where the Greek indices run from 1 to 6, and repeated indices denote summation from 1 to 6. Because $c_{\alpha\beta} = c_{\beta\alpha}$, the number of independent constants in the generalized law of Hooke (13) for an arbitrary anisotropic body is 21.

The function

$$U = \frac{1}{2}\sigma_{ij}\varepsilon_{ij} = \frac{1}{2}\sigma_\alpha\varepsilon_\alpha = \frac{1}{2}c_{ijkl}\varepsilon_{ij}\varepsilon_{kl} = \frac{1}{2}c_{\alpha\beta}\varepsilon_\alpha\varepsilon_\beta \tag{14}$$

is called the density of the potential energy of the deformation, or the **elastic potential**. So the potential energy of the deformation is

$$U = \frac{1}{2}c_{ijkl} \int_V \varepsilon_{ij}\varepsilon_{kl}dV = \int_V U dV$$

and

$$\sigma_\alpha = \frac{\partial U}{\partial \varepsilon_\alpha}.$$

There is a physical reason to require that the elastic potential be a positive definite form, because then, in any given small strain from an unstressed state, the stress must do positive work. Assuming that the elastic potential is positive definite, it follows that the constants $c_{\alpha\beta}$ satisfy the following restrictions: $c_{11} > 0, \dots, \det|c_{\alpha\beta}| > 0$. In the case of isotropic body these inequalities acquire the form $\lambda + 2\mu > 0, 4\mu(\lambda + \mu) > 0, \dots, 4\mu^5(3\lambda + 2\mu) > 0$. Hence the necessary and sufficient condition for these inequalities to be satisfied is:

$$3\lambda + 2\mu > 0, \quad \mu > 0. \quad (15)$$

The elastic potential and its resulting potential energy of the deformation are due to Green, who published them in 1839 [11], and in 1841 [12]. He proposed that the work done by stress in a deformation depends only upon the strain and is recoverable work. In his original papers, Green defines the **stored energy** Σ by

$$\Sigma(\varepsilon) = \frac{1}{2} \sigma_{km} \varepsilon_{km},$$

(later renamed the elastic potential U , which we defined with (14)). Thus, in Green's theory the number of independent elasticities is 21. He derives that

$$\sigma_{km} = \frac{\partial \Sigma}{\partial \varepsilon_{km}}. \quad (16)$$

By the representation theorem for isotropic scalar functions, it follows that the stored energy can be expressed in terms of the first and second invariants of the tensor ε as

$$\Sigma = \frac{1}{2}(\lambda + 2\mu)I_\varepsilon^2 - 2\mu II_\varepsilon.$$

A body is called **hyperelastic** if it obeys Green's theory, based upon the use of Σ as a stress potential according to (16). This theory has some remarkable results, which we review next.

The fact that $\det|c_{\alpha\beta}| > 0$ guarantees that the equations of the generalized law of Hooke (13) can be solved for the deformations, obtaining

$$\varepsilon_\alpha = s_{\alpha\beta} \sigma_\beta,$$

where the matrix $|s_{\alpha\beta}|$ is the inverse of the matrix of elastic constants $|c_{\alpha\beta}|$, and is called **the matrix of stiffnesses**. In the isotropic case the deformations ε_{ij} can be expressed with the stresses, if we take in consideration that for $i = j$ from the generalized Hooke's law (12), namely, $\sigma_{ij} = \lambda \varepsilon_{kk} \delta_{ij} + 2\mu \varepsilon_{ij}$, we obtain

$$\sigma_{ii} = (3\lambda + 2\mu)\varepsilon_{ii}.$$

Then

$$\begin{aligned} \varepsilon_{ii} &= \frac{1}{2\mu}(\sigma_{ij} - \lambda \varepsilon_{kk} \delta_{ij}) \\ &= \frac{1}{2\mu} \sigma_{ij} - \frac{\lambda}{2\mu(3\lambda + 2\mu)} \sigma_{kk} \delta_{ij} = \frac{1 + \nu}{E} \sigma_{ij} - \frac{\nu}{E} \sigma_{kk} \delta_{ij}, \end{aligned} \quad (17)$$

where the constant

$$E = \frac{\mu(3\lambda + 2\mu)}{(\lambda + \mu)}$$

is always positive and is called the module of Jung. For metals the module of Jung is of the order of 10^{11} N/m². The constant ν is called the coefficient of Poisson and is $\nu = \lambda/(2\lambda + 2\mu)$. It is clear that $-1 < \nu < 1/2$. For all known materials Poisson's coefficient is positive. For metals it varies usually in the interval $[1/4, 1/3]$.

Let us now consider a couple of special cases. Let us assume that $f_i = 0$ and that the problem is static, i.e., the components u_i of the displacement do not depend on the time t . In this case the initial conditions of the system of differential equations are no longer present and only the boundary conditions $u_i(\mathbf{x}, t) = g_i(\mathbf{x}, t)$ for $x \in S_u$ and $\sigma_{ij}(\mathbf{x}, t)n_j(\mathbf{x}, t) = h_i(\mathbf{x}, t)$ for $x \in S_\sigma$ remain, because the medium is elastic and not thermoelastic.

1. Simple Shear

Simple shear is characterized by the following stresses:

$$\sigma_{23} = \text{constant} \neq 0, \quad \text{the rest of } \sigma_{ij} = 0. \quad (18)$$

These stresses satisfy the equations of equilibrium $\sigma_{ij,j} = 0$. The deformations that correspond to them are:

$$\varepsilon_{23} = \frac{1}{2\mu}\sigma_{23}, \quad \text{the rest of } \varepsilon_{ij} = 0.$$

The geometric interpretation of the tensor of deformations, which was explained in the first part of this review, follows that a cube with sides parallel to the coordinate planes, will deform under a simple shear in such a way that the right angle between the edges of the cube, that are parallel to the axes x_2 and x_3 decreases (if $\sigma_{23} > 0$) or increases (if $\sigma_{23} < 0$) with the angle $\gamma_{23} = 2\varepsilon_{23}$. From equations (10) follows that

$$\mu = \frac{\sigma_{23}}{\gamma_{23}}.$$

Thus, μ has the meaning of the ratio between the so called shearing stress σ_{23} to the resulting from it change γ_{23} of the right angle. The constant μ is called **module of shearing**, it is often denoted in the technical literature by G .

2. Hydrostatic Pressure

We consider an elastic body with an arbitrary shape. Its boundary is subjected to stresses, that are applied perpendicularly to the surface, toward the body, and have a constant intensity $p > 0$. Then

$$\sigma_i = -pn_i, \quad (19)$$

where n_i are the components of the outward unit normal to the surface of the body. The stresses

$$\sigma_{11} = \sigma_{22} = \sigma_{33} = -p, \quad \text{and } \sigma_{ij} = 0 \text{ if } i \neq j$$

satisfy the equations of equilibrium with $f_i = 0$ and boundary conditions given by (19). From equations (17) it follows that

$$\varepsilon_{11} = \varepsilon_{22} = \varepsilon_{33} = -\frac{p}{3\lambda + 2\mu}, \quad \text{and } \varepsilon_{ij} = 0 \quad \text{if } i \neq j. \quad (20)$$

Both in this and in the previously considered example, the deformations are constant, and thus satisfy the conditions for compatibility of St. Venant, discussed in detail in the first paper of this review. Hence from them the displacements u_i can be calculated, that correspond to the stresses in consideration. From equations (12) one calculates the relative change in the volume (expansion if $p < 0$) and (contraction if $p > 0$). Let $\varepsilon \equiv \varepsilon_{ii}$. Then

$$\varepsilon = -\frac{p}{k}, \quad (21)$$

where

$$k = \lambda + 2\mu/3 = E/(3 - 6\nu) \quad (22)$$

is called the **module of contraction**. It is the ratio of the hydrostatic pressure to the relative change of volume. From the inequalities (21) it follows that $k > 0$.

The elastic material is called noncompressible if under pressure the relative change ε of the volume remains zero. In that case from (21) and (22) we calculate that $\nu = 1/2$.

3. THE LAW OF CONSERVATION OF MECHANICAL ENERGY

We considered the law of conservation of mechanical energy in part I of this review and showed that it has the form

$$\frac{dK}{dt} + \int_V \sigma_{ij} d_{ij} dV = W,$$

where $K = \int_V \rho v_i v_i / 2 dV$ is the kinetic energy,

$$W = \int_V \rho f_i v_i dV + \int_S \sigma_i v_i dS \quad (23)$$

is the power of the external forces and $d_{ij} \equiv (v_{i,j} + v_{j,i})/2 = d_{ji}$ is the tensor of the rate of deformations, introduced by Euler in 1769 [9], §§ 9-12. In the case of small deformations $d_{ij} = \partial\varepsilon_{ij}/\partial t$. The total time-derivative of the potential energy U of the deformations is

$$\frac{dU}{dt} = \frac{1}{2} c_{ijkl} \int_V \frac{\partial}{\partial t} (\varepsilon_{ij} \varepsilon_{kl}) dV = c_{ijkl} \int_V \varepsilon_{ij} \frac{\partial}{\partial t} \varepsilon_{kl} dV.$$

Then it follows that in the linear theory of elasticity the law of conservation of mechanical energy acquires the form

$$\frac{d}{dt} (K + U) = W. \quad (24)$$

If we denote by

$$A(\tau) = \int_0^\tau W dt \quad (25)$$

the work done by the external forces during the interval of time $[0, \tau]$ and assume that at $t = 0$ the body was in an undeformed state and at rest, i.e. $K(0) = U(0) = 0$, then from equation (24) it follows that

$$K(\tau) + U(\tau) = A(\tau), \quad (26)$$

where the argument of the functions K , U and A determines the moment at which they are evaluated. Equation (26) shows that the sum of the kinetic and the potential energies at a given moment equals the work done by the mass forces and the surface forces upto that moment.

4. THE STATIC PROBLEM

In a static problem we are not interested in the process of deformation, but only in the final state, which we regard as an equilibrium. The static theory is a linear one: uniformly doubled displacements always result from uniformly doubled loads, and, more generally, from displacements \mathbf{u}^1 , \mathbf{u}^2 corresponding to stresses σ^1 , σ^2 , assigned forces \mathbf{f}^1 , \mathbf{f}^2 , and assigned surface loads σ_N^1 , σ_N^2 we construct a displacement $\mathbf{u} \equiv \mathbf{u}^1 - \mathbf{u}^2$ answering to the stress $\sigma = \sigma^1 - \sigma^2$, force $\mathbf{f} = \mathbf{f}^1 - \mathbf{f}^2$, and surface load $\sigma_N = \sigma_N^1 - \sigma_N^2$.

Let us assume that such an equilibrium state is reached in the moment $t = \tau$. Then $K(\tau) = 0$ and hence $U(\tau) = A(\tau)$. Since the potential energy $U(\tau)$ does not depend on the "path" of the deformation, but only on the final deformation, we may choose an arbitrary "path" of deformation. Let us choose the mass force components f_i , the stress components σ_i and the components u_i of the displacement in the following way:

In the time interval $0 \leq t \leq \varepsilon$: $f_i(\mathbf{x}, t) = 0$, $\sigma_i(\mathbf{x}, t) = 0$ and $u_i(\mathbf{x}, t) = 0$,
in the interval $\varepsilon \leq t \leq \tau - \varepsilon$:

$$f_i(\mathbf{x}, t) = f_i \frac{t - \varepsilon}{\tau - 2\varepsilon}, \quad \sigma_i(\mathbf{x}, t) = \sigma_i \frac{t - \varepsilon}{\tau - 2\varepsilon}, \quad \text{and} \quad u_i(\mathbf{x}, t) = u_i \frac{t - \varepsilon}{\tau - 2\varepsilon},$$

in the interval $\tau - \varepsilon \leq t \leq \tau$: $f_i(\mathbf{x}, t) = f_i$, $\sigma_i(\mathbf{x}, t) = \sigma_i$ and $u_i(\mathbf{x}, t) = u_i$,
where by f_i , σ_i and u_i we denote the values of these functions at the moment $t = \tau$ and depend only on the position \mathbf{x} . They satisfy the equations of equilibrium and so the functions $f_i(\mathbf{x}, t)$, $\sigma_i(\mathbf{x}, t)$ and $u_i(\mathbf{x}, t)$, defined above satisfy the equations of motion. Then from equations (25) and (23) we obtain

$$U(\tau) = \int_0^\varepsilon W dt + \int_\varepsilon^{\tau-\varepsilon} W dt + \int_{\tau-\varepsilon}^\tau W dt = \int_\varepsilon^{\tau-\varepsilon} W dt = \frac{1}{2} \int_V \rho f_i u_i dV + \frac{1}{2} \int_S \sigma_i u_i dS,$$

because the velocity $\partial u_i / \partial t = 0$ outside the interval $[\varepsilon, \tau - \varepsilon]$, as a consequence of the choice we made on t in the definitions of the functions $f_i(\mathbf{x}, t)$, $\sigma_i(\mathbf{x}, t)$ and $u_i(\mathbf{x}, t)$ above. In this way we arrive at the **formula of Clapeyron** [4] from 1834, asserting that *the potential energy of the deformation equals half of the work which the external forces (mass forces and surface forces) would have done, if they had from the beginning the values which they acquire at the deformed equilibrium stage.*

Solving even equilibrium problems of the linear theory of elasticity often brings significant difficulties. This is due primarily to the form of the boundary conditions. The **principle of St Venant** is helpful in many such situations. *This principle applies to the difference in the stresses and the difference in the deformations inside the body, which result from two different, but statically equivalent systems of surface forces, applied at some portion of the boundary. According to this principle, in domains sufficiently far from this part of the boundary, the difference in the stresses and that in the deformations is ignorably small.*

In 1859 Kirchhoff [19] establishes the uniqueness of the solution to boundary value problems of equilibrium where the stress vector and the displacement are prescribed upon disjoint surfaces S_1 and S_2 , respectively, such that the closure of $S_1 + S_2$ is the complete boundary of a finite body V . The displacement \mathbf{u} is determined uniquely to within an infinitesimal rigid displacement. He published these results also in 1876 in [20].

There is a remarkable variational principle enabling us, in the case of equilibrium subject to given surface displacements and vanishing assigned force in the interior, to select among all kinematically possible deformations that one which satisfies the equations of the theory of elasticity, when a positive definite elastic potential is given. The first to recognize its significance was Kelvin, who in 1863 expressed it as *“the elementary condition of stable equilibrium”*. As a proved theorem of linear three-dimensional elasticity, it was first given by Love [26] in 1906 : *“The displacement that satisfies the equations of equilibrium as well as the conditions at the boundary surface yields a smaller value for the total stored energy that does any other displacement satisfying the same conditions at the bounding surface.”*

For a review of the two-dimensional linear elastic problem and that for cylindrical bodies the reader is referred to Ivanov [18].

5. THE PROPAGATION OF WAVES

Having given consideration to static problems, let us now consider the propagation of waves.

In Continuum Mechanics waves are described as “singularities” across the two sides of a geometric two-dimensional surface that propagates in space. Such surfaces are called **singular**. To make things specific, consider a family of surfaces given by

$$\mathbf{x} = \mathbf{x}(p_1, p_2, t), \quad (27)$$

where p_1, p_2 are a pair of surface parameters, identifying what we shall call a surface point. The velocity of the surface point, identified in this way, is

$$\left. \frac{\partial \mathbf{x}}{\partial t} \right|_{p_1, p_2 = \text{const.}} \quad (28)$$

By eliminating the parameters, we may write (27) in the form

$$\alpha(\mathbf{x}, t) = 0 \quad (29)$$

for some function α . Define the normal component v_n of the velocity of a moving surface by the scalar product

$$v_n \equiv \left. \frac{\partial \mathbf{x}}{\partial t} \right|_{p_1, p_2 = \text{const.}} \cdot \mathbf{n} = - \frac{\frac{\partial \alpha}{\partial t}}{\sqrt{\alpha_{,i} \alpha_{,i}}} \quad (30)$$

where \mathbf{n} is the unit normal to the surface. v_n is called **the speed of displacement** of the surface. The velocity $v_n \mathbf{n}$ is the **normal velocity** of the surface.

Let Ψ be a function defined on the surface, we may for our purposes consider it scalar, vector or tensor-valued. If Ψ undergoes an abrupt change in its value from one side of the surface to the other, the surface is called **a singular surface with respect to the tensor Ψ** . The jump in value of Ψ is denoted by $[\Psi]$. Hugoniot-Duhem theorem states that: *The speed of displacement of a singular surface across which Ψ and its derivatives of orders $1, \dots, p-1$ are continuous, but at least one p -th derivative of Ψ is discontinuous is determined up to sign by the ratio of the jump of $\partial^p \Psi / \partial t^p$ to that of the normal p -th derivative, $\partial^p \Psi / \partial n^p$.*

Let us now recall the material representation of a moving surface. If we express the Eulerian coordinates \mathbf{x} via the Lagrangian coordinates \mathbf{X} and substitute them in the definition (29) of the surface, we obtain $S(\mathbf{X}, t) \equiv \alpha(\mathbf{x}(\mathbf{X}, t))$. In the latter representation, which we denote by $S(t)$, we may consider the medium particles as stationary and the surface $S(t)$ moving amongst them, being occupied by a different set of particles at each time t . **The speed of propagation** of the wave is

$$V_N \equiv - \frac{\frac{\partial S}{\partial t}}{\sqrt{S_{,i} S_{,i}}}$$

This speed is a measure of the rate at which the moving surface $S(t)$ traverses the material.

A surface that is singular with respect to some quantity and that has a nonzero speed of propagation is called **a propagating singular surface** or **a wave**.

Above we defined a singular surface with respect to an arbitrary quantity Ψ . Duhem proposed to regard all quantities associated with a motion as functions of the material variables \mathbf{X} and t and to define the **order** of a singular surface with respect to Ψ as the order of the derivative of Ψ of the lowest order suffering a non-zero jump upon the surface.

Some of the most interesting singularities are included in the case when

$$\Psi \equiv \mathbf{x}(\mathbf{X}, t),$$

i.e., are surfaces across which the motion itself, or one of its derivatives, is discontinuous. Surfaces across which at least one of the functional relations $\mathbf{x} = \mathbf{x}(\mathbf{X}, t)$ or $\mathbf{X} = \mathbf{X}(\mathbf{x}, t)$ defining the motion itself is discontinuous are singularities of order zero; those across which some of the first derivatives of \mathbf{x} are discontinuous are of first order, etc.

For a singular surface of order 1, we put $\Psi = x_i$ and obtain

$$[x_{i,k}] = s_i N_k \quad s_i = [N_m x_{i,m}], \quad [\dot{x}_i] = -V_N s_i.$$

Here N_k are the components of the unit normal \mathbf{N} to the surface $\mathbf{x} = \mathbf{S}(\mathbf{X}, t)$ defining the motion and equal

$$N_k = \frac{S_{,k}}{\sqrt{S_{,m} S_{,m}}}.$$

The vector \mathbf{N} is the **normal velocity** of the material. The vector \mathbf{s} , with components s_i , is the **singularity vector**. It is parallel to the jump of velocity, its magnitude varies with the choice of the initial state and thus does not furnish a measure of the strength of the singularity. The jump in the speed of propagation of a singular surface is the negative of the jump in the normal velocity of the material.

For a singular surface of order 2:

$$[x_{i,km}] = s_i N_k N_m \quad s_i = [N_k N_m x_{i,km}]. \quad (31)$$

Also

$$[\dot{x}_{i,k}] = -V_N s_i N_k, \quad [\ddot{x}_i] = V_N^2 s_i. \quad (32)$$

The formulae (31) and (32) show that a singular surface of order 2 is completely determined by a vector \mathbf{s} and the speed of propagation V_N . They show that every wave of second order carries jumps in the velocity gradient and the acceleration. Waves of second order are therefore called **acceleration waves**.

For a body of continuous constant elasticity \mathbf{C} , putting $\sigma_{mk} = C_{kmpq} \varepsilon_{pq}$ into the equations of motion $\rho \ddot{\mathbf{x}}^k = \sigma_{km,m} + \rho f^k$ yields

$$[\rho \ddot{\mathbf{x}}^k] = C_{kmpq} [\mathbf{u}_{p,qm}], \quad (33)$$

where we have supposed that $\rho \mathbf{f}$ is continuous. In linear elasticity $[u_{p,qm}] = \delta_{\alpha q} \delta_{\beta m} [x_{p,\alpha\beta}]$. By applying the general identities (31) and (32) for an acceleration wave, when the present configuration is taken as the initial one, from (33) we obtain

$$\rho V^2 s^k = C_{kmpq} n_q n_m s_p, \quad (34)$$

or

$$(C_{kmpq} n_q n_m - \rho V^2 \delta_{pk}) s_p = 0.$$

From this it follows that in order for an acceleration wave with normal \mathbf{n} to exist and propagate, the jump \mathbf{s} which it carries must be an eigenvector of $C_{kmpq} n_q n_m$ corresponding to the eigenvalue ρV^2 . For a body such that the work of the stress in any deformation is positive, as in the case for a hyperelastic body with positive definite stored energy, the tensor $C_{kmpq} n_q n_m$ is positive definite, therefore all eigenvalues ρV^2 are positive, and therefore all possible speeds are real. In the general case, *in any linearly elastic body such that the work of the stress is positive for arbitrary deformations, a wave with given normal \mathbf{n} may carry a discontinuity of the acceleration parallel to any one of three uniquely determined, mutually orthogonal directions, and corresponding to each of these directions there is a speed of propagation determined uniquely by the elasticities of the material and by \mathbf{n} .*

When the eigenvalues ρV^2 are not distinct, the above conclusion must be modified, as is seen most easily by considering the isotropic case, for then (33) assumes the more special form

$$[\rho \ddot{\mathbf{x}}^k] = (\lambda + \mu) [\mathbf{u}_{p,pk}] + \mu [\mathbf{u}_{k,pp}],$$

so that for an acceleration wave we have

$$\rho V^2 s_k = (\lambda + \mu) s_p n_p n_k + \mu s_k,$$

specializing (34). Taking the scalar and vector products of this equation by \mathbf{n} yields

$$(\rho V^2 - (\lambda + 2\mu)) \mathbf{s} \cdot \mathbf{n} = 0, \quad (\rho V^2 - \mu) \mathbf{s} \times \mathbf{n} = 0.$$

If $\mathbf{s} \cdot \mathbf{n} \neq 0$, the first of these equations yields $\rho V^2 = \lambda + 2\mu$, and the second, if we exclude the case when $\lambda + \mu = 0$, yields $\mathbf{s} \times \mathbf{n} = 0$. If $\mathbf{s} \cdot \mathbf{n} = 0$, but $\mathbf{s} \times \mathbf{n} \neq 0$, the second equation yields $\rho V^2 = \mu$. Summarizing these results, we see that *in an isotropic linearly elastic body for which $\lambda + \mu \neq 0$, a necessary and sufficient condition that the acceleration waves be propagated at positive speeds is $\lambda + 2\mu > 0$, $\mu > 0$. This condition is satisfied when the stored energy is positive definite. Two kinds of acceleration waves are possible: longitudinal waves, whose speed of propagation is given by*

$$V^2 = (\lambda + 2\mu)/\rho,$$

and transverse waves, for which

$$V^2 = \mu/\rho.$$

The foregoing results were first obtained by Christoffel [5] in 1877 and independently by Hugoniot [16] in 1886. These results demonstrate the far-reaching effect of isotropy: instead of three speeds of propagation, for an isotropic body there are only two, but instead of there being only three possible directions for the discontinuity, there are infinitely many, though the possible directions are still far from arbitrary.

6. PERFECT FLUIDS

A continuous medium is a **perfect fluid** if it can support no shearing stress and no couple stress. As a consequence of these restrictions, the stress tensor σ is hydrostatic, $\sigma = -p\mathbf{1}$, and from Cauchy's first law of motion the **dynamical equation of Euler** is obtained

$$\rho \ddot{\mathbf{x}} = -\text{grad } p + \rho \mathbf{f}. \quad (35)$$

Euler published this equation in 1757, see [8]. Cauchy's second law is satisfied automatically, in other words, balance of linear momentum in a perfect fluid implies balance of moment of momentum, as long as there are no extrinsic couples, while if there are such present, the perfect fluid is incompatible with the principles of mechanics. Hugoniot in 1887 [17] Part I, Hadamard in 1903 [13] and Duhem in 1901 [7] Part II, Chap. IV, proved that a perfect fluid admits only longitudinal waves. Hadamard [13] and Duhem [7] Part II, Chap. I, proved that in an isochoric motion of a perfect fluid wave propagation of any kind is impossible.

In 1869 Kelvin proved that: *"A flow of a perfect fluid subject to lamellar assigned force is circulation preserving if and only if there exists a functional relation*

$$f(p, \rho, t) = 0; \quad (36)$$

alternatively, if and only if, for each fixed time, the pressure is constant, or the density is constant, or the surfaces $p = \text{const.}$ coincide with the surfaces $\rho = \text{const.}$ " Kelvin's theorem is regarded as the fundamental theorem of classical hydrodynamics. Flows satisfying (36) are called **barotropic**.

A perfect fluid may be such that all its flows are barotropic; this is the case for homogeneous incompressible fluids, for which $\rho = \text{const.}$ in space and time, and for piezotropic fluids, for which there is an equation of state of the form $p = f(\rho)$. But these conditions are merely sufficient, not necessary for barotropic flow. For example, in a fluid having equations of state $p = F(\rho, \theta) = G(\rho, \eta)$, special conditions may lead to a flow for which $\theta = \text{const.}$ or for which $\eta = \text{const.}$ Any such flow is barotropic, but the functional form of f in (36) depends upon the particular conditions giving rise to the flow.

When (36) holds, all the numerous theorems appropriate to circulation preserving motion may be applied: the Helmholtz vorticity theorems, the Bernoullian theorems and the Helmholtz theorem of conservation of energy. Indeed, all general theorems of classical hydrodynamics follow from the circulation preserving property.

It should be noted also that *In the case of a barotropic flow, the speed of propagation of acceleration waves*

$$[\ddot{x}_n] = -c^2 \left[\frac{d \log \rho}{dn} \right]$$

is c , where

$$c^2 \equiv \frac{\partial p}{\partial \rho}.$$

This is Hugoniot's theorem, published in 1885, see [15] and [17], Part I.

For barotropic flows in which neither the pressure nor the density is uniform, a necessary and sufficient condition for wave propagation to be possible is that the pressure be an increasing function of the density; this being so, waves of all orders greater than 1 propagate with the unique speed c . Since c is the common speed of propagation of so many kinds of waves, it is called the **speed of sound**.

7. PROBLEMS

In this section we would like to apply the ideas presented so far in order to solve some concrete problems.

Problem 1. A rectangular tank containing a nonviscous liquid of constant density moves horizontally to the right with a constant acceleration. Gravitational force is the only external force. Find the pressure distribution in the liquid and the geometrical shape of the upper surface of the liquid.

Solution. Choose the positive x -direction of the coordinate system to be the direction in which the tank moves, and the positive z -direction to be the vertical direction upward. Then, $d\mathbf{v}/dt = a\mathbf{e}_1$ and $\mathbf{b} = -g\mathbf{e}_3$, where $a = |d\mathbf{v}/dt|$ is a constant and g is the (constant) acceleration due to gravity. Euler's equation

$$\frac{d\mathbf{v}}{dt} = -\frac{1}{\rho}\nabla p + \mathbf{b},$$

where \mathbf{b} is the body force, yields the following three equations for the three Cartesian components (x, y, z) of the ∇p :

$$\frac{\partial p}{\partial x} = -a\rho$$

$$\frac{\partial p}{\partial y} = 0$$

$$\frac{\partial p}{\partial z} = -g\rho.$$

The second of these three equations shows that p is independent of y , and thus has the form

$$p = -\rho ax + f(z)$$

where $f(z)$ is an arbitrary function of z . From this form of p and the third component of ∇p above, we see that $f(z) = -\rho gz + C$, where C is a constant, thus arriving at

$$p = -\rho(ax + gz) + C.$$

At the point where the z -axis meets the upper surface of the liquid, we have $p = p_a$, where p_a is the atmospheric pressure. If this point is at a height h above the origin, the last equation for p gives $C = p_a + \rho gh$. Thus the pressure distribution of the liquid is

$$p = p_a - \rho(ax + gz - gh).$$

For $p = p_a$, the last equation for p becomes

$$z = -\left(\frac{a}{g}\right)x + h.$$

This is the shape of the upper surface of the liquid. Evidently, this surface is a plane, making an acute angle $\theta = \tan^{-1}(a/g)$ with the horizontal. In the limiting case when $a \rightarrow 0$, the liquid moves with a constant velocity and the upper surface of the liquid becomes a horizontal plane.

The interested reader is invited to apply the method of solution of the last problem in order to solve

Problem 2. A column of a nonviscous liquid of constant density contained in a vertical circular vessel rotates like a rigid body about the axis of the vessel with a constant angular velocity ω . Gravitational force is the only external force. Find the pressure distribution in the liquid and the geometrical form of the upper surface of the liquid.

In the next problem we will use Bernoulli's equation

$$\frac{\partial \mathbf{v}}{\partial t} + \mathbf{w} \times \mathbf{v} = -\nabla H,$$

where \mathbf{w} is the vorticity vector and $H \equiv P + \chi + v^2/2$. Here P is

$$P = \int \frac{1}{\rho} dp,$$

with p being the pressure. This equation is known after Daniel Bernoulli (1738). It is the equation of motion for an elastic fluid moving under conservative body force $-\nabla\chi$. Since Bernoulli's equation holds for an elastic fluid for which $\rho = \rho(p)$, it automatically holds in the special case of $\rho = \text{constant}$.

Problem 3. For a certain flow of a nonviscous fluid of constant density under the Earth's gravitational field, the velocity distribution is given by $\mathbf{v} = \nabla\phi$, where $\phi = x^3 - 3xy^2$. Find the pressure distribution.

Solution. From the given \mathbf{v} , we find that $\text{curl } \mathbf{v} = \mathbf{0}$ and $\partial\mathbf{v}/\partial t = \mathbf{0}$. Thus the fluid is irrotational and steady. Then $\partial\mathbf{v}/\partial t = \mathbf{0}$ and either the vorticity vector $\mathbf{w} = \mathbf{0}$ or $\mathbf{v} \times \mathbf{w} = \mathbf{0}$. Further, since the body force is the gravitational force, it is conservative. With these observations, Bernoulli's equation

$$\frac{\partial \mathbf{v}}{\partial t} + \mathbf{w} \times \mathbf{v} = -\nabla H,$$

where H is the Bernoulli's function $H \equiv P + \chi + v^2/2$, reduces to

$$\nabla H = \mathbf{0}.$$

Since $\partial H/\partial t = 0$, it follows that

$$H \equiv P + \chi + \frac{v^2}{2} = \text{constant}$$

everywhere in the fluid. Thus, under the assumed conditions, the function H is an integral of the equation of motion of the elastic fluid.

Since the body force is conservative, $\chi = gz$, where z is measured vertically upward. Accordingly, from $H = P + \chi + v^2/2 = \text{constant}$ with $P = p/\rho$ we obtain

$$\frac{p}{\rho} + \frac{1}{2}v^2 + gz = C,$$

where C is a constant.

From the given \mathbf{v} we also find

$$v_1 = \frac{\partial \phi}{\partial x} = 3(x^2 - y^2), \quad v_2 = \frac{\partial \phi}{\partial y} = -6xy, \quad v_3 = 0$$

and hence

$$v^2 = v_1^2 + v_2^2 = 9(x^2 + y^2)^2.$$

Substituting this result into the equation relating v^2 , p and z (above), we obtain

$$\frac{p}{\rho} + \frac{9}{2}(x^2 + y^2)^2 + gz = C.$$

From this result it is evident that $C = p^0/\rho$, where p^0 is the pressure at the origin. Thus,

$$p = p^0 - \rho \left(\frac{9}{2}(x^2 + y^2)^2 + gz \right)$$

is the sought pressure distribution.

Many interesting problems in Continuum Mechanics can be found in the book of Chandrasekharaiah and Debnath [3].

8. LINEARLY VISCOUS FLUIDS

Let us now consider a medium which in equilibrium, satisfies Euler's equation (35)

$$\text{grad } p = \rho \mathbf{f},$$

but when in motion can support appropriate shearing stresses. More specifically, let us assume that the stress tensor σ is a linear function of the velocity $\dot{\mathbf{x}}$ and the velocity gradient, namely

$$\sigma = \mathbf{g}(\dot{\mathbf{x}}, \mathbf{w}, \mathbf{d}), \quad (37)$$

where \mathbf{g} is a linear function. Here \mathbf{d} is Euler's stretching tensor and \mathbf{w} is Cauchy's spin tensor. The constitutive equations (37) define **linearly viscous fluid**. By applying the principle that the constitutive equation must have the same form for all observers, one shows that in fact σ is independent of $\dot{\mathbf{x}}$ and \mathbf{w} , i.e.,

$$\sigma = \mathbf{g}(\mathbf{d}), \quad (38),$$

with the function \mathbf{g} being linear. This equation in an internal frame, along with $\mathbf{f}(0) = -p\mathbf{1}$, was taken as the definition of a fluid by Stokes [35] in 1845. If we now use a coordinate system with axes that coincide with the principal directions of \mathbf{d} , so that (38) becomes

$$\sigma_{km} = f_{km}(d_1, d_2, d_3)$$

it is easily seen that *the principal axes of stretching are also principal axes of stress*. Another interesting property of fluids included in the definition (37) is that such fluids are necessarily isotropic.

The most general linear isotropic function σ of a symmetric second order tensor \mathbf{d} may be written in the form of **Navier-Poisson law**:

$$\sigma = -p\mathbf{1} + \lambda \mathbf{I}_d \mathbf{1} + 2\mu \mathbf{d}$$

or in components

$$\sigma_{km} = -p\delta_{km} + \lambda d_{qq}\delta_{km} + 2\mu d_{km}, \quad (39)$$

where a use is made of the requirement that $\sigma = -p\mathbf{1}$ when $\mathbf{d} = 0$. Historically, the simplest case of this law was proposed by Newton [31], Lib. II, Chap. IX. It follows from (39) that σ is symmetric, thus Cauchy's second law is automatically satisfied. Thus, for the fluids in question, ballance of momentum implies ballance of moment of momentum. Substitution of (39) into Cauchy's first law of motion

$$\sigma_{ij,j} + \rho f_i = \rho \frac{dv_i}{dt}, \quad i = 1, 2, 3$$

yields a system of three differential equations, known, when subjected to further simplifying assumptions, as **the Navier-Stokes equations**:

$$\mu \nabla^2 \mathbf{v} + (\lambda + \mu) \nabla(\operatorname{div} \mathbf{v}) - \nabla p + \rho \mathbf{f} = \rho \frac{d\mathbf{v}}{dt}. \quad (40)$$

They are attributed to Navier (1822) and Stokes (1845) and hold for both compressible and incompressible viscous fluid flows; in the incompressible case $\rho = \rho_0$ and $\operatorname{div} \mathbf{v} = 0$.

The coefficients λ and μ are the **viscosities** of the fluid. In the absence of viscosity, that is if λ and μ are negligibly small, this equation reduces to Euler's equation of motion (35) for perfect fluids. Because of that perfect fluids are often called **inviscid**.

The portion $\lambda \mathbf{I}_d \mathbf{1} + 2\mu \mathbf{d}$ of the stress is considered as arising from internal friction. Because

$$\sigma_{km} = 2\mu d_{km} \quad \text{when } k \neq m, \quad (41)$$

μ is the ratio of the shear stress to the corresponding shearing of any two orthogonal elements, and so is called the **shear viscosity**.

The stress power assumes the form

$$P = \sigma_{km} d_{km} = -p d_{kk} + \lambda (d_{kk})^2 + 2\mu d_{km} d_{mk}.$$

In 1850 Stokes had shown in [36] that for the fluids in consideration

$$\mu \geq 0, \quad 3\lambda + 2\mu \geq 0.$$

The same conclusion was reached independently by Duhem in 1901, published in [7], Part I. These inequalities have some significant mechanical consequences. For example, equations (40) with $\mu \geq 0$ imply that the shear stress always opposes the shearing. These consequences show that the effect of the viscous stress $\sigma + p\mathbf{1}$, as given by (39) is always to resist change of shape, and thus is of the nature of frictional resistance.

For an incompressible viscous fluid, the Navier-Stokes equations (40) are rewritten in the form:

$$\frac{\mu}{\rho} \nabla^2 \mathbf{v} + \frac{1}{\rho} \nabla p + \mathbf{f} = \frac{d\mathbf{v}}{dt}. \quad (42)$$

The coefficient μ/ρ is called kinematic viscosity.

The presence of viscosity has the effect of making the propagation of most kinds of waves impossible. In 1926 Kotchine [21] proved that the instantaneous existence of a surface upon which $\dot{\mathbf{x}}$ and p are continuous but $\dot{x}_{k,m}$ suffers a jump discontinuity is incompatible with the law of linear viscosity (39). His result is contained in an earlier one of Duhem from 1901 [7], Part II, Chap. III, who uses a different terminology. In [6] and [7], Part II, Chap. III, Duhem asserts that in a linearly viscous fluid no waves of order greater than 1 are possible.

A summary of the existing knowledge of the theories of non-linear viscosity is given in [38]. An excellent text from the latter part of the 20th Century, which the reader can use to get acquainted with the modern developments of the presented theories, is the book of Timoshenko and Goodier [37].

For the readers privileged to know Russian, we list two excellent texts on hydrodynamics [40], [41]. They can be used to deepen knowledge in the theories presented in this review. Two prominent texts in Bulgarian on hydrodynamics are the book of Zaprianov and that of Shkadov and Zaprianov, [39] and [42].

9. CONDITIONS FOR COMPATABILITY OF THE STRESSES

The conditions for compatability of the stresses are due to Beltrami [1], and were independently discovered by Mitchell in 1899 [27].

We consider the static case in which the equations of motion

$$\sigma_{ij,j} + \rho f_i = \rho \frac{\partial^2 u_i}{\partial t^2}$$

become the **equations for equilibrium**

$$\sigma_{ij,j} + F_i = 0, \tag{43}$$

where $\mathbf{F} = r\mathbf{f}$. These equations form a system of 3 equations for the 6 unknowns σ_{ij} (we assume that the volume forces \mathbf{f} are given). This system has infinitely many solutions, but not every one of them corresponds to a real deformation, from which we can calculate the displacement \mathbf{u} in the medium. As we know, for the deformations, determined using equation (17), it is necessary and sufficient to satisfy the conditions for the compatability of the deformations of St. Venant

$$\varepsilon_{ij,kl} + \varepsilon_{kl,ij} - \varepsilon_{ik,jl} - \varepsilon_{jl,ik} = 0. \tag{44}$$

Let us express these conditions with the stresses. Let's substitute the components of the tensor of deformations, using equations (17), into the conditions (44), and then introduce the notation $\sigma = \sigma_{kk}$. We obtain

$$\sigma_{ij,kl} + \sigma_{kl,ij} - \sigma_{ik,jl} - \sigma_{jl,ik} = \frac{\nu}{1+\nu}(\sigma_{,ij}\delta_{kl} + \sigma_{,kl}\delta_{ij} - \sigma_{,ik}\delta_{jl} - \sigma_{,jl}\delta_{ik}). \tag{45}$$

If we set $k = l$ and sum over the repeated index, we will arrive at the following system of equations:

$$\sigma_{ij,kk} + \sigma_{,ij} - \sigma_{ik,jk} - \sigma_{jk,ik} = \frac{\nu}{1+\nu}(\sigma_{,ij} + \sigma_{,kk}\delta_{ij}). \tag{46}$$

This system consists of 9 equations, from which independent are only 6. We can obtain these 6 independent equations if we let, for example, $i \geq j$, because of the symmetry with respect to the indexes i and j . The system (46), obtained in this manner, is equivalent to the initial system (45), because each system consists of 6 independent equations, and the equations of system (46) are linear combinations of the equations of system (45).

Let us differentiate the equations (43) with respect to x_k . We obtain

$$\sigma_{ij,jk} = -F_{i,k}. \tag{47}$$

Substitute (47) in (46) to obtain

$$\sigma_{ij,kk} + \frac{\nu}{1+\nu}\sigma_{,ij} - \frac{\nu}{1+\nu}\sigma_{,kk}\delta_{ij} = -(F_{i,j} + F_{j,i}). \tag{48}$$

We are now going to simplify this system in the following way: We set in (45) $k = i$ and $l = j$ and after a few calculations obtain

$$\sigma_{ij,ij} = \frac{1 - \nu}{1 + \nu} \sigma_{,ii}. \quad (49)$$

Now using (47), we can write equation (49) as

$$\sigma_{,ii} = -\frac{1 + \nu}{1 - \nu} F_{i,i}.$$

Substituting this result in (48), we finally obtain **the conditions of Beltrami-Mitchell for the compatability of the stresses**:

$$\sigma_{ij,kk} + \frac{1}{1 + \nu} \sigma_{,ij} = -\frac{\nu}{1 - \nu} F_{k,k} \delta_{ij} - (F_{i,j} + F_{j,i}).$$

ACKNOWLEDGEMENTS. I am indebted to Professor Tsolo Ivanov, Professor Emeritus at the Department of Mechatronics, Robotics and Mechanics, Sofia University St Kliment Ohridski, Bulgaria, for valuable conversations in continuum mechanics.

I am also thankful to Professor George Boyadjiev, Head of the Department of Mechatronics, Robotics and Mechanics, Sofia University St Kliment Ohridski, Bulgaria, for encouraging me to write this paper.

10. REFERENCES

- [1] Beltrami, E.: Osservazioni sulla nota prec endente. *Rend. Lincei* (5) **1** (1892), 141–142.
- [2] Bernoulli, D.: *Hydrodynamica sive de Viribus. Motibus Fluidorum Commentarii. Argentorati.* (1738).
- [3] Chandrasekharaiah, D.S., Debnath, L.: *Continuum Mechanics*, Academic Press, 1994.
- [4] Clapeyron, E.: Mèmoire sur la pissance motrice de la chaleur. *J. École Polytech.* **14** (1834), cahier 23, 153–190.
- [5] Cristoffel, E.B.: Über die Fortpflanzung von Stöben durch elastische feste Körper. *Ann. Math.* (2) **8** (1877), 193–244.
- [6] Duchem, P.: Sur les théorèmes d'Hugoniot, les lemmes de Hadamard et la propagation des ondes dans les fluides visqueux. *C. R. Acad. Sci., Paris* **132** (1901), 1163–1167.
- [7] Duchem, P.: Recherches sur lhydrodynamique. *Ann. Toulouse* (2) **3** (1901), 315–377, 379–431; **4** (1902), 101–169; **5** (1903), 5–61, 197–255, 353–404.
- [8] Euler, L.: 1757 Principes généraux du mouvement des fluides. *Mèm. Acad. Sci. Berlin* **11**, 274–315 (1755) = *Opera omnia* (2) **12**, 54–91.

- [9] Euler, L.: Sectio secunda de principiis motus fluidorum. *Novi Comm. Acad. Sci. Petrop.* **14** (1769), 270–386 = *Opera omnia* (2) **13**, 73–153.
- [10] Georgieva, B.: Review of continuum mechanics and its history. Part I. Deformation and stress. Conservation laws. Constitutive equations. *Ann. Sofia Univ., Fac. Math and Inf.* **107** (2020), 29–54.
- [11] Green, G.: On the laws of reflection and refraction of light at the common surface of two non-crystalized media. *Trans. Cambridge Phil. Soc.* **7** (1839), 1–24.
- [12] Green, G.: On the propagation of light in crystalized media. *Trans. Cambridge Phil. Soc.* **7** (1841), 121–140.
- [13] Hadamard, J.: *Leçons sur la propagation des ondes et les équations de l'hydrodynamique.* (Lectures of 1898-1900, Paris).
- [14] Hooke, R.: Lectures de Potentia Restitutiva, or of Spring Explaining the Power of Springing Bodies. London (1678) = R.T. Gunther: *Early Science in Oxford* **8** (1931), 331–356.
- [15] Hugoniot, H.: Sur la propagation du mouvement dans un fluide indéfini. (Première Partie). *C. R. Acad. Sci., Paris* **101** (1885), 1118–1120.
- [16] Hugoniot, H.: Sur un théorème general relatif a la propagation du mouvement. *C.R. Acad. Sci., Paris* **102** (1886), 858–860.
- [17] Hugoniot, H.: Mémoire sur la propagation du mouvement dans un fluide indéfini. *J. Math. Pures Appl.* (4) **3** (1887), 477–492; (4) **4** (1887), 153–167.
- [18] Ivanov, Ts.: *Theory of Elasticity.* Sofia University Publishing House (1995).
- [19] Kirchhoff, G.: Über das Gleichgewicht und die Bewegung eines unendlich dunnen elastischen Stabes. *J. Reine Angew. Math.* **56** (1859), 285–313.
- [20] Kirchhoff, G.: *Vorlesungen über mathematische Physik: Mechanik.* Leipzig, 2nd ed. 1877; 3ed ed. 1883.
- [21] Kotchine, N. E.: Sur la théorie des ondes de choc dans un fluide. *Rend. Circ. Mat. Palermo* **50** (1926), 305–344.
- [22] Lagrange, J. L.: *Mechanique Analytique*, Paris, *Oeuvres* **11**, **12** (1788).
- [23] Lamè, G., Clapeyron, E.: Mémoire sur l'équilibre des corps solides homogènes. *Mem. divers savants* (2) (1833) **4**, 465–562.
- [24] Lamè, G.: Mémoire sur les surfaces isostatiques dans les corps solides homogènes en équilibre d'élasticité. *J. Math. Pures Appl.* **6** (1841), 37–60.
- [25] Lamè, G.: *Leçons sur la Théorie Mathématique de l'Élasticité.* Paris 1852, Second ed.: Paris 1866.
- [26] Love, A. E. H.: *A Treatise on the Mathematical Theory of Elasticity*, vol. 1, 2nd ed., Cambridge 1892.
- [27] Michell, J. H.: On the direct determination of stress in an elastic solid, with applications to the theory of plates. *Proc. London Math. Soc.* **31** (1899), 100–124.
- [28] Navier, C. L. M. H.: Sur les lois de l'équilibre et du mouvement des corps solides élastiques. *Bull. Soc. Philomath.* (1821), 177–181.
- [29] Navier, C. L. M. H.: Mémoire sur les lois de l'équilibre et du mouvement des corps solides élastiques. *Mém. Acad. Sci. Inst. France* (2) **7** (1821), 375–393.

- [30] Navier, C.L.M.H.: Sur les lois des mouvements des fluides, en ayant egard a l'adhesion des molecules. *Ann. Chimie* **19** (1821), 244–260.
- [31] Newton, I.: *Philosophiae Naturalis Principia Mathematica*. London, 3rd ed. (H. Pemberton, Ed.), London 1726.
- [32] Poisson, S. D.: Mèmoire sur l'èquilibre et le mouvement des corps èlastiques . *Mèm. Acad. Sci. Inst. France* (2) **8** (1828), 357–570.
- [33] St Venant, A.-J.-C. B.: Sur les pressions qui se dèveloppent a l'intèrieur des corps solides lorsque les dèplacements de leurs points, sans altèrer l'èlasticitè, ne peuvent cependant pas ètre considèrès comme très petits. *Bull. Soc. Philomath.* **5** (1844), 26–28.
- [34] St Venant, A.-J.-C. B.: Theorie de l'elasticite des solides, ou cinematique de leurs deformations. *L'Institut* **32**¹ (1864), 389–390.
- [35] Stokes, G.G.: On the theories of the internal friction of fluids in motion, and of the equilibrium and motion of elastic solids. *Trans. Cambridge Phil. Soc.* **8** (1844-1849), 287–319.
- [36] Stokes, G.G.: On the effect of the internal friction of fluids on the motion of pendulums . *Trans. Cambridge Phil. Soc.* **9**² (1850), 8–106 = Papers **3**, 1–141.
- [37] Timoshenko, S. P., Goodier, J. N.: *Theory of Elasticity*, Third Ed.. McGraw-Hill, N.Y., 1970.
- [38] Truesdell, C., Toupin, R.: *The Non-linear Field Theories*. Encyclopedia Of Physics, Vol. III/. Springer-Verlag, Berlin, 1965.
- [39] Запрянов, З.: *Хидродинамика*, Университетско издателство "Св. Климент Охридски", София, 1996.
- [40] Ландау, Л. Д., Лифшиц, Е. М.: *Гидродинамика*. Второе издание, Гос. изд. техн.-теор. лит., Москва, 1953. (Превод на български език: Ландау, Л. Д., Лифшиц Е. М.: *Хидродинамика*, Наука и изкуство, София, 1978).
- [41] Седов, Л. И.: *Механика сплошной среды*, том 1, Наука, Москва, 1970.
- [42] Шкадов, В., Запрянов, З.: *Динамика на вискозни флуиди*. Наука и изкуство, София, 1986.

Received on June 25, 2021

BOGDANA A. GEORGIEVA
 Faculty of Mathematics and Informatics
 Dept. of Mechatronics, Robotics and Mechanics
 Sofia University St Kliment Ohridski
 5 James Bourchier Blvd.
 1164 Sofia
 BULGARIA
 E-mail: georgieva@fmi.uni-sofia.bg

ГОДИШНИК НА СОФИЙСКИЯ УНИВЕРСИТЕТ „СВ. КЛИМЕНТ ОХРИДСКИ“

ФАКУЛТЕТ ПО МАТЕМАТИКА И ИНФОРМАТИКА

Том 107

ANNUAL OF SOFIA UNIVERSITY „ST. KLIMENT OHRIDSKI“

FACULTY OF MATHEMATICS AND INFORMATICS

Volume 107

UNIVARIATE POLYNOMIALS AND THE CONTRACTABILITY OF CERTAIN SETS

VLADIMIR PETROV KOSTOV

We consider the set Π_d^* of monic polynomials $Q_d = x^d + \sum_{j=0}^{d-1} a_j x^j$, $x \in \mathbb{R}$, $a_j \in \mathbb{R}^*$, having d distinct real roots, and its subsets defined by fixing the signs of the coefficients a_j . We show that for every choice of these signs, the corresponding subset is non-empty and contractible. A similar result holds true in the cases of polynomials Q_d of even degree d and having no real roots or of odd degree and having exactly one real root. For even d and when Q_d has exactly two real roots which are of opposite signs, the subset is contractible. For even d and when Q_d has two positive (resp. two negative) roots, the subset is contractible or empty. It is empty exactly when the constant term is positive, among the other even coefficients there is at least one which is negative, and all odd coefficients are positive (resp. negative).

Keywords: Real polynomial in one variable, hyperbolic polynomial, Descartes' rule of signs.

2020 Math. Subject Classification: 26C10, 30C15.

1. INTRODUCTION

In the present paper we consider the general family of real monic univariate polynomials $Q_d = x^d + \sum_{j=0}^{d-1} a_j x^j$. It is a classical fact that the subsets of $\mathbb{R}^d \cong Oa_0 \dots a_{d-1}$ of values of the coefficients a_j for which the polynomial Q_d has one and the same number of distinct real roots are contractible open sets. These sets are the $[d/2] + 1$ open parts of $R_{1,d} := \mathbb{R}^d \setminus \Delta_d$, where Δ_d is the *discriminant set* corresponding to the family Q_d .

Remarks 1. (1) One defines the discriminant set by the two conditions:

(a) The set Δ_d^1 is defined by the equality $\text{Res}(Q_d, Q'_d, x) = 0$, where $\text{Res}(Q_d, Q'_d, x)$ is the resultant of the polynomials Q_d and Q'_d , i. e. the determinant of the corresponding Sylvester matrix.

(b) One sets $\Delta_d := \Delta_d^1 \setminus \Delta_d^2$, where Δ_d^2 is the set of values of the coefficients a_j for which there is a multiple complex conjugate pair of roots of Q_d and no multiple real root.

One observes that $\dim(\Delta_d) = \dim(\Delta_d^1) = d - 1$ and $\dim(\Delta_d^2) = d - 2$. Thus Δ_d is the set of values of (a_0, \dots, a_{d-1}) for which the polynomial Q_d has a multiple real root.

(2) The discriminant set is invariant under the one-parameter group of quasi-homogeneous dilatations $a_j \mapsto u^{d-j} a_j$, $j = 0, \dots, d$.

Remark 1. If one considers the subsets of \mathbb{R}^d for which the polynomial Q_d has one and the same numbers of positive and negative roots (all of them distinct) and no zero roots, then these sets will be the open parts of the set $R_{2,d} := \mathbb{R}^d \setminus (\Delta_d \cup \{a_0 = 0\})$. To prove their connectedness one can consider the mapping “roots \mapsto coefficients”. Given two sets of nonzero roots with the same numbers of negative and positive roots (in both cases they are all simple) one can continuously deform the first set into the second one while keeping the absence of zero roots, the numbers of positive and negative roots and their simplicity throughout the deformation. The existence of this deformation implies the existence of a continuous path in the set $R_{2,d}$ connecting the two polynomials Q_d with the two sets of roots.

In the present text we focus on polynomials without vanishing coefficients and we consider the set

$$R_{3,d} := \mathbb{R}^d \setminus (\Delta_d \cup \{a_0 = 0\} \cup \{a_1 = 0\} \cup \dots \cup \{a_{d-1} = 0\}).$$

We discuss the question when its subsets corresponding to given numbers of positive and negative roots of Q_d and to given signs of its coefficients are contractible.

Notation 1. (1) We denote by σ the d -tuple $(\text{sign}(a_0), \dots, \text{sign}(a_{d-1}))$, where $\text{sign}(a_j) = +$ or $-$, by \mathcal{E}_d the set of *elliptic* polynomials Q_d , i. e. polynomials with no real roots (hence d is even and $a_0 > 0$), and by $\mathcal{E}_d(\sigma) \subset \mathcal{E}_d$ the set consisting of elliptic polynomials Q_d with signs of the coefficients defined by σ .

(2) For d odd and for a given d -tuple σ , we denote by $\mathcal{F}_d(\sigma)$ the set of monic real polynomials Q_d with signs of their coefficients defined by the d -tuple σ and having exactly one real (and simple) root.

(3) For d even, we denote by $\mathcal{G}_d(\sigma)$ the set of polynomials Q_d having signs of the coefficients defined by the d -tuple σ and having exactly two simple real roots.

Remark 2. For an elliptic polynomial Q_d , one has $a_0 > 0$, because for $a_0 < 0$, there is at least one positive root. The sign of the real root of a polynomial of $\mathcal{F}_d(\sigma)$ is opposite to $\text{sign}(a_0)$. A polynomial from $\mathcal{G}_d(\sigma)$ has two roots of same (resp. opposite) signs if $a_0 > 0$ (resp. if $a_0 < 0$).

In order to formulate our first result we need the following definition:

Definition 1. (1) For d even and $a_0 < 0$, we set $\mathcal{G}_{d,(1,1)}(\sigma) := \mathcal{G}_d(\sigma)$. For d even and $a_0 > 0$, we set $\mathcal{G}_d(\sigma) := \mathcal{G}_{d,(2,0)}(\sigma) \cup \mathcal{G}_{d,(0,2)}(\sigma)$, where for $Q_d \in \mathcal{G}_{d,(2,0)}$ (resp. $Q_d \in \mathcal{G}_{d,(0,2)}$), Q_d has two positive (resp. two negative) distinct roots and no other real roots. Clearly $\mathcal{G}_{d,(2,0)}(\sigma) \cap \mathcal{G}_{d,(0,2)}(\sigma) = \emptyset$.

(2) For d even, we define two special cases according to the signs of the coefficients of Q_d and the quantities of its positive or negative real roots:

Case 1). The constant term and all coefficients of monomials of odd degrees are positive, there is at least one coefficient of even degree which is negative, and Q_d has 2 positive and no negative roots.

Case 2). The constant term is positive, all coefficients of monomials of odd degrees are negative, there is at least one coefficient of even degree which is negative, and Q_d has 2 negative and no positive roots.

Note that Cases 1) and 2) are exchanged when one performs the change of variable $x \mapsto -x$.

Our first result concerns real polynomials with not more than 2 real roots:

Theorem 1. (1) For d even and for each d -tuple σ , the subset $\mathcal{E}_d(\sigma) \subset \mathcal{E}_d$ is non-empty and convex hence contractible.

(2) For d odd and for each d -tuple σ , the set $\mathcal{F}_d(\sigma)$ is non-empty and contractible.

(3) For d even and for each d -tuple σ with $a_0 < 0$, the set $\mathcal{G}_{d,(1,1)}(\sigma)$ is contractible. For d even and for each d -tuple σ with $a_0 > 0$, each set $\mathcal{G}_{d,(2,0)}(\sigma)$ (resp. $\mathcal{G}_{d,(0,2)}(\sigma)$) is contractible or empty. It is empty exactly in Case 1) (resp. Case 2)).

The theorem is proved in Section 4. The next result of this paper concerns *hyperbolic polynomials*, i. e. polynomials Q_d with d real roots counted with multiplicity.

Notation 2. We denote by Π_d the *hyperbolicity domain*, i. e. the subset of \mathbb{R}^d for which the corresponding polynomial Q_d is hyperbolic. The interior of Π_d is the set of polynomials having d distinct real roots and its border $\partial\Pi_d$ equals $\Delta_d \cap \Pi_d$. We set

$$\Pi_d^* := \Pi_d \setminus (\Delta_d \cup \{a_0 = 0\} \cup \{a_1 = 0\} \cup \dots \cup \{a_{d-1} = 0\}).$$

Thus Π_d^* is the set of monic degree d univariate polynomials with d distinct real roots and with all coefficients non-vanishing. We denote by Π_d^k and Π_d^{*k} the projections of the sets Π_d and Π_d^* in the space $O_{a_{d-k}} \dots a_{d-1}$ (hence $\Pi_d^d = \Pi_d$ and $\Pi_d^{*d} = \Pi_d^*$), by $\partial\Pi_d^k$ the border of Π_d^k and by *pos* and *neg* the numbers of positive and negative roots of a polynomial Q_d having no vanishing coefficients.

We set $a := (a_0, a_1, \dots, a_{d-1})$, $a' := (a_1, \dots, a_{d-1})$, $a'' := (a_2, \dots, a_{d-1})$ and $a^{(k)} := (a_k, \dots, a_{d-1})$. In what follows we use the same notation for functions and for their graphs.

Remarks 2. (1) For a hyperbolic polynomial with no vanishing coefficients, the d -tuple σ defines the numbers *pos* and *neg*. Indeed, by Descartes' rule of signs a real univariate polynomial Q_d with c sign changes in its sequence of coefficients has $\leq c$ positive roots and the difference $c - \text{pos}$ is even, see [13] and [10]. When applying this rule to the polynomial $Q(-x)$ one finds that the number p of sign preservations is $\geq \text{neg}$ and the difference $p - \text{neg}$ is even. For a hyperbolic polynomial one has $\text{pos} + \text{neg} = c + p = d$, so in this case $c = \text{pos}$ and $p = \text{neg}$.

(2) By Rolle's theorem the non-constant derivatives of a hyperbolic polynomial (resp. of a polynomial of the set Π_d^*) are also hyperbolic (resp. are hyperbolic with all roots non-zero and simple). Hence for two hyperbolic polynomials of the same degree and with the same signs of their respective coefficients, their derivatives of the same orders have one and the same numbers of positive and negative roots.

Our next result is the following theorem (proved in Section 5):

Theorem 2. *For each d -tuple σ , there exists exactly one open component of the set Π_d^* the polynomials Q_d from which have exactly *pos* positive simple and *neg* negative simple roots and have signs of the coefficients as defined by σ . This component is contractible.*

One can give more explicit information about the components of the set Π_d^* . Denote by M such a component defined after a d -tuple σ and by M^k its projection in the space $Oa_{d-k} \cdots a_{d-1}$. It is shown in [19] (see Proposition 1 therein) that M is non-empty. In Section 5 we prove the following statement:

Theorem 3. *For $k \geq 3$, the set M^k is the set of all points between the graphs L_{\pm}^k of two continuous functions defined on M^{k-1} :*

$$M^k = \{a^{(d-k)} \in \mathbb{R}^{d-k} \mid L_-^k(a^{(d-k+1)}) < a_{d-k} < L_+^k(a^{(d-k+1)}), a^{(d-k+1)} \in M^{k-1}\}.$$

The functions L_{\pm}^k can be extended to continuous functions defined on $\overline{M^{k-1}}$, whose values might coincide (but this does not necessarily happen) only on ∂M^{k-1} .

Remark 3. Theorem 2 can be deduced from Theorem 3 (but we give in Section 5 a direct proof which is short enough). Indeed, given a component M of the set Π_d^* , one can successively contract it into its projections M^{d-1} , M^{d-2} , ..., M^2 . The latter is one of the sets $\Pi_{d\pm\pm}^{*2}$ defined in Example 2 which are contractible.

In Section 2 we remind some results which are used in the proof of Theorem 2. In Section 3 we introduce some notation and we give examples concerning the sets Π_d and Π_d^* for $d = 1, 2$ and 3 . These examples are used in the proofs of Theorems 2 and 3. In Section 6 we make comments on Theorems 1, 2 and 3 and we formulate open problems.

2. KNOWN RESULTS ABOUT THE HYPERBOLICITY DOMAIN

Before proving Theorems 1, 2 and 3 we remind some results about the set Π_d which are due to V. I. Arnold, A. B. Givental and the author, see [3], [11] and [14] or Chapter 2 of [17] and the references therein.

Notation 3. We denote by K_d the simplicial angle $\{x_1 \geq x_2 \geq \dots \geq x_d\} \subset \mathbb{R}^d$ and by $\tilde{\mathcal{V}}$ the Viète mapping

$$\tilde{\mathcal{V}} : (x_1, \dots, x_d) \mapsto (\varphi_1, \dots, \varphi_d), \quad \varphi_j = \sum_{1 \leq i_1 < i_2 < \dots < i_j \leq d} x_{i_1} x_{i_2} \dots x_{i_j}.$$

Strata of K_d are denoted by their *multiplicity vectors*. E. g. for $d = 5$, the stratum of K_5 defined by the multiplicity vector $(2, 2, 1)$ is the set $\{x_1 = x_2 > x_3 = x_4 > x_5\} \subset \mathbb{R}^5$. The same notation is used for strata of Π_d which is justified by parts (3) and (4) of Theorem 4.

Remark 4. The set $\Delta_d \cap \Pi_d = \Delta_d^1 \cap \Pi_d$ consists of points $a \in \Pi_d \subset \mathbb{R}^d$, for which the hyperbolic polynomial Q_d has at least one root of multiplicity ≥ 2 . That is why $\Pi_d \setminus \Delta_d = \Pi_d \setminus \Delta_d^1 = S_{1^d}$ is the stratum of Π_d with multiplicity vector $1^d = (1, \dots, 1)$ and

$$\Pi_d^* = S_{1^d} \setminus (\{a_0 = 0\} \cup \dots \cup \{a_{d-1} = 0\}).$$

The strata of Π_d^* (they are all of dimension d , so they can also be called *components*) are of the form

$$S_{1^d}(\sigma) := \{a \in S_{1^d} \mid \text{sign}(a_j) = \sigma_j, 0 \leq j \leq d-1\}$$

for some $\sigma = (\sigma_0, \dots, \sigma_{d-1}) \in \{\pm\}^d$.

Theorem 4. (1) For $k \geq 3$, every non-empty fibre \tilde{f}_k of the projection $\pi^k : \Pi_d^k \rightarrow \Pi_d^{k-1}$ is either a segment or a point.

(2) The fibre \tilde{f}_k is a segment (resp. a point) exactly if the fibre is over a point of the interior of Π_d^{k-1} (resp. over $\partial\Pi_d^{k-1}$).

(3) The mapping $\tilde{\mathcal{V}} : K_d \rightarrow \Pi_d$ is a homeomorphism.

(4) The restriction of the mapping $\tilde{\mathcal{V}}$ to (the closure of) any stratum of K_d defines a homeomorphism of the (closure of the) stratum onto its image which is (the closure of) a stratum of Π_d .

(5) A stratum S of Π_d defined by a multiplicity vector with ℓ components is a smooth ℓ -dimensional real submanifold in \mathbb{R}^d . It is the graph of a smooth $(d - \ell)$ -dimensional vector-function defined on the projection of the stratum in $Oa_{d-\ell} \dots a_{d-1}$. Thus S is a real manifold with boundary. The field of tangent spaces to S continuously extends to the strata from the closure of S . The extension is everywhere transversal to the space $Oa_0 \dots a_{d-\ell-1}$. That is, the sum of the two

vector spaces $Oa_0 \dots a_{d-\ell-1}$ and (the extension of) the field of tangent spaces to S is the space $Oa_0 \dots a_{d-1}$.

(6) For $k \geq 3$, the set Π_d^k is the set of points on and between the graphs H_+^k and H_-^k of two locally Lipschitz functions defined on Π_d^{k-1} whose values coincide on and only on $\partial\Pi_d^{k-1}$:

$$\begin{aligned} \Pi_d^k &= \{(a_{d-k}, a^{(d-k+1)}) \in \mathbb{R} \times \Pi_d^{k-1} \mid H_-^k(a^{(d-k+1)}) \leq a_{d-k} \leq H_+^k(a^{(d-k+1)})\}, \\ &(H_-^k(a^{(d-k+1)}) = H_+^k(a^{(d-k+1)})) \Leftrightarrow (a^{(d-k+1)} \in \partial\Pi_d^{k-1}). \end{aligned}$$

(7) For $k \geq 3$, the graph H_+^k (resp. H_-^k) consists of the closures of the strata whose multiplicity vectors are of the form $(r, 1, s, 1, \dots)$ (resp. $(1, r, 1, s, \dots)$) and which have exactly $k-1$ components. (In [17] it is written “ k components” which is wrong.)

(8) For $2 \leq k \leq \ell$, the projection S^k of every ℓ -dimensional stratum S of Π_d in the space $Oa_{d-k} \dots a_{d-1}$ is the set of points on and between the graphs $H_+^k(S)$ and $H_-^k(S)$ of two locally Lipschitz functions defined on the closure S^{k-1} of S^{k-1} whose values coincide on and only on ∂S^{k-1} .

Remarks 3. (1) The projections π^k are defined also for $k=2$. For $k=2$, each fibre \tilde{f}_2 is a half-line and only the graph H_2^+ (but not H_2^-) is defined, see Example 2.

(2) Consider two strata S_1 and S_2 of Π_d defined by their multiplicity vectors $\mu(S_1)$ and $\mu(S_2)$. The stratum S_2 belongs to the topological and algebraic closure of the stratum S_1 if and only if the vector $\mu(S_2)$ is obtained from the vector $\mu(S_1)$ by finitely-many replacings of two consecutive components by their sum.

Remark 5. For $m \geq 2$, consider the fibres f_m° of the projection

$$\pi_*^m : \Pi_d \rightarrow \Pi_d^m, \quad \pi_*^m := \pi^{m+1} \circ \dots \circ \pi^d.$$

In particular, $\tilde{f}_d = f_{d-1}^\circ$. Suppose that such a fibre f_m° is over a point $A := (a_{d-m}^0, \dots, a_{d-1}^0) \in \Pi_d^m$. When non-empty, the fibre f_m° is either a point (when $A \in \partial\Pi_d^m$) or a set homeomorphic to a $(d-m)$ -dimensional cell and its boundary (when $A \in \Pi_d^m \setminus \partial\Pi_d^m$). This follows from part (6) of Theorem 4. The boundary of the cell can be represented as consisting of:

- two 0-dimensional cells (these are the graphs of the functions $H_\pm^{m+1}|_A$),
- two 1-dimensional cells (the graphs of $H_\pm^{m+2}|_{(\pi^{m+1})^{-1}(A)}$),
- two 2-dimensional cells (the graphs of $H_\pm^{m+3}|_{(\pi^{m+1} \circ \pi^{m+2})^{-1}(A)}$),
- ...,
- two $(d-m-1)$ -dimensional cells (the graphs of $H_\pm^d|_{(\pi^{m+1} \circ \pi^{m+2} \circ \dots \circ \pi^{d-1})^{-1}(A)}$).

Remark 6. It is a priori clear that for the functions L_{\pm}^k defined in Theorem 3, one has the inequalities

$$L_{+}^k(a^{(d-k+1)}) \leq H_{+}^k(a^{(d-k+1)}) \quad \text{and} \quad L_{-}^k(a^{(d-k+1)}) \geq H_{-}^k(a^{(d-k+1)})$$

for each value of $a^{(d-k+1)}$, where L_{+}^k or L_{-}^k (hence H_{+}^k or H_{-}^k) is defined. It is also clear that the border of each component of the set Π_d^* consists of parts of the closures of the graphs H_{\pm}^d and of parts of the hyperplanes $\{a_j = 0\}$, $j = 1, \dots, d-1$.

In Chapter 2 of [17] one can find also results concerning the hyperbolicity domain which are exposed in the thesis [21] of I. Méguerditchian.

3. NOTATION AND EXAMPLES

Notation 4. Given a d -tuple $\sigma = (\sigma_0, \dots, \sigma_{d-1})$, where $\sigma_j = +$ or $-$, we denote by $\mathcal{R}(\sigma)$ the subset of $\mathbb{R}^d \cong \mathcal{O}a_0 \cdots a_{d-1}$ defined by the conditions $\text{sign}(a_j) = \sigma_j$, $j = 0, \dots, d-1$, and we set $\Pi_{d,\sigma}^* := \Pi_d^* \cap \mathcal{R}(\sigma)$. For a set $T \subset \mathcal{O}a_0 \cdots a_{d-1}$, we denote by T^k its projection in the space $\mathcal{O}a_{d-k} \cdots a_{d-1}$.

Example 1. For $k = 1$ and for $a_j = 0$, $j = 0, \dots, d-2$, there exists a hyperbolic polynomial of the form $(x + a_{d-1})x^{d-1}$ with any $a_{d-1} \in \mathbb{R}$, so $\Pi_d^1 = \mathbb{R}$. If one chooses any hyperbolic degree d polynomial Q_d^* with distinct roots, the shift $x \mapsto x + g$ results in $a_{d-1} \mapsto a_{d-1} + dg$, so there exist such polynomials Q_d^* with any values of a_{d-1} . In addition, one can perturb the coefficients a_0, \dots, a_{d-2} to make them all non-zero by keeping the roots real and distinct. Thus $\Pi_d^{*1} = \mathbb{R}^* = \mathbb{R} \setminus \{a_{d-1} = 0\}$,

$$\Pi_d^{*1} \cap \{a_{d-1} > 0\} = \{\mathbb{R}_+^* : a_{d-1} > 0\}, \quad \Pi_d^{*1} \cap \{a_{d-1} < 0\} = \{\mathbb{R}_-^* : a_{d-1} < 0\}.$$

Example 2. One can formulate analogs to parts (1), (6) and (7) of Theorem 4 for $k = 2$ by saying that the border of the set Π_d^2 is the set H_{+}^2 while H_{-}^2 is empty, see part (1) of Remarks 3.

The set H_{+}^2 is the projection in $\mathbb{R}^2 \cong \mathcal{O}a_{d-2}a_{d-1}$ of the stratum of Π_d consisting of polynomials having a d -fold real root: $(x + \lambda)^d$. Its multiplicity vector equals (d) . Hence $a_{d-1} = d\lambda$, $a_{d-2} = d(d-1)\lambda^2/2$, so $H_{+}^2 : a_{d-2} = (d-1)a_{d-1}^2/2d$. One can observe that

$$\Pi_d^{*2} = \{a_{d-2} \neq 0 \neq a_{d-1}, a_{d-2} < (d-1)a_{d-1}^2/2d\},$$

$$\Pi_d^{*2} \cap \{a_{d-1} > 0, a_{d-2} > 0\} = \{a_{d-1} > 0, 0 < a_{d-2} < (d-1)a_{d-1}^2/2d\} =: \Pi_{d++}^{*2},$$

$$\Pi_d^{*2} \cap \{a_{d-1} < 0, a_{d-2} > 0\} = \{a_{d-1} < 0, 0 < a_{d-2} < (d-1)a_{d-1}^2/2d\} =: \Pi_{d-+}^{*2},$$

$$\Pi_d^{*2} \cap \{a_{d-1} > 0, a_{d-2} < 0\} = \{a_{d-1} > 0, a_{d-2} < 0\} =: \Pi_{d+-}^{*2} \quad \text{and}$$

$$\Pi_d^{*2} \cap \{a_{d-1} < 0, a_{d-2} < 0\} = \{a_{d-1} < 0, a_{d-2} < 0\} =: \Pi_{d--}^{*2}.$$

To obtain similar formulas for Π_d^2 instead of Π_d^{*2} one has to replace everywhere the inequalities $a_{d-1} < 0$, $a_{d-1} > 0$, $a_{d-2} < 0$, $a_{d-2} > 0$ and $a_{d-2} < (d-1)a_{d-1}^2/2d$ by $a_{d-1} \leq 0$, $a_{d-1} \geq 0$, $a_{d-2} \leq 0$, $a_{d-2} \geq 0$ and $a_{d-2} \leq (d-1)a_{d-1}^2/2d$ respectively.

Example 3. For $d = 3$ (hence $\sigma = (\sigma_0, \sigma_1, \sigma_2)$), we set $a_2 := a$, $a_1 := b$, $a_0 := c$, and we consider the polynomial $Q_3 := x^3 + ax^2 + bx + c$. Taking into account the group of quasi-homogeneous dilatations which preserves the discriminant set (see part (2) of Remarks 1) one concludes that each set $\Pi_{3,\sigma}^*$ is diffeomorphic to the corresponding direct product

$$(\Pi_{3,\sigma}^* \cap \{a = 1\}) \times (0, \infty) \text{ if } \sigma_2 = + \text{ or } (\Pi_{3,\sigma}^* \cap \{a = -1\}) \times (-\infty, 0) \text{ if } \sigma_2 = -.$$

Set $\sigma' := (-\sigma_0, \sigma_1, -\sigma_2)$. Using the same group of dilatations with $u = -1$ one deduces that the set $\Pi_{3,\sigma'}^* \cap \{a = -1\}$ is diffeomorphic to the set $\Pi_{3,\sigma}^* \cap \{a = 1\}$. Therefore in order to prove that all sets $\Pi_{3,\sigma}^*$ are contractible it suffices to show this for the sets $\Pi_{3,\sigma}^* \cap \{a = 1\}$ with $\sigma_2 = +$. The latter sets are shown in Figure 1.

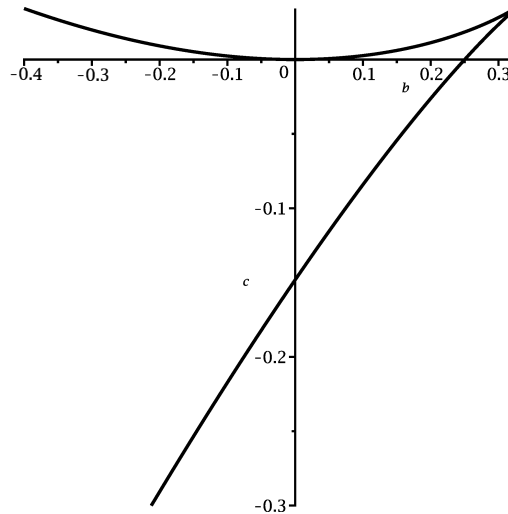


Figure 1: The discriminant set of the family of polynomials $x^3 + x^2 + bx + c$ and the sets $\Pi_{3,\sigma}^* \cap \{a = 1\}$.

The figure represents the discriminant set of the polynomial $Q_3^\bullet := x^3 + x^2 + bx + c$, i. e. the set

$$\text{Res}(Q_3^\bullet, Q_3^{\bullet'}, x) = 4b^3 - b^2 - 18bc + 27c^2 + 4c = 0.$$

(The set Δ_3^2 is empty, because there is not more than one complex conjugate pair of roots, so $\Delta_3 = \Delta_3^1$, see Remarks 1.) This is a curve in $\mathbb{R}^2 := Obc$ having a cusp

at $(b, c) = (1/3, 1/27)$ which corresponds to the polynomial $(x + 1/3)^3$. The four sets $\Pi_{3,\sigma}^* \cap \{a = 1\}$ are the intersections of the interior of the curve with the open coordinate quadrants. The intersections with $\{b > 0, c > 0\}$ and $\{b > 0, c < 0\}$ are bounded curvilinear triangles.

4. PROOF OF THEOREM 1

Part (1). Each set $\mathcal{E}_d(\sigma)$ is non-empty. Indeed, given a polynomial Q_d with $a_0 > 0$ (see Remark 2), for $C > 0$ large enough, the polynomial $Q_d + C$ is elliptic. If the polynomials $Q_{d,1}$ and $Q_{d,2}$ belong to the set $\mathcal{E}_d(\sigma)$, then for $t \in [0, 1]$, the polynomial $Q_d^\sharp := tQ_{d,1} + (1-t)Q_{d,2}$ also belongs to it. Indeed, the signs of the respective coefficients are the same and if $Q_{d,1}(x) > 0$ and $Q_{d,2}(x) > 0$, then $Q_d^\sharp(x) > 0$. Thus the set $\mathcal{E}_d(\sigma)$ is convex hence contractible.

Part (2). Each set $\mathcal{F}_d(\sigma)$ is non-empty. Indeed, for $C > 0$ large enough, the polynomial $Q_d + \text{sign}(a_0)C$ has a single real root which is simple and the sign of this root is opposite to the sign of $Q_d(0)$. For a given polynomial $Q_d \in \mathcal{F}_d(\sigma)$, denote this root by ξ . Hence the polynomial $Q_d^0 := |\xi|^d Q_d(x/|\xi|)$ is in $\mathcal{F}_d(\sigma)$ and has a root at 1 or -1 . Suppose that the root is at 1 (for -1 the proof is similar). We show that the subset $\mathcal{F}_d^0(\sigma)$ of $\mathcal{F}_d(\sigma)$ consisting of such polynomials Q_d^0 is convex hence contractible. On the other hand the set $\mathcal{F}_d(\sigma)$ is diffeomorphic to $\mathcal{F}_d^0(\sigma) \times \mathbb{R}_+^*$ from which contractibility of $\mathcal{F}_d(\sigma)$ follows.

For any two polynomials $Q_d^{0,\dagger}, Q_d^{0,*} \in \mathcal{F}_d^0(\sigma)$, the signs of the coefficients of the polynomial

$$Q_d^{0,b} := tQ_d^{0,\dagger} + (1-t)Q_d^{0,*}, \quad t \in [0, 1],$$

are the same as the signs of the respective coefficients of $Q_d^{0,\dagger}$ and $Q_d^{0,*}$, hence $Q_d^{0,b} \in \mathcal{F}_d^0(\sigma)$. This proves that $\mathcal{F}_d^0(\sigma)$ is convex.

Part (3).

A) *Contractibility of the sets $\mathcal{G}_{d,(2,0)}(\sigma)$ and $\mathcal{G}_{d,(0,2)}(\sigma)$.*

The two real roots of Q_d have the same sign (i. e. $a_0 > 0$). We assume that they are positive, i. e. we prove contractibility only of $\mathcal{G}_{d,(2,0)}(\sigma)$; otherwise one can consider the polynomial $Q_d(-x)$ with the d -tuple $\tilde{\sigma}$ resulting from σ via $x \mapsto -x$ (this mapping induces a bijection of the set of d -tuples onto itself) and contractibility of $\mathcal{G}_{d,(0,2)}(\tilde{\sigma})$ will be proved in the same way. Denote the real roots of Q_d by $0 < \xi < \eta$.

We can assume that at least one coefficient of odd degree of Q_d is negative. Indeed, if all coefficients of Q_d^0 of odd degree are positive, then by Descartes' rule of signs the polynomial Q_d^0 can have two real positive roots only if there is at least one coefficient of even degree which is negative. However in this case (and this is Case 1)) the set $\mathcal{G}_{d,(2,0)}(\sigma)$ is empty, see Proposition 4 in [7].

Next, we assume that $\eta = 1$ (hence $\xi \in (0, 1)$). Indeed, if one considers instead of $Q_d \in \mathcal{G}_{d,(2,0)}(\sigma)$ the polynomial $Q_d^0 := \eta^d Q_d(x/\eta)$, one has $Q_d^0 \in \mathcal{G}_{d,(2,0)}(\sigma)$ and $Q_d^0(1) = 0$. We denote the set of such polynomials Q_d^0 by $\mathcal{G}_{d,(2,0)}^0(\sigma)$. As $\mathcal{G}_{d,(2,0)}(\sigma)$ is diffeomorphic to $\mathcal{G}_{d,(2,0)}^0(\sigma) \times \mathbb{R}_+^*$, contractibility of $\mathcal{G}_{d,(2,0)}^0(\sigma)$ implies the one of $\mathcal{G}_{d,(2,0)}(\sigma)$.

For $\xi^* \in (0, 1)$, we denote by $\mathcal{G}_{d,(2,0)}^{0,\xi^*}(\sigma)$ the subset of polynomials of $\mathcal{G}_{d,(2,0)}^0(\sigma)$ with $\xi = \xi^*$. If $Q_d^{0,1}$ and $Q_d^{0,2}$ are two polynomials of $\mathcal{G}_{d,(2,0)}^{0,\xi^*}(\sigma)$, then for $t \in [0, 1]$, one has $tQ_d^{0,1} + (1-t)Q_d^{0,2} \in \mathcal{G}_{d,(2,0)}^{0,\xi^*}(\sigma)$. Therefore for each $\xi \in (0, 1)$, the set $\mathcal{G}_{d,(2,0)}^{0,\xi}(\sigma)$ is convex hence contractible, and to prove contractibility of $\mathcal{G}_{d,(2,0)}^0(\sigma)$ (and hence of $\mathcal{G}_{d,(2,0)}(\sigma)$) it suffices to find for each $\xi \in (0, 1)$ a polynomial $Q_d^{0,\xi} \in \mathcal{G}_{d,(2,0)}^{0,\xi}(\sigma)$ depending continuously on ξ .

Suppose m is odd, $1 \leq m \leq d-1$, and that the coefficient of $Q_d \in \mathcal{G}_{d,(2,0)}(\sigma)$ of x^m must be negative. There exists a unique polynomial of the form

$$R := x^d - Ax^m + B, \quad A > 0, \quad B > 0, \quad \text{such that} \quad R(\xi) = R(1) = 0. \quad (4.1)$$

Indeed, the conditions

$$\xi^d - A\xi^m + B = 1 - A + B = 0 \quad (4.2)$$

imply

$$A = (1 - \xi^d)/(1 - \xi^m) > 0 \quad \text{and} \quad B = -1 + A = \xi^m(1 - \xi^{d-m})/(1 - \xi^m) > 0. \quad (4.3)$$

Remarks 4. (1) The fractions for A , B and B/ξ^m can be extended by continuity for $\xi = 0$ and $\xi = 1$. For $\xi \in [0, 1]$, one has

$$\begin{aligned} A &\in [1, \frac{d}{m}], & \lim_{\xi \rightarrow 0^+} A &= 1, & \lim_{\xi \rightarrow 1^-} A &= \frac{d}{m}, \\ B &\in [0, \frac{d-m}{m}], & \lim_{\xi \rightarrow 0^+} B &= 0^+, & \lim_{\xi \rightarrow 1^-} B &= \frac{d-m}{m}, \\ B/\xi^m &\in [0, \max(\frac{d-m}{m}, 1)], & \lim_{\xi \rightarrow 0^+} B/\xi^m &= 1, & \lim_{\xi \rightarrow 1^-} B/\xi^m &= \frac{d-m}{m}. \end{aligned} \quad (4.4)$$

(2) The function R has a global minimum at some point $x_M = x_M(\xi) \in (0, 1)$. One has

$$\lim_{\xi \rightarrow 0^+} x_M(\xi) = x_{M,0} = (m/d)^{1/(d-m)} \in (0, 1), \quad R(x_{M,0}) < 0 \quad \text{and} \quad \lim_{\xi \rightarrow 1^-} x_M(\xi) = 1.$$

For $m \geq 3$, the tangent line to the graph of R for $x = 0$ is horizontal and $(0, R(0))$ is an inflection point. There is also another inflection point $x_I = x_I(\xi) \in (0, x_M)$.

Set

$$\mathcal{I} := \{1, 2, \dots, m-1, m+1, m+2, \dots, d-1\}. \quad (4.5)$$

We construct a polynomial $\Psi := \sum_{j=0}^{d-1} \psi_j x^j$ with signs of the coefficients ψ_j , $j \in \mathcal{I}$, as defined by the d -tuple σ and satisfying the conditions

$$\Psi(\xi) = \Psi(1) = 0. \quad (4.6)$$

The latter conditions can be considered as a linear system with unknown variables ψ_0 and ψ_m . Its determinant equals $\xi^m - 1 \neq 0$, so for given ψ_j , $j \in \mathcal{I}$, these conditions define a unique couple (ψ_0, ψ_m) whose signs are not necessarily the ones defined by the d -tuple σ . So to construct Ψ it suffices to fix ψ_j for $j \in \mathcal{I}$.

For each $\xi \in (0, 1)$ fixed and for $\varepsilon > 0$ sufficiently small, one has $R + \varepsilon\Psi \in \mathcal{G}_{d,(2,0)}^{0,\xi}(\sigma)$. Indeed, for $m \neq j \neq 0$, the coefficients of $R + \varepsilon\Psi$ have the signs defined by the d -tuple σ , so one has to check two things:

1) If ε is small enough, then

$$-A + \varepsilon\psi_m < 0 \quad \text{and} \quad B + \varepsilon\psi_0 > 0. \quad (4.7)$$

To obtain these two conditions simultaneously for all $\xi \in (0, 1)$, one has to choose ε as a function of ξ .

The conditions (4.6) can be given the form

$$\xi^m \psi_m + \psi_0 = U, \quad \psi_m + \psi_0 = V,$$

where U and V are polynomials in ξ of degree $\leq d - 1$. Hence

$$\psi_0 = (U - \xi^m V)/(1 - \xi^m) \quad \text{and} \quad \psi_m = (V - U)/(1 - \xi^m). \quad (4.8)$$

Formulas (4.8) imply that Ψ is of the form

$$K(x, \xi)/(1 - \xi)^m, \quad K \in \mathbb{R}[x, \xi], \quad \deg_x K \leq d - 1. \quad (4.9)$$

As $\xi \rightarrow 0^+$, the quantity B decreases as ξ^m , see (4.3) and (4.4). As $\xi \rightarrow 1^-$, the quantities $|\psi_0|$ and $|\psi_m|$ increase not faster than $C/(1 - \xi)$ for some $C > 0$. So to obtain $\varepsilon = \varepsilon(\xi)$ such that conditions (4.7) hold for $\xi \in (0, 1)$, it suffices to set $\varepsilon := c\xi^{m+1}(1 - \xi)^3$ for some $c > 0$ small enough.

2) For $\xi \in (0, 1)$, one must have

$$R + \varepsilon\Psi > 0 \quad \text{for} \quad x \in (-\infty, \xi) \cup (1, \infty), \quad \text{and} \quad R + \varepsilon\Psi < 0 \quad \text{for} \quad x \in (\xi, 1). \quad (4.10)$$

Lemma 1. *It is possible to choose $c > 0$ so small that conditions (4.7) and (4.10) hold true simultaneously.*

The lemma implies that for such $c > 0$, $R + \varepsilon(\xi)\Psi \in \mathcal{G}_{d,(2,0)}^{0,\xi}(\sigma)$. So one can set $Q_d^{0,\xi} := R + \varepsilon(\xi)\Psi$ from which contractibility of $\mathcal{G}_{d,(2,0)}(\sigma)$ follows.

Proof of Lemma 1. Conditions (4.7) were already discussed, so we focus on conditions (4.10). Lowercase indices ξ indicate differentiations w.r.t. ξ .

a) To obtain the condition $R + \varepsilon\Psi > 0$ for $x > 1$, it suffices to get $(R + \varepsilon\Psi)' > 0$ for $x \geq 1$. For $x \geq 1$, one has

$$R' = dx^{d-1} - mA x^{m-1} = x^{m-1}(dx^{d-m} - mA) \geq x^{m-1}(d - mA) \quad (4.11)$$

(as $R'(1) > 0$, one knows that $d - mA > 0$). Next,

$$d - mA = \Lambda/(1 - \xi^m), \quad \Lambda := d - m + m\xi^d - d\xi^m.$$

There exists $\alpha > 0$ such that for $\xi \in [0, 1]$, $\Lambda \geq \alpha(1 - \xi)^2$. Indeed, $\Lambda_\xi = dm(\xi^{d-1} - \xi^{m-1}) \leq 0$, with equality only for $\xi = 0$ and $\xi = 1$, so Λ is strictly decreasing on $[0, 1]$. The existence of α follows from

$$\Lambda(0) = d - m > 0, \quad \Lambda(1) = \Lambda_\xi(1) = 0 \quad \text{and}$$

$$\Lambda_{\xi\xi} = dm((d-1)\xi^{d-1} - (m-1)\xi^{m-1}) \quad \text{hence}$$

$$\Lambda_{\xi\xi}(1) = dm(d - m) > 0.$$

Thus for $\xi \in (0, 1)$ and $x > 1$, one has

$$R \geq (x^m/m)\alpha(1 - \xi)^2/(1 - \xi^m) \quad \text{and} \quad \varepsilon(\xi)\Psi \leq c\xi^{m+1}(1 - \xi)^3 K(x, \xi)/(1 - \xi^m),$$

see (4.9). One can choose $c > 0$ sufficiently small so that for $x \in (1, 2]$, $R + \varepsilon\Psi > 0$. There exists $\beta > 0$ such that for $x \geq 2$, $dx^{d-m} - mA > \beta x^{d-m}$ (see (4.11) and (4.4)), so $R \geq \beta x^d/d$ and for $c > 0$ small enough, $R + \varepsilon\Psi > 0$.

b) For $x \leq -1$ (resp. for $x \in [-1, 0]$), one has

$$R \geq |x|^m(|x|^{d-m} + A) \quad (\text{resp.} \quad R \geq B \geq (\max((d-m)/m, 1))\xi^m)$$

(see (4.1) and (4.4)) which for $c > 0$ small enough is larger than $|\varepsilon(\xi)\Psi|$ and (4.10) holds true.

c) Suppose that $x \in (0, \xi)$. Then $R \geq \min(h(x, \xi), q(x, \xi))$, where

$$\tau : y = h(x, \xi) := R'(\xi)(x - \xi) \quad \text{and} \quad \chi : y = q(x, \xi) := B - Bx/\xi$$

are the tangent line to the graph of R at the point $(\xi, 0)$ and the line joining the points $(0, B)$ and $(\xi, 0)$ respectively. Indeed, if $x_I \in [\xi, 1]$ (see part (2) of Remarks 4), then the graph of R is concave for $x \in [0, \xi]$, so it is situated above the line χ . If $x_I \in (0, \xi)$, then for $x \in [x_I, \xi]$, one has $R \geq h(x, \xi)$ and for $x \in (0, x_I]$, one has $R \geq q_1(x, \xi)$, where

$$\chi_1 : y = q_1(x, \xi) := R(x_I) + (x - x_I)(R(x_I) - B)/x_I$$

is the line joining the points $(0, B)$ and $(x_I, R(x_I))$. The line χ_1 is above the line χ for $x \in (0, x_I)$.

Consider the smaller in absolute value of the slopes of the lines τ and χ , i.e. $\mu := \min(|R'(\xi)|, B/\xi)$. One finds that

$$R'(\xi) = \xi^{m-1}g(\xi)/(1 - \xi^m), \quad g := d\xi^{d-m} - m - (d-m)\xi^d,$$

with $g_\xi = d(d-m)(\xi^{d-m-1} - \xi^{d-1}) \geq 0$, with equality only for $\xi = 0$ and $\xi = 1$. As $g(0) = -m < 0$, $g(1) = 0$,

$$g_{\xi\xi} = d(d-m)((d-m-1)\xi^{d-m-2} - (d-1)\xi^{d-2}), \quad \text{so } g_{\xi\xi}(1) = -md(d-m) < 0,$$

there exists $\tilde{\beta} > 0$ such that for $\xi \in (0, 1)$, $|R'(\xi)| \geq \tilde{\beta}\xi^{m-1}(1-\xi)^2/(1-\xi^m)$. On the other hand $B/\xi = \xi^{m-1}(1-\xi^{d-m})/(1-\xi^m)$. Thus $\mu \geq \mu_0 := \gamma\xi^{m-1}(1-\xi)$ for some $\gamma > 0$. Hence for $x \in (0, \xi)$, the graph of R is above the line $\delta : y = -\mu_0(x - \xi)$.

There exists $D_0 > 0$ such that for $\xi \in (0, 1)$ and $x \in [0, 1]$, one has $|(1-\xi)\Psi'| \leq D_0$, see (4.8). Hence if $c > 0$ is sufficiently small, the graph of $\varepsilon\Psi$ is below the line δ for $x \in [0, \xi)$, so $R + \varepsilon\Psi > 0$.

d) Suppose that $m \geq 3$ and that $\xi > 0$ is close to 0. Then for $x > \xi$, the line $\tilde{\tau}$, which is tangent to the graph of R at the point $(\xi, 0)$, is above the straight line $\tilde{\rho}$ joining the points $(\xi, 0)$ and $(x_M, R(x_M))$. Indeed,

$$R'(\xi) = \xi^{m-1}(d\xi^{d-m} - m - (d-m)\xi^d)/(1 - \xi^m) = O(\xi^{m-1})$$

whereas the slope of $\tilde{\rho}$ is close to $R(x_{M,0})/x_{M,0} < 0$. Therefore for $x \in (\xi, x_M]$, the graph of R is below the line $\tilde{\tau}$.

For $x \in [x_M, 1)$, the graph of R is below the line $\tilde{\chi}$ joining the points $(x_M, R(x_M))$ and $(1, 0)$ whose slope $-R(x_M)/(1-x_M)$ is close to $-R(x_{M,0})/(1-x_{M,0}) > 0$. On the other hand one has $|(1-\xi)\Psi'| \leq D_0$ (see c)), so $|\varepsilon(\xi)\Psi'| \leq c\xi^{m+1}(1-\xi)^2D_0$. Thus the graph of $\varepsilon(\xi)\Psi$ is above the line $\tilde{\tau}$ for $x \in (\xi, x_M]$ and above $\tilde{\chi}$ for $x \in [x_M, 1)$, hence it is between the graph of R and the x -axis for $x \in (\xi, 1)$, so $R + \varepsilon\Psi < 0$.

e) For $m \geq 3$, we fix $\theta_0 > 0$ small enough such that for $\xi \in (0, \theta_0]$, $R + \varepsilon\Psi < 0$, see d). For $m \geq 3$, $\xi \in [\theta_0, 1]$, $x \in (\xi, 1)$, and for $m = 1$, $\xi \in [0, 1]$, $x \in (\xi, 1)$, one has $R + \varepsilon(\xi)\Psi < 0$ if $c > 0$ is small enough. Indeed, one can write

$$R = (x-1)(x-\xi)R_1 \quad \text{and} \quad \Psi = (x-1)(x-\xi)\Psi_1, \quad R_1, \Psi_1 \in \mathbb{R}[x, \xi].$$

Then $R_1(x, \xi) > 0$. In particular, for $\xi = 1$, one obtains

$$R = x^d - (d/m)x^m + (d-m)/m, \quad R' = dx^{d-1} - dx^{m-1}, \quad \text{so } R'(1) = 0,$$

and $R'' = d((d-1)x^{d-2} - (m-1)x^{m-2})$ hence $R''(1) = d(d-m) > 0$, i. e. R is divisible by $(x-1)^2$, but not by $(x-1)^3$.

For $m = 1$, $\xi = 0$, one has $R'(0) < 0$ (whereas for $m = 3$, $\xi = 0$, one has $R'(0) = 0$), this why for $m = 1$ our reasoning is valid for $\xi \in [0, 1]$, not only for $\xi \in [\theta_0, 1]$.

Denote by $R_{1,0} > 0$ the minimal value of R_1 and by $\Psi_{1,0}$ the maximal value of Ψ_1 for $x \in [0, 1]$. One can choose $c > 0$ so small that for $x \in (\xi, 1)$ and for the values of ξ mentioned at the beginning of e),

$$R_1 - \varepsilon\Psi_1 \geq R_{1,0} - \varepsilon\Psi_{1,0} > 0, \quad \text{so} \quad R + \varepsilon\Psi < 0, \quad \text{because} \quad (x-1)(x-\xi) < 0.$$

The proof of the lemma results from a) – e). □

B) *Contractibility of the set $\mathcal{G}_{d,(1,1)}(\sigma)$.*

The two real roots of Q_d have opposite signs (hence $a_0 < 0$). Denote them by $-\eta < 0 < \xi$. We define the sets

$$\mathcal{K} := \mathcal{G}_{d,(1,1)}(\sigma) \cap \{\xi > \eta\}, \quad \mathcal{L} := \mathcal{G}_{d,(1,1)}(\sigma) \cap \{\xi < \eta\} \quad \text{and} \quad \mathcal{M} := \mathcal{G}_{d,(1,1)}(\sigma) \cap \{\xi = \eta\}.$$

Lemma 2. *Set $\sigma := (\sigma_0, \dots, \sigma_{d-1})$, $\sigma_j = +$ or $-$.*

(1) *Suppose that $\sigma_{2j+1} = +$, $j = 0, 1, \dots, (d/2) - 1$. Then $\mathcal{K} = \mathcal{M} = \emptyset$.*

(2) *Suppose that $\sigma_{2j+1} = -$, $j = 0, 1, \dots, (d/2) - 1$. Then $\mathcal{L} = \mathcal{M} = \emptyset$.*

(3) *Suppose that there exist two odd integers $j_1 \neq j_2$, $1 \leq j_1, j_2 \leq d-1$, such that $\sigma_{j_1} = -\sigma_{j_2}$. Then all three sets \mathcal{K} , \mathcal{L} and \mathcal{M} are non-empty. There exists an open d -dimensional ball $\mathcal{B} \subset \mathcal{G}_{d,(1,1)}(\sigma)$ centered at a point in \mathcal{M} and such that $\mathcal{B} \cap \mathcal{K} \neq \emptyset$ and $\mathcal{B} \cap \mathcal{L} \neq \emptyset$.*

Proof. Parts (1) and (2). If $\sigma_{2j+1} = +$ (resp. $\sigma_{2j+1} = -$), $j = 0, 1, \dots, (d/2) - 1$, then for a polynomial $Q_d \in \mathcal{G}_{d,(1,1)}(\sigma)$, one has $Q_d(0) < 0$ and $Q_d(a) > Q_d(-a)$ (resp. $Q_d(0) < 0$ and $Q_d(a) < Q_d(-a)$) for $a > 0$. Hence $\xi < \eta$ (resp. $\xi > \eta$).

Part (3). We construct a polynomial $Q_d^\circ \in \mathcal{M}$. Set $u := \xi^{j_1 - j_2}$ and

$$Q_d^\circ := x^d - \xi^d + \sigma_{j_1}(x^{j_1} - ux^{j_2}) + \varepsilon(Q_d^{\circ,o} + Q_d^{\circ,e}),$$

where

$$Q_d^{\circ,e} = b + \sum_{j=1}^{d/2} \sigma_{2j} x^{2j}, \quad b \in \mathbb{R}, \quad Q_d^{\circ,o} = rx^{j_1} + \sum_{j=0}^{d/2-1} \sigma_{2j+1} x^{2j+1}$$

and $\varepsilon > 0$ is small enough. We choose b and r such that $Q_d^{\circ,e}(\pm\xi) = 0$ and $Q_d^{\circ,o}(\pm\xi) = 0$ respectively. Then $Q_d^\circ(\pm\xi) = 0$ and for $j \neq 0$ and $j_1 \neq j \neq j_2$, the sign of the coefficients of x^j of Q_d° is as defined by σ . For $\varepsilon > 0$ small enough, one has $\text{sign}(Q_d^\circ(0)) = \text{sign}(-\xi^d + \varepsilon b) = -$. The coefficient of x^{j_1} (resp. x^{j_2}) of Q_d° equals $\sigma_{j_1} \times (1 + \varepsilon(1+r))$ (resp. $\sigma_{j_2} \times (u + \varepsilon(1+r))$), so it has the same sign as σ_{j_1} (resp. as σ_{j_2}).

Consider a d -dimensional ball \mathcal{B} centered at a point $Q_d^\circ \in \mathcal{M}$, with $\xi = \eta = \xi_0$ and belonging to $\mathcal{G}_{d,(1,1)}(\sigma)$. Perturb the real root ξ of Q_d° so that it takes values smaller and values larger than ξ_0 . The signs of the coefficients of Q_d° do not change. Hence \mathcal{B} intersects \mathcal{K} and \mathcal{L} . \square

We show first that each of the two sets \mathcal{K} and \mathcal{L} , when nonempty, is contractible. If we are in the conditions of part (1) or (2) of Lemma 2, then this implies contractibility of $\mathcal{G}_{d,(1,1)}(\sigma)$. When we are in the conditions of part (3), then one can contract \mathcal{K} and \mathcal{L} into points of \mathcal{B} and then contract \mathcal{B} into a point, so in this case $\mathcal{G}_{d,(1,1)}(\sigma)$ is also contractible.

We prove contractibility only of \mathcal{K} (when non-empty). The one of \mathcal{L} is performed by complete analogy (the change of variable $x \mapsto -x$ exchanges the roles of \mathcal{K} and \mathcal{L} and changes the d -tuple σ accordingly). So we suppose that $\xi > \eta$. As in the proof of A) we reduce the proof of the contractibility of \mathcal{K} to the one of the contractibility of $\mathcal{K} \cap \{\xi = 1\}$. As in A) we observe that if

$$Q_d^\ddagger, \quad Q_d^\Delta \in \mathcal{K}^{\eta^*} := \mathcal{K} \cap \{\xi = 1, \eta = \eta^* \in (0, 1)\},$$

then $tQ_d^\ddagger + (1-t)Q_d^\Delta \in \mathcal{K}^{\eta^*}$, so \mathcal{K}^{η^*} is convex hence contractible and contractibility of $\mathcal{K} \cap \{\xi = 1\}$ (and also of \mathcal{K}) will be proved if we construct for each $\eta \in (0, 1)$ a polynomial $Q_d \in \mathcal{K}^\eta$ depending continuously on η .

Suppose that there is a negative coefficient of Q_d of odd degree m (otherwise \mathcal{K} is empty). For $\eta \in (0, 1)$, we construct a polynomial

$$S := x^d - \tilde{A}x^m - \tilde{B}, \quad \tilde{A} > 0, \quad \tilde{B} > 0, \quad \text{such that} \quad S(1) = S(-\eta) = 0.$$

The latter two equalities imply

$$\tilde{A} = (1 - \eta^d)/(1 + \eta^m) > 0 \quad \text{and} \quad \tilde{B} = \eta^m(1 + \eta^{d-m})/(1 + \eta^m) > 0. \quad (4.12)$$

Remarks 5. (1) Thus for $\eta \in [0, 1]$, there exist constants $0 < B_{\min} \leq B_{\max}$ such that $\tilde{B}/\eta^m \in [B_{\min}, B_{\max}]$. Moreover one has

$$\begin{aligned} \tilde{A} &\in [0, 1], & \lim_{\eta \rightarrow 0^+} \tilde{A} &= 1, & \lim_{\eta \rightarrow 1^-} \tilde{A} &= 0^+, \\ \tilde{B} &\in [0, 1], & \lim_{\eta \rightarrow 0^+} \tilde{B} &= 0^+, & \lim_{\eta \rightarrow 1^-} \tilde{B} &= 1, \\ \lim_{\eta \rightarrow 0^+} \tilde{B}/\eta^m &= 1 & \text{and} & & \lim_{\eta \rightarrow 1^-} \tilde{B}/\eta^m &= 1. \end{aligned} \quad (4.13)$$

(2) The derivative S' has a unique root \tilde{x}_M (which is simple) in $(0, 1)$. All non-constant derivatives of S are increasing for $x > \tilde{x}_M$, have one or two roots (depending on m) in $[0, \tilde{x}_M)$ and no root outside this interval.

We construct a polynomial $\Phi := \sum_{j=0}^{d-1} \varphi_j x^j$, where for $j \in \mathcal{I}$ (see (4.5)), the sign of φ_j is defined by the d -tuple σ . This polynomial must satisfy the condition

$$\Phi(-\eta) = \Phi(1) = 0$$

which can be regarded as a linear system with known quantities φ_j , $j \in \mathcal{I}$, and with unknown variables φ_0 and φ_m :

$$-\eta^m \varphi_m + \varphi_0 = W, \quad \varphi_m + \varphi_0 = T, \quad W, T \in \mathbb{R}[\eta], \quad \text{so} \quad (4.14)$$

$$\varphi_0 = (\eta^m T + W)/(1 + \eta^m), \quad \varphi_m = (T - W)/(1 + \eta^m).$$

One must also have $S + \varepsilon_1(\eta)\Phi \in \mathcal{K}^\eta$, $\eta \in (0, 1)$, for some suitably chosen positive-valued continuous function $\varepsilon_1(\eta)$. For $\varepsilon_1(\eta) > 0$ small enough, the sign of the coefficient of x^j , $j \in \mathcal{I}$, of the polynomial $S + \varepsilon_1(\eta)\Phi$ is as defined by the d -tuple σ . So one needs to choose $\varepsilon_1(\eta)$ such that

$$-\tilde{A} + \varepsilon_1(\eta)\varphi_m < 0, \quad -\tilde{B} + \varepsilon_1(\eta)\varphi_0 < 0 \quad (4.15)$$

and

$$S + \varepsilon_1(\eta)\Phi > 0 \text{ for } x \in (-\infty, -\eta) \cup (1, \infty), \quad S + \varepsilon_1(\eta)\Phi < 0 \text{ for } x \in (-\eta, 1). \quad (4.16)$$

We set $\varepsilon_1 := \tilde{c}\eta^m(1 - \eta)^2$, $\tilde{c} > 0$. If one chooses \tilde{c} small enough, conditions (4.15) will hold true.

Lemma 3. *For $\tilde{c} > 0$ small enough, conditions (4.16) hold true.*

Contractibility of \mathcal{K} follows from the lemma.

Proof of Lemma 3. All derivatives of S of order $\leq d - 1$ are increasing functions in x for $x \geq 1$ (see Remarks 5). As

$$S'(1) = (d + d\eta^m - m + m\eta^d)/(1 + \eta^m) \geq (d - m)/2,$$

one can choose \tilde{c} small enough so that for $x \in [1, 2]$, $S' + \varepsilon_1(\eta)\Phi' > 0$. Hence $S + \varepsilon_1(\eta)\Phi > 0$ for $x \in (1, 2]$. If $x \geq 2$, then for some positive constants k_1 and k_2 , one has $S' \geq k_1 x^{d-1}$ and $\Phi' \leq k_2 x^{d-2}$, so if $\tilde{c} > 0$ is small enough, then for $x \geq 2$ (hence for $x > 1$), $S' + \varepsilon_1(\eta)\Phi' > 0$ and $S + \varepsilon_1(\eta)\Phi > 0$.

One has

$$S'(-\eta) = -(d\eta^{d-1} + (d - m)\eta^{d+m-1} + m\eta^{m-1})/(1 + \eta^m) = O(\eta^{m-1}),$$

$S'(-\eta) < 0$ and S is convex for $x < 0$. Hence one can choose $\tilde{c} > 0$ so small that for $x \in [-2, -\eta]$, $S' + \varepsilon_1(\eta)\Phi' < 0$ hence $S + \varepsilon_1(\eta)\Phi > 0$. Indeed, for $\eta \in [0, 1]$ and $x \in [-2, 0]$, Φ' is bounded. For $x \leq -2$, one has

$$S' \leq k_1^* x^{d-1} \quad \text{and} \quad |\Phi'| \leq k_2^* x^{d-2}$$

for some positive constants k_1^* , k_2^* , so $S + \varepsilon_1(\eta)\Phi < 0$ (thus this holds true for $x < -\eta$).

The function S is convex on $[-\eta, 0]$, see Remarks 5. Hence for $x \in [-\eta, 0]$, the graph of S is below the line ζ joining the points $(-\eta, 0)$ and $(0, -\tilde{B})$. Its slope is $-\tilde{B}/\eta$, with $|\tilde{B}/\eta| = O(\eta^{m-1})$. Hence for $x \in [-\eta, 0]$ and for $\tilde{c} > 0$ sufficiently small, the graph of Φ is above the line ζ (because Φ' is bounded for $x \in [-1, 0]$, $\eta \in [0, 1]$) and one has $S + \varepsilon_1(\eta)\Phi < 0$.

Suppose that $x \in [0, \tilde{x}_M]$. The function S is decreasing, see Remarks 5, hence $S(x) \leq S(0) = -\tilde{B} = O(\eta^m)$. As there exists $k_3 > 0$ such that for $x \in [0, 1]$, $|\Phi| \leq k_3$, for $\tilde{c} > 0$ sufficiently small, one has $S + \varepsilon_1(\eta)\Phi < 0$.

For $x \in [\tilde{x}_M, 1]$, the function S is convex, hence its graph is below the line $\tilde{\zeta}$ joining the points $(\tilde{x}_M, S(\tilde{x}_M))$ and $(1, 0)$. Recall that $S(\tilde{x}_M) \leq S(0) = -\tilde{B} = O(\eta^m)$. There exists $k_4 > 0$ such that for $x \in [0, 1]$ and $\eta \in [0, 1]$, $|\Phi'| \leq k_4$. Thus the slope of $\tilde{\zeta}$ is

$$\geq \tilde{B}/(1 - \tilde{x}_M) > \tilde{B} = O(\eta^m)$$

while $|\varepsilon\Phi'| \leq \tilde{c}\eta^m(1 - \eta)^2 k_4$. Hence for sufficiently small values of $\tilde{c} > 0$, the graph of $\varepsilon\Phi$ is above the line $\tilde{\zeta}$ and $S + \varepsilon_1(\eta)\Phi < 0$. \square

5. PROOFS OF THEOREMS 2 AND 3

Proof of Theorem 2. In the proof we assume that the polynomials of Π_d are of the form $Q_d := x^d + a_{d-1}x^{d-1} + \dots + a_2x^2 + a_1x + a_0$ and the ones of Π_{d-1} are of the form $Q_{d-1} := x^{d-1} + a_{d-1}x^{d-2} + \dots + a_2x + a_1$. Thus the intersection $\Pi_d \cap \{a_0 = 0\}$ can be identified with Π_{d-1} .

We show that every polynomial $Q_d \in \Pi_d^*$ can be continuously deformed so that it remains in Π_d^* , the signs of its coefficients do not change throughout the deformation except the one of a_0 which vanishes at the end of the deformation. Therefore

1) throughout the deformation the quantities of positive and negative roots do not change;

2) at the end of the deformation exactly one root vanishes and a polynomial of the form xQ_{d-1} is obtained with $Q_{d-1} \in \Pi_d \cap \{a_0 = 0\}$.

Moreover, we show that throughout and at the end of the deformation one obtains polynomials with distinct real roots. Thus any given component of the set Π_d^* can be retracted into a component of the set Π_{d-1}^* ; the latter is defined by the $(d - 1)$ -tuple obtained from σ by deleting its first component. For $d = 2$, all components of the set Π_2^* are contractible, see Example 2.

This means that for every given d and σ , there exists exactly one component of Π_d^* , and which is contractible. The deformation mentioned above is defined like this:

$$Y_d := (Q_d + txQ'_d)/(1 + td) = \sum_{j=0}^d ((1 + tj)/(1 + td))a_jx^j, \quad t \geq 0.$$

It is clear that the polynomial Y_d is monic, with $\text{sign}(a_j) = \text{sign}((1 + tj)a_j/(1 + td))$ and $\lim_{t \rightarrow +\infty} ((1 + tj)a_j/(1 + td)) = ja_j/d$. There remains to prove only that Y_d has d distinct real roots.

Denote the roots of Q_d by $\eta_1 < \dots < \eta_s < 0 < \xi_1 < \dots < \xi_{d-s}$. The polynomial Q'_d has exactly one root in each of the intervals $(\eta_1, \eta_2), \dots, (\eta_{s-1}, \eta_s), (\eta_s, \xi_1), (\xi_1, \xi_2), \dots, (\xi_{d-s-1}, \xi_{d-s})$. We denote these roots by $\tau_1 < \dots < \tau_{d-1}$.

For each $t \geq 0$, the polynomial Y_d changes sign in each of the intervals (η_j, τ_j) , $j = 1, \dots, s - 1$, and in each of the intervals (τ_{s+i-1}, ξ_i) , $i = 2, \dots, d - s$, so it has a root there. This makes not less than $d - 2$ distinct real roots.

If $\tau_s > 0$ (resp. $\tau_s < 0$), then Y_d changes sign in each of the intervals $(\eta_s, 0)$ and (τ_s, ξ_1) (resp. (η_s, τ_s) and $(0, \xi_1)$), so it has two more real distinct roots. Hence for any $t \geq 0$, Y_d is hyperbolic, with d distinct roots. \square

Proof of Theorem 3. We remind that we denote by H_{\pm}^k not only the graphs mentioned in Theorem 4, but also the corresponding functions.

A) We prove Theorem 3 by induction on d . The induction base are the cases $d = 2$ and $d = 3$, see Examples 2 and 3.

Suppose that Theorem 3 holds true for $d = d_0 \geq 3$. Set $d := d_0 + 1$. As in the proof of Theorem 2 we set $Q_d := x^d + a_{d-1}x^{d-1} + \dots + a_2x^2 + a_1x + a_0$ and $Q_{d-1} := x^{d-1} + a_{d-1}x^{d-2} + \dots + a_2x + a_1$, so that the intersection $\Pi_d \cap \{a_0 = 0\}$ can be identified with Π_{d-1} .

B) We remind that any stratum (or component) U of Π_{d-1}^* is of the form (see Notation 2 and Remark 4)

$$U = S_{1^{d-1}}(\sigma_1, \dots, \sigma_{d-1}) = \{a' \in S_{1^{d-1}} \mid \text{sign}(a_j) = \sigma_j, 1 \leq j \leq d - 1\}.$$

Starting with such a component U (hence $U = U^{d-1}$), we construct in several steps the components U_+ and U_- of the set Π_d^* sharing with U the signs of the coefficients a_{d-1}, \dots, a_1 . One has $a_0 > 0$ in U_+ and $a_0 < 0$ in U_- .

At the first step we construct the sets $U_{1,\pm}$ as follows. We remind that the projections π^k and their fibres \tilde{f}_k were defined in part (1) of Theorem 4. Each fibre \tilde{f}_d of the projection π^d which is over a point of U is a segment, see part (1) of Theorem 4. If $Q_{d-1} \in U$, then for $\varepsilon > 0$ small enough, both polynomials $xQ_{d-1} \pm \varepsilon$ are hyperbolic. Indeed, all roots of Q_{d-1} are real and simple. The set $U_{1,+}$ (resp.

$U_{1,-}$) is the union of the interior points of these fibres \tilde{f}_d which are with positive (resp. with negative) a_0 -coordinates. Thus

$$U_{1,+} = \{a \in \tilde{f}_d \mid a' \in U, 0 < a_0 < H_+^d(a')\} \quad \text{and}$$

$$U_{1,-} = \{a \in \tilde{f}_d \mid a' \in U, H_-^d(a') < a_0 < 0\},$$

see part (6) of Theorem 4). Hence the sets $U_{1,\pm}$ are open, non-empty and contractible.

For $d \geq 2$, the intersection $\Pi_d \cap \{a_0 = 0\}$ is strictly included in the projection Π_d^{d-1} of Π_d in $Oa_1 \cdots a_{d-1}$. Therefore one can expect that the sets $U_{1,\pm}$ are not the whole of two components of Π_d^* . We construct contractible sets $U_{1,\pm} \subset U_{2,\pm} \subset \cdots \subset U_{d-1,\pm}$, where for $1 \leq j \leq d-1$, the signs of the coordinates a_j of each point of $U_{k,+}$ (resp. $U_{k,-}$) are defined by σ , and $U_{d-1,\pm}$ are components of Π_d^* . One has $a_0 > 0$ in $U_{k,+}$ and $a_0 < 0$ in $U_{k,-}$.

C) Recall that the set U consists of all the points between the graphs L_{\pm}^{d-1} of two continuous functions defined on U^{d-2} :

$$U = \{a' \mid L_-^{d-1}(a'') < a'' < L_+^{d-1}(a''), a'' \in U^{d-2}\},$$

see Notation 2. Thus $(L_+^{d-1} \cup L_-^{d-1}) \subset \partial U$. Depending on the sign of a_1 in U , for each of these graphs, part or the whole of it could belong to the hyperplane $a_1 = 0$.

Consider a fibre \tilde{f}_d over a point of one of the graphs L_{\pm}^{d-1} and not belonging to the hyperplane $a_1 = 0$. A priori the two endpoints of the fibre cannot have a_0 -coordinates with opposite signs. Indeed, if this were the case for the fibre over $a' = a^{*'}$ (see Notation 2), then for all fibres over a' close to $a^{*'}$, these signs would also be opposite, because the functions L_{\pm}^d , whose values are the values of the a_0 -coordinates of the endpoints, are continuous. Hence all these fibres \tilde{f}_d intersect the hyperplane $a_0 = 0$ (see part (1) of Theorem 4), but not the hyperplane $a_1 = 0$. Hence the point $a^{*'}$ is an interior point of Π_d (hence of U as well) and not a point of ∂U which is a contradiction, see part (2) of Theorem 4.

Both endpoints cannot have non-zero coordinates of the same sign, because then in the same way the fibres \tilde{f}_d over all points a' close to $a^{*'}$ would not intersect the hyperplane $a_0 = 0$ hence $a^{*'} \notin \bar{U}$, so $a^{*'} \notin \partial U$.

Hence the following three possibilities remain:

- a) both endpoints have zero a_0 -coordinates;
- b) one endpoint has a zero and the other endpoint has a positive a_0 -coordinate;
- c) one endpoint has a zero and the other endpoint has a negative a_0 -coordinate.

D) Consider the points of the graph L_+^{d-1} which do not belong to the hyperplane $a_1 = 0$ (for L_-^{d-1} the reasoning is similar). If for $B \in (L_+^{d-1} \setminus \{a_1 = 0\})$, possibility a) takes place, then there is nothing to do.

Suppose that possibility b) takes place. Denote by $a_{j,B}$ the coordinates of the point B (hence $a_{0,B} = 0$). For each such point B , fix the coordinates $a_j = a_{j,B}$ for $j \neq 1$ and increase a_1 . The interior points of the corresponding fibres \tilde{f}_d (when non-void) have the same signs of their a_0 -coordinates, hence these signs are positive. Then for some $a_1 = a_{1,C} > a_{1,B}$, one has either $a_{1,C} = 0$ (this can happen only when $a_{1,B} < 0$) or the point C belongs to the graph H_+^{d-1} and for $a_1 > a_{1,C}$, the fibres \tilde{f}_d are void, see Theorem 4.

In both these situations we add to the set $U_{1,+}$ the points of the interior of all fibres \tilde{f}_d over the interval $[a_{1,B}, a_{1,C}]$ (with $a_j = a_{j,B}$ for $j \neq 1$), over all points $B \in (L_+^{d-1} \setminus \{a_1 = 0\})$.

If possibility c) takes place, then we fix again $a_{j,B}$ for $j \neq 1$ and increase a_1 . The interior points of the corresponding fibres \tilde{f}_d (when non-void) have negative sign of their a_0 -coordinates. We add to the set $U_{1,-}$ the interior points of all fibres \tilde{f}_d over the interval $[a_{1,B}, a_{1,C}]$ (with $a_j = a_{j,B}$ for $j \neq 1$), over all points $B \in (L_+^{d-1} \setminus \{a_1 = 0\})$.

E) We perform a similar reasoning and construction with L_-^{d-1} (in which the role of H_+^{d-1} is played by H_-^{d-1}). In this case a_1 is to be decreased, one has $a_{1,C} < a_{1,B}$ and the interval $[a_{1,B}, a_{1,C}]$ is to be replaced by the interval $(a_{1,C}, a_{1,B}]$.

F) Thus we have enlarged the sets $U_{1,\pm}$; the new sets are denoted by $U_{2,\pm}$:

$$\begin{aligned}
 U_{2,+} &= U_{1,+} \cup \{a \in \Pi_d^* \mid a'' \in U^{d-2}, L_+^{d-1}(a'') < a_1 < H_+^{d-1}(a''), \\
 &\quad \text{if } L_+^{d-1}(a'') \geq 0, L_+^{d-1}(a'') < a_1 < \min(0, H_+^{d-1}(a'')), \text{ if } L_+^{d-1}(a'') < 0\}, \\
 U_{2,-} &= U_{1,-} \cup \{a \in \Pi_d^* \mid a'' \in U^{d-2}, H_-^{d-1}(a'') < a_1 < L_-^{d-1}(a''), \\
 &\quad \text{if } L_-^{d-1}(a'') \leq 0, \max(0, H_-^{d-1}(a'')) < a_1 < L_-^{d-1}(a''), \text{ if } L_-^{d-1}(a'') > 0\}.
 \end{aligned}$$

The sets $U_{1,\pm}$ and $U_{2,\pm}$ satisfy the conclusion of Theorem 3. We denote the graphs L_{\pm}^k defined for the sets $U_{1,\pm}$ and $U_{2,\pm}$ by $L_{1,\pm}^k$ and $L_{2,\pm}^k$. The construction of these graphs implies that they are graphs of continuous functions (because such are the graphs H_{\pm}^k). The set $U_{1,+} \cup U_{1,-}$ (resp. $U_{2,+} \cup U_{2,-}$) contains all points of the set $(\pi^d)^{-1}(U) \cap \Pi_{d,\sigma}^*$ (resp. $(\pi^{d-1} \circ \pi^d)^{-1}(U^{d-2}) \cap \Pi_{d,\sigma}^*$).

G) We remind that $\tilde{f}_d = f_{d-1}^\circ$, see Remark 5. Suppose that the sets $U_{s,\pm}$, $2 \leq s \leq d-3$, are constructed such that they satisfy the conclusion of Theorem 3 (the graphs L_{\pm}^k are denoted by $L_{s,\pm}^k$) and that the set $U_{s,+} \cup U_{s,-}$ contains all points of the set $(\pi^{d-s+1} \circ \dots \circ \pi^d)^{-1}(U^{d-s}) \cap \Pi_{d,\sigma}^*$.

Consider a point $D \in L_+^{d-s}$ which does not belong to the hyperplane $a_s = 0$. For the fibre f_{d-s}° of the projection $\pi^{d-s+1} \circ \dots \circ \pi^d$ which is over D (see Remark 5) one of the three possibilities takes place:

a') the minimal and the maximal possible value of the a_s -coordinate of the points of the fibre are zero;

b') the minimal possible value is 0 and the maximal possible value is positive;

c') the minimal possible value is negative and the maximal possible value is 0.

It is not possible to have both the maximal and minimal possible value of a_s non-zero, because in this case the point D does not belong to the set ∂U^{d-s} . This is proved by analogy with C). With regard to Remark 5, when the fibre f_{d-s}^\diamond is not a point, then the maximal (resp. the minimal) value of a_s is attained at one of the 0-dimensional cells (resp. at the other 0-dimensional cell) and only there. This can be deduced from part (2) of Theorem 4.

H) When possibility a') takes place, then there is nothing to do. Suppose that possibility b') takes place. Denote by $a_{j,D}$ the coordinates of the point D (hence $a_{0,D} = \dots = a_{s-1,D} = 0$). Fix $a_{j,D}$ for $j \neq s$ and increase a_s . Then for some $a_s = a_{s,E} > a_{s,D}$, one has either $a_{s,E} = 0$ (which is possible only if $a_{s,D} < 0$) or the point E belongs to the graph H_+^{d-s} . In this case we add to the set $U_{s,+}$ the points of the interior of all fibres f_{d-s}^\diamond over the interval $[a_{s,D}, a_{s,E}]$ (with $a_j = a_{j,D}$ for $j \neq s$), over all points $D \in (L_+^{d-s} \setminus \{a_s = 0\})$. The a_{s-1} -coordinates of all points thus added are positive.

If possibility c') takes place, then we fix again $a_{j,D}$ for $j \neq s$ and increase a_s . We add to the set $U_{s,-}$ the points of the interior of all fibres f_{d-s}^\diamond over the interval $[a_{s,D}, a_{s,E}]$ (with $a_j = a_{j,D}$ for $j \neq s$), over all points $D \in L_+^{d-s} \setminus \{a_s = 0\}$. The a_{s-1} -coordinates of all points thus added are negative.

We consider in a similar way the graph L_-^{d-s} in which case the role of H_+^{d-s} is played by H_-^{d-s} , a_s is to be decreased, one has $a_{s,E} < a_{s,D}$ and the interval $[a_{s,D}, a_{s,E}]$ is to be replaced by the interval $(a_{s,E}, a_{s,D}]$.

I) We have thus constructed the sets $U_{s+1,\pm}$ which satisfy the conclusion of Theorem 3:

$$\begin{aligned} U_{s+1,+} &= U_{s,+} \cup \{a \in \Pi_d^* \mid a^{(s+1)} \in U^{d-s-1}, \\ &\quad L_{s,+}^{d-s}(a^{(s+1)}) < a_s < H_+^{d-s}(a^{(s+1)}), \text{ if } L_{s,+}^{d-s}(a^{(s+1)}) \geq 0, \\ &\quad L_{s,+}^{d-s}(a^{(s+1)}) < a_s < \min(0, H_+^{d-s}(a^{(s+1)})), \text{ if } L_{s,+}^{d-s}(a^{(s+1)}) < 0 \}, \\ U_{s+1,-} &= U_{s,-} \cup \{a \in \Pi_d^* \mid a^{(s+1)} \in U^{d-s-1}, \\ &\quad H_-^{d-s}(a^{(s+1)}) < a_s < L_{s,-}^{d-s}(a^{(s+1)}), \text{ if } L_{s,-}^{d-s}(a^{(s+1)}) \leq 0, \\ &\quad \max(0, H_-^{d-s}(a^{(s+1)})) < a_s < L_{s,-}^{d-s}(a^{(s+1)}), \text{ if } L_{s,-}^{d-s}(a^{(s+1)}) > 0 \}. \end{aligned}$$

The set $U_{s+1,+} \cup U_{s+1,-}$ contains all points of the set $(\pi^{d-s} \circ \dots \circ \pi^d)^{-1}(U^{d-s-1}) \cap \Pi_{d,\sigma}^*$. It should be noticed that as the fibres f_{d-s}^\diamond contain cells of dimension from 0 to

s , all graphs $L_{s,\pm}^k$ would have to be changed when passing from $L_{s,\pm}^k$ to $L_{s+1,\pm}^k$. The new graphs are graphs of continuous functions; this follows from the construction and from the fact that such are the graphs H_{\pm}^k .

J) One can construct the sets $U_{d-1,\pm}$ in a similar way. The only difference is the fact that there is a graph H_+^2 , but not a graph H_-^2 , see Example 2:

$$\begin{aligned} U_{d-1,+} &= U_{d-2,+} \cup \{a \in \Pi_d^* \mid a^{(d-1)} \in U^1, \\ &L_{d-2,+}^2(a^{(d-1)}) < a_{d-2} < H_+^2(a^{(d-1)}), \text{ if } L_{d-2,+}^2(a^{(d-1)}) \geq 0, \\ &L_{d-2,+}^2(a^{(d-1)}) < a_{d-2} < \min(0, H_+^2(a^{(d-1)})), \text{ if } L_{d-2,+}^2(a^{(d-1)}) < 0 \}, \\ U_{d-1,-} &= U_{d-2,-} \cup \{a \in \Pi_d^* \mid a^{(d-1)} \in U^1, \\ &a_{d-2} < L_{d-2,-}^2(a^{(d-1)}), \text{ if } L_{d-2,-}^2(a^{(d-1)}) \leq 0, \\ &0 < a_{d-2} < L_{d-2,-}^2(a^{(d-1)}), \text{ if } L_{d-2,-}^2(a^{(d-1)}) > 0 \}. \end{aligned}$$

We set $U_{\pm} := U_{d-1,\pm}$. The set $U_+ \cup U_-$ contains all points from the set $(\pi^2 \circ \dots \circ \pi^d)^{-1}(U^1) \cap \Pi_{d,\sigma}^*$. The sets U_{\pm} satisfy the conclusion of Theorem 3. Hence they are contractible.

K) The functions L_{\pm}^k encountered throughout the proof of the theorem can be extended by continuity on the closures of the sets on which they are defined, because this is the case of the functions H_{\pm}^k . Moreover, fibres \tilde{f}_k which are points appear only in case they are over points of the graphs H_{\pm}^{k-1} . Hence this describes the only possibility for the values of the functions L_{\pm}^k to coincide. \square

6. COMMENTS AND OPEN PROBLEMS

One could try to generalize Theorem 2 by considering instead of the set Π_d^* the set $R_{3,d}$, i. e. by dropping the requirement the polynomial Q_d to be hyperbolic. So an open problem can be formulated like this:

Open problem 1. *For a given degree d , consider the triples $(\sigma, \text{pos}, \text{neg})$ compatible with Descartes' rule of signs. Is it true that for each such triple, the corresponding subset of the set $R_{3,d}$ is either contractible or empty?*

The difference between this open problem and Theorem 2 is the necessity to check whether the subset is empty or not (see part (3) of Theorem 1). For instance, if $d = 4$, then for neither of the triples

$$((+, -, +, +), 2, 0) \quad \text{and} \quad ((-, -, -, +), 0, 2)$$

(both compatible with Descartes' rule of signs) does there exist a polynomial $x^4 + a_3x^3 + a_2x^2 + a_1x + a_0$ with signs of the coefficients a_j as defined by σ and with 2 positive and 0 negative or with 0 positive and 2 negative roots respectively, see [12] (all roots are assumed to be simple).

The question of realizability of triples (σ, pos, neg) has been asked in [2]. The exhaustive answer to this question is known for $d \leq 8$. For $d = 4$, it is due to D. Grabiner ([12]), for $d = 5$ and 6, to A. Albouy and Y. Fu ([1]), for $d = 7$ and partially for $d = 8$, to J. Forsgård, V. P. Kostov and B. Shapiro ([7] and [8]) and for $d = 8$ the result was completed in [15]. Other results in this direction can be found in [4], [6] and [16].

Remarks 6. (1) It is not easy to imagine how one could prove that all components of $R_{3,d}$ are either contractible or empty without giving an exhaustive answer to the question which triples (σ, pos, neg) are realizable and which are not. Unfortunately, at present, giving such an answer for any degree d is out of reach.

(2) If one can prove not contractibility of the non-empty components, but only that they are (simply) connected, would also be of interest.

For a degree d univariate real monic polynomial Q_d without vanishing coefficients, one can define the couples (pos_ℓ, neg_ℓ) of the numbers of positive and negative roots of $Q_d^{(\ell)}$, $\ell = 0, 1, \dots, d-1$. One can observe that the d couples (pos_ℓ, neg_ℓ) define the signs of the coefficients of Q_d and that their choice must be compatible not only with Descartes' rule of signs, but also with Rolle's theorem. We call such d -tuples of couples *compatible* for short. We assume that for $\ell = 0, 1, \dots, d-1$, all real roots of $Q_d^{(\ell)}$ are simple and non-zero.

To have a geometric idea of the situation we define the discriminant sets $\tilde{\Delta}_j$, $j = 1, \dots, d$ as the sets Δ_j defined in the spaces $Oa_{d-j} \dots a_{d-1}$ for the polynomials $Q_d^{(d-j)}$. In particular, $\tilde{\Delta}_d = \Delta_d$. For $j = 1, \dots, d-1$, we set $\Delta_j := \tilde{\Delta}_j \times Oa_0 \dots a_{d-j-1}$. We define the set $R_{4,d}$ as

$$R_{4,d} := \mathbb{R}^d \setminus ((\cup_{j=1}^d \Delta_j) \cup (\cup_{j=0}^{d-1} \{a_j = 0\})) .$$

For $d \leq 5$, the question when a subset of $R_{4,d}$ defined by a given compatible d -tuple of couples (pos_ℓ, neg_ℓ) is empty is considered in [5].

Open problem 2. *Given the d compatible couples (pos_ℓ, neg_ℓ) , is it true that the subset of $R_{4,d}$ defined by them is either connected (eventually contractible) or empty? In other words, is it true that each d -tuple of such couples defines either exactly one or none of the components of the set $R_{4,d}$?*

Some problems connected with comparing the moduli of the positive and negative roots of hyperbolic polynomials are treated in [18], [20] and [19]. Other problems concerning hyperbolic polynomials are to be found in [17]. A tropical analog of Descartes' rule of signs is discussed in [9].

ACKNOWLEDGEMENT. B. Z. Shapiro from the University of Stockholm attracted the author's attention to Open Problem 1 and suggested the proof of part (1) of Theorem 1. The remarks of the anonymous referee allowed to improve the clarity of the proofs of the theorems.

7. REFERENCES

- [1] Albouy, A., Fu, Y.: Some remarks about Descartes' rule of signs. *Elem. Math.* **69** (2014), 186–194. Zbl 1342.12002, MR3272179
- [2] Anderson, B., Jackson, J., and Sitharam, M.: Descartes' rule of signs revisited. *Amer. Math. Monthly* **105** (1998), 447–451. Zbl 0913.12001, MR1622513
- [3] Arnold, V.I.: Hyperbolic polynomials and Vandermonde mappings. *Funct. Anal. Appl.* **20** (1986), 52–53.
- [4] Cheriha, H., Gati, Y., and Kostov, V.P.: A nonrealization theorem in the context of Descartes' rule of signs. *Ann. Sofia Univ. St. Kliment Ohridski, Fac. Math. and Inf.* **106** (2019), 25–51.
- [5] Cheriha, H., Gati, Y., and Kostov, V.P.: Descartes' rule of signs, Rolle's theorem and sequences of compatible pairs. *Studia Scientiarum Mathematicarum Hungarica* **57:2** (2020), 165–186, DOI: <https://doi.org/10.1556/012.2020.57.2.1463>
- [6] Cheriha, H., Gati, Y., and Kostov, V.P.: On Descartes' rule for polynomials with two variations of sign. *Lithuanian Math. J.* **60** (2020), 456–469.
- [7] Forsgård, J., Kostov, V.P., and Shapiro, B.: Could René Descartes have known this? *Exp. Math.* **24** (4) (2015), 438–448. Zbl 1326.26027, MR3383475
- [8] Forsgård, J., Kostov, V.P., and Shapiro, B.: Corrigendum: "Could René Descartes have known this?". *Exp. Math.* **28** (2) (2019), 255–256.
- [9] Forsgård, J., Novikov, D., and Shapiro, B.: A tropical analog of Descartes' rule of signs. *Int. Math. Res. Not. IMRN* 2017, no. 12, 3726–3750. arXiv:1510.03257 [math.CA].
- [10] Fourier, J.: Sur l'usage du théorème de Descartes dans la recherche des limites des racines. *Bulletin des sciences par la Société philomatique de Paris (1820)* 156–165, 181–187; œuvres 2, 291–309, Gauthier- Villars, 1890.
- [11] Givental, A.B.: Moments of random variables and the equivariant Morse lemma (in Russian). *Uspekhi Mat. Nauk* **42** (1987), 221–222.
- [12] Grabiner, D.J.: Descartes' rule of signs: another construction. *Amer. Math. Monthly* **106** (1999), 854–856. Zbl 0980.12001, MR1732666
- [13] Jullien, V.: Descartes La "Geometrie" de 1637.
- [14] Kostov, V.P.: On the geometric properties of Vandermonde's mapping and on the problem of moments. *Proc. Royal Soc. Edinburgh* **112A** (1989), 203–211.
- [15] Kostov, V.P.: On realizability of sign patterns by real polynomials. *Czech. Math. J.* **68 (143)** (2018), no. 3, 853–874.
- [16] Kostov, V.P.: Polynomials, sign patterns and Descartes' rule of signs. *Math. Bohemica* **144** (2019), no. 1, 39–67.

- [17] Kostov, V. P.: Topics on hyperbolic polynomials in one variable. *Panoramas et Synthèses* **33** (2011), vi + 141 p. SMF.
- [18] Kostov, V. P.: Descartes' rule of signs and moduli of roots. *Publicationes Mathematicae Debrecen* **96**(1-2) (2020) 161–184, DOI: 10.5486/PMD.2020.8640.
- [19] Kostov, V. P.: Hyperbolic polynomials and canonical sign patterns. *Serdica Math. J.* **46**(2) (2020) , 135–150, arXiv:2006.14458.
- [20] Kostov, V. P.: Hyperbolic polynomials and rigid moduli orders. arXiv:2008.11415.
- [21] Méguerditchian, I.: *Géométrie du discriminant réel et des polynômes hyperboliques*. Thesis defended in 1991 at the University Rennes 1.

Received on July 9, 2021

VLADIMIR PETROV KOSTOV
Université Côte d'Azur
CNRS, LJAD
FRANCE
E-mail: vladimir.kostov@unice.fr

ГОДИШНИК НА СОФИЙСКИЯ УНИВЕРСИТЕТ „СВ. КЛИМЕНТ ОХРИДСКИ“

ФАКУЛТЕТ ПО МАТЕМАТИКА И ИНФОРМАТИКА

Том 107

ANNUAL OF SOFIA UNIVERSITY „ST. KLIMENT OHRIDSKI“

FACULTY OF MATHEMATICS AND INFORMATICS

Volume 107

EXAMPLES OF HNN-EXTENSIONS WITH NONTRIVIAL QUASI-KERNELS

NIKOLAY A. IVANOV

We introduce some examples of HNN-extensions motivated by the problems of C^* -simplicity and unique trace property. Moreover, we prove that our examples are not inner amenable and identify a relatively large, simple, normal subgroup in each one.

Keywords: C^* -simplicity, HNN-extensions, inner amenability.

2020 Math. Subject Classification: Primary: 22D25, 20E06; Secondary: 46L05, 43A07, 20E08.

1. INTRODUCTION AND PRELIMINARIES

1.1. INTRODUCTION

The questions of C^* -simplicity and unique trace property for a discrete group have been studied extensively. By definition, a discrete group G is C^* -simple if the C^* -algebra associated to the left regular representation, $C_r^*(G)$, is simple; likewise it has the unique trace property if $C_r^*(G)$ has a unique tracial state. An extensive introduction to that topic was given by de la Harpe ([6]). Recently, Kalantar and Kennedy ([10]) gave a necessary and sufficient condition for C^* -simplicity in terms of action on the Furstenberg boundary of the group in question. Later, Breuillard, Kalantar, Kennedy, and Ozawa ([2]) studied further the question of C^* -simplicity

and also showed that a group has the unique trace property if and only if its amenable radical is trivial. They also showed that C^* -simplicity implies the unique trace property. The reverse implication was disproven by examples given by Le Boudec ([11]). In the case of group amalgamations and HNN-extensions, the kernel controls the uniqueness of trace, and the quasi-kernels control the C^* -simplicity.

The notion of inner amenability for discrete groups was introduced by Effros ([5]) as an analogue to Property Γ for II_1 factors that was introduced by Murray and von Neumann ([12]). By definition, a discrete group G is inner amenable if there exist a conjugation invariant, positive, finitely additive, probability measure on $G \setminus \{1\}$. Effros showed that Property Γ implies inner amenability, but the reverse implication doesn't hold, as demonstrated by Vaes ([14]).

Our examples (all of which being HNN-extensions) stem from the questions of C^* -simplicity and the unique trace properties for groups. In particular, all of our examples have the unique trace property, and we also determine the C^* -simple ones and the non- C^* -simple ones. The examples of section 2 generalize the example given in [3, Section 5] (which corresponds to the group $\Lambda[\text{Sym}(2), \text{Sym}(2)]$ of section 2). There is a resemblance to the groups introduced by Le Boudec in [11] since they all act on trees. The main benefit is that our groups are given concretely by generators and relations, which makes them more tractable to investigate some further properties they possess.

We study some additional analytic properties of our examples. We show that they are all non-inner-amenable by showing that they are finitely fledged - a property that we introduce in [8].

We also explore some of the group-theoretical properties of our groups. We remark that they are not finitely presented. Also, under some mild natural assumptions, we show that each group has a relatively large, simple, normal subgroup.

1.2. PRELIMINARIES

For a group Γ acting on a set X , we denote the set-wise stabilizer of a subset $Y \subset X$ by

$$\Gamma_{\{Y\}} \equiv \{ g \in \Gamma \mid gY = Y \}$$

and the point-wise stabilizer of a subset $Y \subset X$ by

$$\Gamma_{(Y)} \equiv \{ g \in \Gamma \mid gy = y, \forall y \in Y \}.$$

For a point $x \in X$, we denote its stabilizer by

$$\Gamma_x = \{ g \in \Gamma \mid gx = x \}.$$

Note that, $\Gamma_{\{Y\}}$, $\Gamma_{(Y)}$, and Γ_x are all subgroups of Γ . Also note that,

$$g\Gamma_{\{Y\}}g^{-1} = \Gamma_{\{gY\}}, \quad g\Gamma_xg^{-1} = \Gamma_{gx}, \quad \text{and} \quad g\Gamma_{(Y)}g^{-1} = \Gamma_{(gY)}.$$

For a group G and its subgroup H , by $\langle\langle H \rangle\rangle_G$ or by $\langle\langle H \rangle\rangle$, we denote the normal closure of H in G .

For some general references on group amalgamations and HNN-extensions see, e.g., [1], [4], [13], [7], etc.

Let $G = \langle X \mid R \rangle$ be a group; let H be a subgroup of G ; and let $\theta : H \hookrightarrow G$ be a monomorphism. Then an HNN-extension of this data (named after G. Higman, B. Neumann, H. Neumann) is the group

$$HNN(G, H, \theta) \equiv G *_\theta \equiv \langle X \sqcup \{\tau\} \mid R \sqcup \{\theta(h) = \tau^{-1}h\tau \mid h \in H\} \rangle.$$

It is convenient to denote $H_{-1} \equiv H$ and $H_1 \equiv \theta(H)$. Every element $\gamma \in HNN(G, H, \theta)$ can be written in reduced form as

$$\begin{aligned} \gamma &= g_1 \tau^{\varepsilon_1} \cdots g_n \tau^{\varepsilon_n} g_{n+1}, \text{ where } n \in \mathbb{N}, g_1, \dots, g_{n+1} \in G, \varepsilon_1, \dots, \varepsilon_n = \pm 1, \\ &\text{and where if } \varepsilon_{i+1} = -\varepsilon_i \text{ for } 1 \leq i \leq n-1, \text{ then } g_{i+1} \notin H_{\varepsilon_i}. \end{aligned}$$

If S_ε is a set of left coset representatives for G/H_ε , where $\varepsilon = \pm 1$, satisfy $S_{-1} \cap S_1 = \{1\}$, then every element $\gamma \in HNN(G, H, \theta)$ can be uniquely written in normal form as

$$\begin{aligned} \gamma &= s_1 \tau^{\varepsilon_1} s_2 \tau^{\varepsilon_2} \cdots s_n \tau^{\varepsilon_n} g, \text{ where } n \in \mathbb{N}_0, g \in G, \varepsilon_i = \pm 1, s_i \in S_{-\varepsilon_i}, \forall 1 \leq i \leq n, \\ &\text{and where if } \varepsilon_{i-1} = -\varepsilon_i \text{ for } 2 \leq i \leq n, \text{ then } s_i \neq 1. \end{aligned}$$

The HNN-extension $HNN(G, H, \theta)$ is called nondegenerate if either $H \neq G$ or $\theta(H) \neq G$ and is called non-ascending if $H \neq G \neq \theta(G)$.

The Bass-Serre tree $T(HNN(G, H, \theta))$ of $HNN(G, H, \theta)$ is the graph, that can be shown to be a tree, consisting of a vertex set

$$\begin{aligned} \text{Vertex}(HNN(G, H, \theta)) &= \\ &\{G\} \cup \{s_1 \tau^{\varepsilon_1} s_2 \tau^{\varepsilon_2} \cdots s_n \tau^{\varepsilon_n} G \mid n \in \mathbb{N}, s_1 \tau^{\varepsilon_1} s_2 \tau^{\varepsilon_2} \cdots s_n \tau^{\varepsilon_n} \text{ is in normal form}\} \end{aligned}$$

and an edge set

$$\begin{aligned} \text{Edge}(HNN(G, H, \theta)) &= \\ &\{H\} \cup \{s_1 \tau^{\varepsilon_1} s_2 \tau^{\varepsilon_2} \cdots s_n \tau^{\varepsilon_n} s_{n+1} H \mid n \in \mathbb{N}, s_1 \tau^{\varepsilon_1} s_2 \tau^{\varepsilon_2} \cdots s_n \tau^{\varepsilon_n} \text{ is in normal form}\}. \end{aligned}$$

The group $HNN(G, H, \theta)$ acts on $T(HNN(G, H, \theta))$ by left multiplication.

The vertex $v = s_1 \tau^{\varepsilon_1} s_2 \tau^{\varepsilon_2} \cdots s_n \tau^{\varepsilon_n} G$ is adjacent to the vertex $w = s_1 \tau^{\varepsilon_1} s_2 \tau^{\varepsilon_2} \cdots s_n \tau^{\varepsilon_n} s_{n+1} \tau^{\varepsilon_{n+1}} G$ with connecting edge

$$e = \begin{cases} s_1 \tau^{\varepsilon_1} s_2 \tau^{\varepsilon_2} \cdots s_n \tau^{\varepsilon_n} s_{n+1} \tau^{\varepsilon_{n+1}} H & \text{if } \varepsilon_{n+1} = -1, \\ s_1 \tau^{\varepsilon_1} s_2 \tau^{\varepsilon_2} \cdots s_n \tau^{\varepsilon_n} s_{n+1} H & \text{if } \varepsilon_{n+1} = 1. \end{cases}$$

To see the reason for this, we need to look at the stabilizers. The stabilizer of v is

$$HNN(G, H, \theta)_v = s_1 \tau^{\varepsilon_1} s_2 \tau^{\varepsilon_2} \cdots s_n \tau^{\varepsilon_n} G (s_1 \tau^{\varepsilon_1} s_2 \tau^{\varepsilon_2} \cdots s_n \tau^{\varepsilon_n})^{-1}$$

and the stabilizer of w is

$$HNN(G, H, \theta)_w = s_1 \tau^{\varepsilon_1} s_2 \tau^{\varepsilon_2} \cdots s_n \tau^{\varepsilon_n} s_{n+1} \tau^{\varepsilon_{n+1}} G (s_1 \tau^{\varepsilon_1} s_2 \tau^{\varepsilon_2} \cdots s_n \tau^{\varepsilon_n} s_{n+1} \tau^{\varepsilon_{n+1}})^{-1}.$$

Therefore the stabilizer of e is

$$\begin{aligned} HNN(G, H, \theta)_e &= HNN(G, H, \theta)_v \cap HNN(G, H, \theta)_w = \\ &= s_1 \tau^{\varepsilon_1} s_2 \tau^{\varepsilon_2} \cdots s_n \tau^{\varepsilon_n} s_{n+1} [G \cap \tau^{\varepsilon_{n+1}} G \tau^{-\varepsilon_{n+1}}] (s_1 \tau^{\varepsilon_1} s_2 \tau^{\varepsilon_2} \cdots s_n \tau^{\varepsilon_n} s_{n+1})^{-1} = \\ &= s_1 \tau^{\varepsilon_1} s_2 \tau^{\varepsilon_2} \cdots s_n \tau^{\varepsilon_n} s_{n+1} H_{-\varepsilon_{n+1}} (s_1 \tau^{\varepsilon_1} s_2 \tau^{\varepsilon_2} \cdots s_n \tau^{\varepsilon_n} s_{n+1})^{-1} = \\ &= \begin{cases} s_1 \tau^{\varepsilon_1} s_2 \tau^{\varepsilon_2} \cdots s_n \tau^{\varepsilon_n} s_{n+1} H (s_1 \tau^{\varepsilon_1} s_2 \tau^{\varepsilon_2} \cdots s_n \tau^{\varepsilon_n} s_{n+1})^{-1} & \text{if } \varepsilon_{n+1} = 1, \\ s_1 \tau^{\varepsilon_1} s_2 \tau^{\varepsilon_2} \cdots s_n \tau^{\varepsilon_n} s_{n+1} \tau^{\varepsilon_{n+1}} H \tau^{-\varepsilon_{n+1}} (s_1 \tau^{\varepsilon_1} s_2 \tau^{\varepsilon_2} \cdots s_n \tau^{\varepsilon_n} s_{n+1})^{-1} & \text{if } \varepsilon_{n+1} = -1. \end{cases} \end{aligned}$$

Finally, since $HNN(G, H, \theta)$ can be expressed as

$$HNN(G, H, \theta) = (G * \langle \tau \rangle) / \langle \langle \tau^{-1} h \tau \theta(h^{-1}) \mid h \in H \rangle \rangle,$$

it has the following universal property (see, e.g., [4], page 36):

Remark 1.1. *Let C be a group; let $\alpha : G \rightarrow C$ be a group homomorphism; and let $t \in C$ be an element for which the following holds: $t^{-1} \alpha(h) t = \alpha(\theta(h))$ for each $h \in H$. Then there is a unique group homomorphism $\beta : HNN(G, H, \theta) \rightarrow C$ satisfying $\beta|_G = \alpha$ and $\beta(\tau) = t$.*

To conclude this section, we recall that we called a group amenable if it has no nontrivial C^* -simple quotients ([9, Definition 7.1]). We showed in [9] that the class on amenable groups is a radical class, so every group has a unique maximal normal amenable subgroup, the amenable radical. Also, the class of amenable groups is closed under extensions. The amenable radical 'detects' C^* -simplicity the same way as the amenable radical 'detects' the unique trace property (see [9, Corollary 7.3] and [2, Theorem 1.3]).

2. HNN-EXTENSIONS

2.1. NOTATION, DEFINITIONS, QUASI-KERNELS

We use the following notations, some of which appear in [3]:

$$\begin{aligned} T_\varepsilon &= \{\gamma = g_0 \tau^\varepsilon g_1 \tau^{\varepsilon_1} \cdots g_n \tau^{\varepsilon_n} g_{n+1} \mid n \geq 0, \gamma \in \Lambda \text{ is reduced}\}, \\ T_\varepsilon^\dagger &= \{\gamma = \tau^\varepsilon g_1 \tau^{\varepsilon_1} \cdots g_n \tau^{\varepsilon_n} g_{n+1} \mid n \geq 0, \gamma \in \Lambda \text{ is reduced}\}. \end{aligned}$$

For $\varepsilon = \pm 1$, consider also the quasi-kernels defined in [3]:

$$K_\varepsilon \equiv \bigcap_{r \in \Lambda \setminus T_\varepsilon^\dagger} rHr^{-1}. \quad (1)$$

They satisfy the relation $\ker \Lambda = K_1 \cap K_{-1}$, where, by definition,

$$\ker \Lambda \equiv \bigcap_{r \in \Lambda} rHr^{-1}.$$

It follows from [3, Theorem 4.19] that Λ has the unique trace property if and only if $\ker \Lambda$ has the unique trace property. It also follows from [3, Theorem 4.20] that Λ is C^* -simple if and only if K_{-1} or K_1 is trivial or non-amenable provided Λ is a non-ascending HNN-extension and $\ker \Lambda$ is trivial.

We need the following results.

Remark 2.1. Consider the Bass-Serre tree $\Theta = \Theta[\Lambda]$ of the group

$$\Lambda = \text{HNN}(G, H, \theta) = \langle G, \tau \mid \tau^{-1}h\tau = \theta(h) \text{ for all } h \in H \rangle,$$

and consider the edge H connecting vertices G and τG . Denote by Θ_1 the full subtree of Θ consisting of all vertices $v \in \Theta$ satisfying $\text{dist}(v, G) < \text{dist}(v, \tau G)$. Also, denote by $\bar{\Theta}_1$ the full subtree of Θ consisting of all vertices $v \in \Theta$ satisfying $\text{dist}(v, G) > \text{dist}(v, \tau G)$. Likewise, consider the edge $\tau^{-1}H$ connecting vertices G and $\tau^{-1}G$. Then, denote by Θ_{-1} the full subtree of Θ consisting of all vertices $v \in \Theta$ satisfying $\text{dist}(v, G) < \text{dist}(v, \tau^{-1}G)$, and denote by $\bar{\Theta}_{-1}$ the full subtree of Θ consisting of all vertices $v \in \Theta$ satisfying $\text{dist}(v, G) > \text{dist}(v, \tau^{-1}G)$.

It is easy to see that $\bar{\Theta}_\varepsilon = \tau^\varepsilon \Theta_{-\varepsilon}$,

$$\Theta_\varepsilon = \{G\} \cup \{t_\varepsilon G \mid t_\varepsilon \in \Lambda \setminus T_\varepsilon^\dagger\}, \text{ and } \bar{\Theta}_\varepsilon = \{t_\varepsilon^\dagger G \mid t_\varepsilon^\dagger \in T_\varepsilon^\dagger\}.$$

Proposition 2.2. With the notation from the previous Remark, the following hold for each $\varepsilon = \pm 1$:

- (i) $K_\varepsilon = \Lambda_{(\Theta_\varepsilon)}$.
- (ii) $K_\varepsilon < H \cap \theta(H)$.
- (iii) $\gamma K_\varepsilon \gamma^{-1} = \Lambda_{(\gamma \Theta_\varepsilon)}$ for every $\gamma \in \Lambda$.

In particular $\Lambda_{(\bar{\Theta}_\varepsilon)} = \tau^\varepsilon K_{-\varepsilon} \tau^{-\varepsilon}$.

Proof. (i)

$$\begin{aligned} g \in K_\varepsilon &\iff r^{-1}gr \in H, \quad \forall r \in \Lambda \setminus T_\varepsilon^\dagger \iff gr \in rH, \quad \forall r \in \Lambda \setminus T_\varepsilon^\dagger \\ &\iff grH = rH, \quad \forall r \in \Lambda \setminus T_\varepsilon^\dagger \iff g \text{ fixes every edge of } \Theta_\varepsilon \\ &\iff g \in \Lambda_{(\Theta_\varepsilon)}. \end{aligned}$$

(ii) From (i), we know that every element $g \in K_\varepsilon$ fixes all vertices adjacent to G except for the vertex $\tau^\varepsilon G$, eventually. Therefore it also fixes $\tau^\varepsilon G$, so g fixes all edges around G . In particular, g fixes the edge H , so $g \in H$. Likewise, g fixes the edge $\tau^{-1}H$, so $g \in \tau^{-1}H\tau = \theta(H)$.

(iii) As in (i), we have

$$\begin{aligned} g \in \gamma K_\varepsilon \gamma^{-1} &\iff \gamma^{-1}g\gamma \in K_\varepsilon &\iff \gamma^{-1}g\gamma \in \Lambda_{(\Theta_\varepsilon)} \\ &\iff g \in \gamma \Lambda_{(\Theta_\varepsilon)} \gamma^{-1} &\iff g \in \Lambda_{(\gamma\Theta_\varepsilon)}. \end{aligned}$$

□

Lemma 2.3. *For $\varepsilon = \pm 1$, K_ε is a normal subgroup of $H_{-\varepsilon}$, and a normal subgroup of $H \cap \theta(H)$. Moreover, if $\ker \Lambda$ is trivial, then K_{-1} and K_1 have a trivial intersection and mutually commute.*

Proof. From Proposition 2.2 (ii), it follows that K_1 and K_{-1} are subgroups of $H \cap \theta(H)$. Take $h \in H_{-\varepsilon}$. Then

$$\begin{aligned} h \cdot T_\varepsilon^\dagger &= \{h\tau^\varepsilon g_1\tau^{\varepsilon_1} \cdots g_n\tau^{\varepsilon_n} g_{n+1} \mid n \geq 0, \tau^\varepsilon g_1\tau^{\varepsilon_1} \cdots g_n\tau^{\varepsilon_n} g_{n+1} \text{ is reduced}\} = \\ &= \{\tau^\varepsilon \theta^\varepsilon(h)g_1\tau^{\varepsilon_1} \cdots g_n\tau^{\varepsilon_n} g_{n+1} \mid n \geq 0, \tau^\varepsilon g_1\tau^{\varepsilon_1} \cdots g_n\tau^{\varepsilon_n} g_{n+1} \text{ is reduced}\} = T_\varepsilon^\dagger. \end{aligned}$$

This gives the first assertion. For the second assertion, take $k_\varepsilon \in K_\varepsilon$ for each $\varepsilon = \pm 1$. Then, from $K_\varepsilon \triangleleft H \cap \theta(H)$, it follows that $k_{-1}k_1^{-1}k_{-1}^{-1} \in K_1$ and $k_1k_{-1}k_1^{-1} \in K_{-1}$. Thus

$$K_{-1} \ni (k_1k_{-1}k_1^{-1})k_{-1}^{-1} = k_1(k_{-1}k_1^{-1}k_{-1}^{-1}) \in K_1,$$

and therefore $k_1k_{-1}k_1^{-1}k_{-1}^{-1} \in K_1 \cap K_{-1} = \ker \Lambda = \{1\}$. □

Lemma 2.4.

- (i) *Let $\gamma = \tau^{\varepsilon_n} g_n \cdots g_2 \tau^{\varepsilon_2} g_1 \tau^{\varepsilon_1} \in \Lambda$ be reduced. Then $\gamma \cdot T_{-\varepsilon}^\dagger \supset T_{-\varepsilon_n}^\dagger$. In particular, $K_{-\varepsilon_n} < \gamma K_{-\varepsilon} \gamma^{-1}$.*
- (ii) *Let $\gamma \in G \setminus H_\varepsilon$. Then $\gamma T_{-\varepsilon}^\dagger \cap T_{-\varepsilon}^\dagger = \emptyset$. In particular, $\gamma K_{-\varepsilon} \gamma^{-1} \cap K_{-\varepsilon} = \ker \Lambda$.*
- (iii) *Let $\gamma \in \Lambda$ be a reduced word starting and ending with τ^ε . Then $T_{-\varepsilon}^\dagger \cap \gamma T_\varepsilon^\dagger = \emptyset$. In particular, $K_{-\varepsilon} \cap \gamma K_\varepsilon \gamma^{-1} = \ker \Lambda$.*

Proof. (i) Observe that

$$\begin{aligned} \gamma \cdot T_{-\varepsilon} &= \\ &\supset \{\gamma \cdot \tau^{-\varepsilon} g_1^{-1} \tau^{-\varepsilon_1} \cdots g_n^{-1} \tau^{-\varepsilon_n} \cdot \tau^{-\varepsilon_n} \cdot g_{n+1} \tau^{\varepsilon_{n+1}} g_{n+2} \tau^{\varepsilon_{n+2}} \cdots g_{n+m} \tau^{\varepsilon_{n+m}} g_{n+m+1} \mid \\ &\quad m \geq 0, \tau^{-\varepsilon_n} g_{n+1} \tau^{\varepsilon_{n+1}} g_{n+2} \cdots g_{n+m} \tau^{\varepsilon_{n+m}} g_{n+m+1} \text{ is reduced}\} \\ &= \{\lambda = \tau^{-\varepsilon_n} g_{n+1} \tau^{\varepsilon_{n+1}} g_{n+2} \cdots g_{n+m} \tau^{\varepsilon_{n+m}} g_{n+m+1} \mid m \geq 0, \lambda \text{ is reduced}\} \\ &= T_{-\varepsilon_n}. \end{aligned}$$

The second statement follows from the observation

$$\gamma \cdot (\Lambda \setminus T_{-\varepsilon}^\dagger) = \Lambda \setminus \gamma T_{-\varepsilon}^\dagger \subset \Lambda \setminus T_{-\varepsilon_n}^\dagger.$$

(ii) and (iii) follow easily. \square

Lemma 2.5. *Let $\gamma = g_{n+1}\tau^{\varepsilon_n}g_n \cdots g_2\tau^{\varepsilon_1}g_1\tau^\varepsilon$, $\gamma' = g'_{n+1}\tau^{\varepsilon'_n}g'_n \cdots g'_2\tau^{\varepsilon'_1}g'_1\tau^\varepsilon$, and $\gamma'' = g''_{n+1}\tau^{\varepsilon''_n}g''_n \cdots g''_2\tau^{\varepsilon''_1}g''_1\tau^{-\varepsilon}$ be reduced, where $n \geq 0$ and $\varepsilon = \pm 1$. Then:*

- (i) *If $(\gamma')^{-1}\gamma \in H_{-\varepsilon}$, then $\gamma K_\varepsilon \gamma^{-1} = \gamma' K_\varepsilon (\gamma')^{-1}$.*
- (ii) *If $\ker \Lambda$ is trivial and if $(\gamma')^{-1}\gamma \notin H_{-\varepsilon}$, then $\gamma K_\varepsilon \gamma^{-1}$ and $\gamma' K_\varepsilon (\gamma')^{-1}$ have a trivial intersection and mutually commute.*
- (iii) *If $\ker \Lambda$ is trivial, then $\gamma K_\varepsilon \gamma^{-1}$ and $\gamma'' K_{-\varepsilon} (\gamma'')^{-1}$ have a trivial intersection and mutually commute.*

Proof. (i) $(\gamma')^{-1}\gamma K_\varepsilon \gamma^{-1}\gamma' = K_\varepsilon$ by Lemma 2.3.

(ii) If $(\gamma')^{-1}\gamma$ is an element of $G \setminus H_{-\varepsilon}$, then the assertion follows from Lemma 2.4 (ii). If $(\gamma')^{-1}\gamma$ starts with $\tau^{-\varepsilon}$ and ends with τ^ε , then, by Lemma 2.4 (i), it follows that

$$(\gamma')^{-1}\gamma K_\varepsilon \gamma^{-1}\gamma' < K_{-\varepsilon},$$

which, combined with $K_\varepsilon \cap K_{-\varepsilon} = \ker \Lambda = \{1\}$, proves the assertion.

(iii) Observe that the reduced form of $(\gamma'')^{-1}\gamma$ starts and ends with τ^ε , therefore the assertion follows from Lemma 2.4 (iii). \square

Assume that $\ker \Lambda = \{1\}$. Let S_ε be a left coset representatives of G/H_ε for $\varepsilon = \pm 1$.

It follows from Lemma 2.5 that, for two reduced words

$$\gamma = s_{n+1}\tau^{\varepsilon_n}s_n \cdots s_2\tau^{\varepsilon_1}s_1\tau^\varepsilon \text{ and } \gamma' = t_{n+1}\tau^{\varepsilon'_n}t_n \cdots t_2\tau^{\varepsilon'_1}t_1\tau^\varepsilon$$

with $s_i, t_i \in S_{-1} \cup S_1$ and $\varepsilon, \varepsilon_i, \varepsilon'_i \in \{-1, 1\}$,

$$\gamma K_\varepsilon \gamma^{-1} = \gamma' K_\varepsilon (\gamma')^{-1}$$

if and only if $\gamma = \gamma'$, and this happens if and only if $\varepsilon_i = \varepsilon'_i$ and $s_i = t_i$, $\forall i$. In the case $\gamma \neq \gamma'$, $\gamma K_\varepsilon \gamma^{-1}$ and $\gamma' K_\varepsilon (\gamma')^{-1}$ have a trivial intersection and mutually commute.

If $\gamma'' = r_{n+1}\tau^{\varepsilon''_n}r_n \cdots r_2\tau^{\varepsilon''_1}r_1\tau^{-\varepsilon}$ is another reduced word, where $r_i \in S_{-1} \cup S_1$ and $\varepsilon''_i \in \{-1, 1\}$, then $\gamma K_\varepsilon \gamma^{-1}$ and $\gamma'' K_{-\varepsilon} (\gamma'')^{-1}$ have a trivial intersection and mutually commute.

From these considerations, it follow that

$$\mathcal{K}(0) \equiv \bigoplus_{s \in S_{-1}} sK_1s^{-1} \oplus \bigoplus_{t \in S_1} tK_{-1}t^{-1} \quad (2)$$

and, for $n \geq 0$,

$$\mathcal{K}(n+1) \equiv \bigoplus_{\substack{\varepsilon=\pm 1 \\ s_i \in S_{-1} \cup S_1, \varepsilon_i=\pm 1 \\ s_{n+1}\tau^{\varepsilon_n} s_n \cdots s_2 \tau^{\varepsilon_1} s_1 \tau^\varepsilon \text{ reduced}}} s_{n+1}\tau^{\varepsilon_n} s_n \cdots s_2 \tau^{\varepsilon_1} s_1 \tau^\varepsilon K_\varepsilon \tau^{-\varepsilon} s_1^{-1} \tau^{-\varepsilon_1} s_2^{-1} \cdots s_n^{-1} \tau^{-\varepsilon_n} s_{n+1}^{-1} \quad (3)$$

are normal subgroups of G . Also, consider the groups

$$\mathcal{K}(0, \varepsilon) \equiv \bigoplus_{s \in S_{-\varepsilon}} s K_1 s^{-1} \oplus \bigoplus_{t \in S'_\varepsilon} t K_{-1} t^{-1},$$

which are normal in H_ε for $\varepsilon = \pm 1$.

Remark 2.6. *The group G acts transitively on the vertices $s\tau G$, where $s \in S_{-1}$. It also acts transitively on the vertices $s\tau^{-1}G$, where $s \in S_1$. This fact is an important ingredient in the examples below.*

Remark 2.7. *It follows from Lemma 2.4 that K_{-1} is isomorphic to a subgroup of K_1 and vice-versa. Consequently, $K_{-1} = \{1\}$ if and only if $K_1 = \{1\}$. In this situation, $\mathcal{K}(n) = \{1\} \forall n \geq 0$.*

2.2. A FAMILY OF EXAMPLES

For $\varepsilon = \pm 1$, consider nonempty sets I'_ε , and let $I_\varepsilon \equiv I'_\varepsilon \sqcup \{\iota_\varepsilon\}$. Also, let Σ_ε be transitive permutation groups on I_ε , and let $\Gamma = \Sigma_{-1} \cdot \Sigma_1$ be the corresponding permutation group on $I_{-1} \sqcup I_1$. Let $\Sigma'_\varepsilon \equiv (\Sigma_\varepsilon)_{\iota_\varepsilon}$ be the respective stabilizer groups, and define $\Gamma_\varepsilon \equiv \Gamma_{\iota_\varepsilon} = \Sigma'_\varepsilon \cdot \Sigma_{-\varepsilon}$. Define

$$\begin{aligned} \Lambda[\Sigma_{-1}, \Sigma_1] &\equiv \Lambda[I_{-1}, I_1, \iota_{-1}, \iota_1; \Sigma_{-1}, \Sigma_1] \\ &\equiv \text{HNN}(G, H, \theta) = \langle G, \tau \mid \tau^{-1}h\tau = \theta(h) \text{ for all } h \in H \rangle, \end{aligned}$$

where

$$\underline{H} \equiv \langle \{h(i_1, \varepsilon_1 \dots, i_n, \varepsilon_n; \sigma_n) \mid n \in \mathbb{N}, \varepsilon_t \in \{-1, 1\}, i_t \in I_{-\varepsilon_t}, \text{ and } \sigma_n \in \Gamma_{\varepsilon_n} \text{ satisfy } i_t \in I'_{-\varepsilon_t} \text{ whenever } \varepsilon_t \varepsilon_{t-1} = -1; \} \rangle \text{ and}$$

$$H_\varepsilon = \langle \underline{H} \cup \{h(\sigma_\varepsilon) \mid \sigma_\varepsilon \in \Gamma_\varepsilon\} \rangle, \varepsilon = \pm 1.$$

Finally, define

$$G = \langle H_{-1}, H_1 \rangle = \langle \underline{H} \cup \{h(\sigma) \mid \sigma \in \Gamma\} \rangle,$$

where the following relations hold (there are redundancies):

(R1) Elements $h(\sigma_{-1})$'s and $h(\sigma_1)$'s commute for all $\sigma_\varepsilon \in \Sigma_\varepsilon$, where $\varepsilon = \pm 1$.

(R2) Let $1 \leq m < n$, $\sigma_n \in \Gamma_{\varepsilon_n}$, and $\sigma'_m \in \Gamma_{\varepsilon_m}$. If $(i_1, \varepsilon_1 \dots, i_m, \varepsilon_m) \neq (j_1, \varepsilon_1 \dots, j_m, \varepsilon_m)$, the elements

$$h(j_1, \varepsilon_1 \dots, j_m, \varepsilon_m; \sigma'_m) \text{ and } h(i_1, \varepsilon_1 \dots, i_m, \varepsilon_m, \dots, i_n, \varepsilon_n; \sigma_n)$$

commute.

(R3) For $1 \leq m < n$ and $\sigma_t \in \Gamma_{\varepsilon_t}$, the following holds

$$h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m; \sigma_m) h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m, i_{m+1}, \varepsilon_{m+1}, \dots, i_n, \varepsilon_n; \sigma_n) h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m; \sigma_m)^{-1} \\ = h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m, \sigma_m(i_{m+1}), \varepsilon_{m+1}, \dots, i_n, \varepsilon_n; \sigma_n).$$

(R4) For $\sigma_m, \sigma'_m \in \Gamma_{\varepsilon_m}$, the following holds

$$h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m; \sigma_m) h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m; \sigma'_m) = h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m; \sigma_m \sigma'_m).$$

(R5) For $\sigma, \sigma' \in \Gamma$, the following holds

$$h(\sigma) h(\sigma') = h(\sigma \sigma').$$

(R6) For $n \in \mathbb{Z}$, $\sigma \in \Gamma$, and $\sigma_n \in \Gamma_{\varepsilon_n}$, the following holds

$$h(\sigma) h(i_1, \varepsilon_1, \dots, i_n, \varepsilon_n; \sigma_n) h(\sigma)^{-1} = h(\sigma(i_1), \varepsilon_1, i_2, \varepsilon_2, \dots, i_n, \varepsilon_n; \sigma_n).$$

(R7) For $\varepsilon = \pm 1$ and $\sigma_\varepsilon \in \Gamma_\varepsilon$, the following holds

$$\theta^{-\varepsilon}(h(\sigma_\varepsilon)) = (\tau^\varepsilon h(\sigma_\varepsilon) \tau^{-\varepsilon}) = h(\iota_{-\varepsilon}, \varepsilon; \sigma_\varepsilon).$$

(R8) For $\varepsilon = \pm 1$, $n \in \mathbb{N}$, and $\sigma_n \in \Gamma_{\varepsilon_n}$, the following holds

$$\theta^{-\varepsilon}(h(i_1, \varepsilon, i_2, \varepsilon_2, \dots, i_n, \varepsilon_n; \sigma_n)) = (\tau^\varepsilon h(i_1, \varepsilon, i_2, \varepsilon_2, \dots, i_n, \varepsilon_n; \sigma_n) \tau^{-\varepsilon}) \\ = h(\iota_{-\varepsilon}, \varepsilon, i_1, \varepsilon, i_2, \varepsilon_2, \dots, i_n, \varepsilon_n; \sigma_n).$$

(R9) For $\varepsilon = \pm 1$, $n \in \mathbb{N}$, and $\sigma_n \in \Gamma_{\varepsilon_n}$, the following holds

$$\theta^\varepsilon(h(i_1, \varepsilon, \dots, i_n, \varepsilon_n; \sigma_n)) = (\tau^{-\varepsilon} h(i_1, \varepsilon, \dots, i_n, \varepsilon_n; \sigma_n) \tau^\varepsilon) \\ = \begin{cases} h(i_2, \varepsilon_2, \dots, i_n, \varepsilon_n; \sigma_n), & \text{if } i_1 = \iota_{-\varepsilon}, \\ h(\iota_\varepsilon, -\varepsilon, i_1, \varepsilon, \dots, i_n, \varepsilon_n; \sigma_n), & \text{if } i_1 \neq \iota_{-\varepsilon}. \end{cases}$$

2.3. SOME BASIC PROPERTIES OF THE EXAMPLES AND THEIR QUASI-KERNELS

In this subsection we fix a group $\Lambda = \Lambda[I_{-1}, I_1, \iota_{-1}, \iota_1; \Sigma_{-1}, \Sigma_1]$.

First, let's note that $\text{Index}[G : H_\varepsilon] = \#(I_\varepsilon)$ for $\varepsilon = \pm 1$. To see this, recall that Σ_ε acts transitively on I_ε , and for $i \in I_\varepsilon$, choose $\mu_\varepsilon^i \in \Sigma_\varepsilon$ satisfying $\mu_\varepsilon^i(\iota_\varepsilon) = i$. Let's denote $\lambda_\varepsilon^i = h(\mu_\varepsilon^i)$. If $\sigma \in \Sigma_\varepsilon \setminus \Sigma'_\varepsilon$ satisfies $\sigma(\iota_\varepsilon) = i$, then $(\mu_\varepsilon^i)^{-1} \circ \sigma(\iota_\varepsilon) = \iota_\varepsilon$. Therefore $(\mu_\varepsilon^i)^{-1} \circ \sigma \in \Sigma'_\varepsilon$, so $h((\mu_\varepsilon^i)^{-1} \circ \sigma) \in H_\varepsilon$. It follows that $h(\sigma) \in h(\mu_\varepsilon^i) H_\varepsilon = \lambda_\varepsilon^i H_\varepsilon$. Consequently, for each $\varepsilon = \pm 1$,

$$G = H_\varepsilon \sqcup \bigsqcup_{i \in I'_\varepsilon} \lambda_\varepsilon^i H_\varepsilon. \quad (4)$$

It is easy to see in these notations that for $\varepsilon = \pm 1$, the set

$$S_\varepsilon = \{ \lambda_\varepsilon^i \mid i \in I'_\varepsilon \} \cup \{ 1 \}$$

is a left coset representative of H_ε in G .

Next, consider the action of Λ on its Bass-Serre tree $\Theta = \Theta[\Lambda]$. The set of all adjacent vertices to the vertex G is

$$\{ \tau G \} \cup \{ \lambda_{-1}^i \tau G \mid i \in I'_{-1} \} \cup \{ \tau^{-1} G \} \cup \{ \lambda_1^i \mid i \in I'_1 \}.$$

This set can be indexed by the set $I_{-1} \cup I_1$ in the obvious way: Denote by $v(\emptyset)$ the vertex G , by $v(\iota_{-1}, 1)$ the vertex τG , by $v(\iota_1, -1)$ the vertex $\tau^{-1} G$, by $v(i_{-1}, 1)$ the vertex $\lambda_{-1}^{i_{-1}} \tau G$, where $i_{-1} \in I'_{-1}$, and by $v(i_1, -1)$ the vertex $\lambda_1^{i_1} \tau^{-1} G$, where $i_1 \in I'_1$. Denote a general vertex

$$\lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \dots \lambda_{-\varepsilon_n}^{i_n} \tau^{\varepsilon_n} G$$

by $v(i_1, \varepsilon_1, \dots, i_n, \varepsilon_n)$ for an element $\lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \dots \lambda_{-\varepsilon_n}^{i_n} \tau^{\varepsilon_n} \in \Lambda$ in its normal form, i.e., $i_t \in I_{-\varepsilon_t}$ and if $\varepsilon_{t-1} \cdot \varepsilon_t = -1$, then $i_t \in I'_{-\varepsilon_t}$.

With the notation of Remark 2.1, for $\varepsilon = \pm 1$, Θ_ε is the full subtree of Θ containing the vertex $v(\emptyset) = G$ and vertices $v(i_1, \varepsilon_1, \dots, i_n, \varepsilon_n)$, where $n \geq 1$ and $(i_1, \varepsilon_1) \neq (\iota_{-\varepsilon}, \varepsilon)$, and $\bar{\Theta}_\varepsilon$ is the full subtree of Θ containing the vertices $v(\iota_{-\varepsilon}, \varepsilon, i_1, \varepsilon_1, \dots, i_n, \varepsilon_n)$, where $n \geq 0$.

Remark 2.8. *It follows from [1, Exercise VI.3] that our examples are never finitely presented since H is never finitely generated.*

We continue with

Lemma 2.9. (i) *Let $m \geq 1$, $\sigma_m \in \Gamma_{\varepsilon_m}$, $i_t \in I_{-\varepsilon_t}$, and $\varepsilon \in \{-1, 1\}$ satisfy $\varepsilon_t \varepsilon_{t-1} = -1 \Rightarrow i_t \in I'_{-\varepsilon_t}$. Then*

$$\begin{aligned} h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m; \sigma_m) \\ = \lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \dots \lambda_{-\varepsilon_m}^{i_m} \tau^{\varepsilon_m} h(\sigma_m) \tau^{-\varepsilon_m} (\lambda_{-\varepsilon_m}^{i_m})^{-1} \dots \tau^{-\varepsilon_1} (\lambda_{-\varepsilon_1}^{i_1})^{-1}. \end{aligned}$$

(ii) *Every element h of G can be written as*

$$h = h(\sigma) \prod_{k=1}^m h(i_1^k, \varepsilon_{k,1}, \dots, i_{n_k}^k, \varepsilon_{k,n_k}; \sigma_k),$$

where $m \geq 1$, $\sigma_k \in \Gamma_{\varepsilon_{k,n_k}}$, $1 \leq n_1 \leq \dots \leq n_m$, and $\sigma \in \Gamma$ satisfy the condition: if $n_k = n_{k+a}$ for some $1 \leq k \leq m$ and some $a \geq 1$, then

$$(i_1^k, \varepsilon_{k,1}, \dots, i_{n_k}^k, \varepsilon_{k,n_k}) \neq (i_1^{k+a}, \varepsilon_{k+a,1}, \dots, i_{n_{k+a}}^{k+a}, \varepsilon_{k+a,n_{k+a}}).$$

(ii) Every element $g \in T_\varepsilon$ can be written as

$$g = \lambda_{-\varepsilon}^i \tau^\varepsilon \lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \dots \lambda_{-\varepsilon_m}^{i_m} \tau^{\varepsilon_m} h,$$

where $h \in G$ and $m \geq 0$.

Proof. (i) follows by repeated applications of relations (R7), (R8), and (R6).

(ii) follows by repeated applications of relations (R3) and (R6).

(iii) follows by equation (4) and the structure of HNN-extensions. \square

Lemma 2.10. Let $n > m \geq 1$ and $\sigma_k \in \Gamma_{\varepsilon_k}$. Then the following hold

$$(i) \quad h(i_1, \varepsilon_1 \dots, i_m, \varepsilon_m; \sigma_m) v(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m, i_{m+1}, \varepsilon_{m+1}, \dots, i_n, \varepsilon_n) \\ = v(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m, \sigma_m(i_{m+1}), \varepsilon_{m+1}, \dots, i_n, \varepsilon_n).$$

$$(ii) \quad h(i_1, \varepsilon_1 \dots, i_m, \varepsilon_m; \sigma_m) \in \Lambda_{v(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m)} \text{ and } h(\sigma) \in \Lambda_{v(\emptyset)} \text{ for } \sigma \in \Gamma.$$

$$(iii) \quad \text{If } \sigma_\varepsilon \in \Gamma_\varepsilon, \text{ then } h(\sigma_\varepsilon) \in \Lambda_{(\bar{\Theta}_{-\varepsilon})} = \tau^{-\varepsilon} K_\varepsilon \tau^\varepsilon.$$

(iv) Let $m \leq n$ and let $h(i_1, \varepsilon_1 \dots, i_n, \varepsilon_n; \sigma_n), h(j_1, e_1 \dots, j_m, e_m; \delta_m) \in \Lambda$. If $(i_1, \varepsilon_1 \dots, i_m, \varepsilon_m) \neq (j_1, e_1 \dots, j_m, e_m)$, then $h(i_1, \varepsilon_1 \dots, i_n, \varepsilon_n; \sigma_n) \in \Lambda_{v(j_1, e_1, \dots, j_m, e_m)}$ and $h(j_1, e_1 \dots, j_m, e_m; \delta_m) \in \Lambda_{v(i_1, \varepsilon_1, \dots, i_n, \varepsilon_n)}$.

$$(iv) \quad h(i_1, \varepsilon_1 \dots, i_n, \varepsilon_n; \sigma_n) \in \Lambda_{(\bar{\Theta}_\varepsilon)} \iff (i_1, \varepsilon_1) \neq (\iota_{-\varepsilon}, \varepsilon).$$

Proof. (i) First, note that

$$\sigma \equiv (\lambda_{-\varepsilon_{m+1}}^{\sigma_m(i_{m+1})})^{-1} \circ h(\sigma_m) \lambda_{-\varepsilon_{m+1}}^{i_{m+1}} \in \Gamma_{-\varepsilon_{m+1}}$$

since it fixes $\iota_{-\varepsilon_{m+1}}$. It follows by Lemma 2.9 (i) and (iii) that there are $k_t \in I_{\varepsilon_t}$ and a $\chi \in H_{\varepsilon_n}$ that satisfy $(\tau^{\varepsilon_{m+1}} \dots \lambda_{-\varepsilon_n}^{i_n} \tau^{\varepsilon_n})^{-1} = \chi \tau^{-\varepsilon_n} \lambda_{\varepsilon_{n-1}}^{k_{n-1}} \dots \lambda_{\varepsilon_{m+1}}^{k_{m+1}} \tau^{-\varepsilon_{m+1}}$. Therefore

$$(\tau^{\varepsilon_{m+1}} \dots \lambda_{-\varepsilon_n}^{i_n} \tau^{\varepsilon_n})^{-1} h(\sigma) \tau^{\varepsilon_{m+1}} \dots \lambda_{-\varepsilon_n}^{i_n} \tau^{\varepsilon_n} \\ = \chi \tau^{-\varepsilon_n} \lambda_{\varepsilon_{n-1}}^{k_{n-1}} \dots \lambda_{\varepsilon_{m+1}}^{k_{m+1}} \tau^{-\varepsilon_{m+1}} h(\sigma) \tau^{\varepsilon_{m+1}} (\lambda_{\varepsilon_{m+1}}^{k_{m+1}})^{-1} \dots (\lambda_{\varepsilon_{n-1}}^{k_{n-1}})^{-1} \tau^{\varepsilon_n} \chi^{-1} \\ = \chi h(\iota_{\varepsilon_n}, -\varepsilon_n, k_{n-1}, -\varepsilon_{n-1}, \dots, k_{m+2}, -\varepsilon_{m+2}, \iota_{\varepsilon_{m+1}}, -\varepsilon_{m+1}; \sigma) \chi^{-1}.$$

Then Lemma 2.9 (i) implies

$$h(i_1, \varepsilon_1 \dots, i_m, \varepsilon_m; \sigma_m) v(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m, i_{m+1}, \varepsilon_{m+1}, \dots, i_n, \varepsilon_n) \\ = \lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \dots \lambda_{-\varepsilon_m}^{i_m} \tau^{\varepsilon_m} h(\sigma_m) \tau^{-\varepsilon_m} (\lambda_{-\varepsilon_m}^{i_m})^{-1} \dots \tau^{-\varepsilon_1} (\lambda_{-\varepsilon_1}^{i_1})^{-1} \cdot \lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \dots \lambda_{-\varepsilon_n}^{i_n} \tau^{\varepsilon_n} G \\ = \lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \dots \lambda_{-\varepsilon_m}^{i_m} \tau^{\varepsilon_m} h(\sigma_m) \lambda_{-\varepsilon_{m+1}}^{i_{m+1}} \tau^{\varepsilon_{m+1}} \dots \lambda_{-\varepsilon_n}^{i_n} \tau^{\varepsilon_n} G \\ = \lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \dots \lambda_{-\varepsilon_m}^{i_m} \tau^{\varepsilon_m} \lambda_{-\varepsilon_{m+1}}^{\sigma_m(i_{m+1})} h(\sigma) \tau^{\varepsilon_{m+1}} \dots \lambda_{-\varepsilon_n}^{i_n} \tau^{\varepsilon_n} G$$

$$\begin{aligned}
&= \lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \dots \lambda_{-\varepsilon_m}^{i_m} \tau^{\varepsilon_m} \lambda_{-\varepsilon_{m+1}}^{\sigma_m(i_{m+1})} \tau^{\varepsilon_{m+1}} \dots \lambda_{-\varepsilon_n}^{i_n} \tau^{\varepsilon_n} \\
&\quad \cdot (\tau^{\varepsilon_{m+1}} \dots \lambda_{-\varepsilon_n}^{i_n} \tau^{\varepsilon_n})^{-1} h(\sigma) \tau^{\varepsilon_{m+1}} \dots \lambda_{-\varepsilon_n}^{i_n} \tau^{\varepsilon_n} G \\
&= \lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \dots \lambda_{-\varepsilon_m}^{i_m} \tau^{\varepsilon_m} \lambda_{-\varepsilon_{m+1}}^{\sigma_m(i_{m+1})} \tau^{\varepsilon_{m+1}} \dots \lambda_{-\varepsilon_n}^{i_n} \tau^{\varepsilon_n} \\
&\quad \cdot \chi h(\iota_{\varepsilon_n}, -\varepsilon_n, k_{n-1}, -\varepsilon_{n-1}, \dots, k_{m+2}, -\varepsilon_{m+2}, \iota_{\varepsilon_{m+1}}, -\varepsilon_{m+1}; \sigma) \chi^{-1} G \\
&= \lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \dots \lambda_{-\varepsilon_m}^{i_m} \tau^{\varepsilon_m} \lambda_{-\varepsilon_{m+1}}^{\sigma_m(i_{m+1})} \tau^{\varepsilon_{m+1}} \dots \lambda_{-\varepsilon_n}^{i_n} \tau^{\varepsilon_n} G \\
&= v(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m, \sigma_m(i_{m+1}), \varepsilon_{m+1}, \dots, i_n, \varepsilon_n).
\end{aligned}$$

(ii) The second claim is obvious. For the first claim,

$$\begin{aligned}
&h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m; \sigma_m) v(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m) \\
&= \lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \dots \lambda_{-\varepsilon_m}^{i_m} \tau^{\varepsilon_m} h(\sigma_m) \tau^{-\varepsilon_m} (\lambda_{-\varepsilon_m}^{i_m})^{-1} \dots \tau^{-\varepsilon_1} (\lambda_{-\varepsilon_1}^{i_1})^{-1} \cdot \lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \dots \lambda_{-\varepsilon_n}^{i_n} \tau^{\varepsilon_n} G \\
&= \lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \dots \lambda_{-\varepsilon_m}^{i_m} \tau^{\varepsilon_m} h(\sigma_m) G = v(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m).
\end{aligned}$$

(iii) The fact $\Lambda_{(\bar{\Theta}_{-\varepsilon})} = \tau^{-\varepsilon} K_\varepsilon \tau^\varepsilon$ is stated in Proposition 2.2. Let $n \geq 0$ and let $v(\iota_\varepsilon, -\varepsilon, i_1, \varepsilon_1, \dots, i_n, \varepsilon_n) \in \bar{\Theta}_{-\varepsilon}$. By the argument at the beginning of the proof of (i), there are $k_t \in I_{\varepsilon_t}$ and a $\chi \in H_{\varepsilon_n}$ satisfying

$$\begin{aligned}
&(\tau^{-\varepsilon} \lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \dots \lambda_{-\varepsilon_n}^{i_n} \tau^{\varepsilon_n})^{-1} h(\sigma_\varepsilon) \tau^{-\varepsilon} \lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \dots \lambda_{-\varepsilon_n}^{i_n} \tau^{\varepsilon_n} \\
&= \chi h(\iota_{\varepsilon_n}, -\varepsilon_n, k_{n-1}, -\varepsilon_{n-1}, \dots, i_{\varepsilon_1}, -\varepsilon_1, \varepsilon, \iota_{-\varepsilon}; \sigma_\varepsilon) \chi^{-1}.
\end{aligned}$$

Therefore

$$\begin{aligned}
&h(\sigma_\varepsilon) v(\iota_\varepsilon, -\varepsilon, i_1, \varepsilon_1, \dots, i_n, \varepsilon_n) = h(\sigma_\varepsilon) \tau^{-\varepsilon} \lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \dots \lambda_{-\varepsilon_m}^{i_m} \tau^{\varepsilon_m} G \\
&= \tau^{-\varepsilon} \lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \dots \lambda_{-\varepsilon_m}^{i_m} \tau^{\varepsilon_m} \cdot (\tau^{-\varepsilon} \lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \dots \lambda_{-\varepsilon_m}^{i_m} \tau^{\varepsilon_m})^{-1} h(\sigma_\varepsilon) \tau^{-\varepsilon} \lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \dots \lambda_{-\varepsilon_m}^{i_m} \tau^{\varepsilon_m} G \\
&= \tau^{-\varepsilon} \lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \dots \lambda_{-\varepsilon_m}^{i_m} \tau^{\varepsilon_m} \cdot \chi h(\iota_{\varepsilon_n}, -\varepsilon_n, k_{n-1}, -\varepsilon_{n-1}, \dots, i_{\varepsilon_1}, -\varepsilon_1, \varepsilon, \iota_{-\varepsilon}; \sigma_\varepsilon) \chi^{-1} G \\
&= v(\iota_\varepsilon, -\varepsilon, i_1, \varepsilon_1, \dots, i_n, \varepsilon_n).
\end{aligned}$$

Consequently $h(\sigma_\varepsilon) \in \bar{\Theta}_{-\varepsilon}$.

(iv) Note that the element $\gamma = \tau^{-e_m} (\lambda_{-e_m}^{j_m})^{-1} \dots \tau^{-e_1} (\lambda_{-e_1}^{j_1})^{-1} \lambda_{-e_1}^{i_1} \tau^{\varepsilon_1} \dots \lambda_{-\varepsilon_n}^{i_n} \tau^{\varepsilon_n}$ belongs to $T_{-e_m}^\dagger$ because of the condition $(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m) \neq (j_1, e_1, \dots, j_m, e_m)$. It follows from Lemma 2.9 (iii) that $\gamma = \tau^{-e_m} \lambda_{-l_1}^{k_1} \tau^{l_1} \lambda_{-l_2}^{k_2} \tau^{l_2} \dots \lambda_{-l_s}^{k_s} \tau^{l_s} h$, where $h \in G$ and where $k_t \in I_{-l_t}$, $\forall t$. Then

$$\begin{aligned}
&h(j_1, e_1, \dots, j_m, e_m; \delta_m) \in \Lambda_{v(i_1, \varepsilon_1, \dots, i_n, \varepsilon_n)} \\
&\iff \lambda_{-e_1}^{j_1} \tau^{e_1} \dots \lambda_{-e_m}^{j_m} \tau^{e_m} h(\delta_m) \tau^{-e_m} (\lambda_{-e_m}^{j_m})^{-1} \dots \tau^{-e_1} (\lambda_{-e_1}^{j_1})^{-1} \in \Lambda_{v(i_1, \varepsilon_1, \dots, i_n, \varepsilon_n)} \\
&\iff h(\delta_m) \in \tau^{-e_m} (\lambda_{-e_m}^{j_m})^{-1} \dots \tau^{-e_1} (\lambda_{-e_1}^{j_1})^{-1} \Lambda_{v(i_1, \varepsilon_1, \dots, i_n, \varepsilon_n)} \lambda_{-e_1}^{j_1} \tau^{e_1} \dots \lambda_{-e_m}^{j_m} \tau^{e_m} \\
&\iff h(\delta_m) \in \Lambda_{\tau^{-e_m} (\lambda_{-e_m}^{j_m})^{-1} \dots \tau^{-e_1} (\lambda_{-e_1}^{j_1})^{-1} v(i_1, \varepsilon_1, \dots, i_n, \varepsilon_n)} \\
&\iff h(\delta_m) \in \Lambda_{\tau^{-e_m} (\lambda_{-e_m}^{j_m})^{-1} \dots \tau^{-e_1} (\lambda_{-e_1}^{j_1})^{-1} \lambda_{-e_1}^{i_1} \tau^{\varepsilon_1} \dots \lambda_{-\varepsilon_n}^{i_n} \tau^{\varepsilon_n} G} \\
&\iff h(\delta_m) \in \Lambda_{\tau^{-e_m} \lambda_{-l_1}^{k_1} \tau^{l_1} \lambda_{-l_2}^{k_2} \tau^{l_2} \dots \lambda_{-l_s}^{k_s} \tau^{l_s} h G} \\
&\iff h(\delta_m) \in \Lambda_{v(\iota_{e_m}, -e_m, k_1, l_1, \dots, k_s, l_s)}.
\end{aligned}$$

The last equivalence holds according to (iii). The inclusion $h(i_1, \varepsilon_1, \dots, i_n, \varepsilon_n; \sigma_n) \in \Lambda_{v(j_1, e_1, \dots, j_m, e_m)}$ is proven analogously.

(v) Every vertex of $\Lambda_{(\bar{\Theta}_\varepsilon)}$ is of the form $v(\iota_{-\varepsilon}, \varepsilon, j_1, e_1, \dots, j_m, e_m)$, so if tuples $(i_1, \varepsilon_1, \dots, i_n, \varepsilon_n)$ and $(\iota_{-\varepsilon}, \varepsilon, j_1, e_1, \dots, j_m, e_m)$ satisfy the assumptions of (iv), then $h(i_1, \varepsilon_1, \dots, i_n, \varepsilon_n; \sigma_n) \in \Lambda_{(\bar{\Theta}_\varepsilon)}$. By (i), $h(\iota_{-\varepsilon}, \varepsilon, j_1, e_1, \dots, j_m, e_m; \sigma_m) \notin \Lambda_{(\bar{\Theta}_\varepsilon)}$, and the statement follows. \square

Proposition 2.11. *For a group $\Lambda = \Lambda[I_{-1}, I_1, \iota_{-1}, \iota_1; \Sigma_{-1}, \Sigma_1]$ and for $\varepsilon = \pm 1$, the following hold*

- (i) $\Lambda_{(\bar{\Theta}_\varepsilon)} = \langle \{ h(\sigma_{-\varepsilon}) \mid \sigma_{-\varepsilon} \in \Gamma_{-\varepsilon} \} \cup \{ h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m; \sigma_m) \mid m \geq 1, h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m; \sigma_m) \in H_{-\varepsilon}, \text{ and } (i_1, \varepsilon_1) \neq (\iota_{-\varepsilon}, \varepsilon) \} \rangle$;
- (ii) $|K_\varepsilon| = \langle \{ h(\iota_\varepsilon, -\varepsilon; \sigma_{-\varepsilon}) \mid \sigma_{-\varepsilon} \in \Gamma_{-\varepsilon} \} \sqcup \{ h(\iota_\varepsilon, -\varepsilon, i_1, \varepsilon_1, \dots, i_n, \varepsilon_n; \sigma_n) \mid n \geq 1, \sigma_n \in \Gamma_{\varepsilon_n} \} \rangle$;
- (iii) $\ker \Lambda = \{1\}$.

Proof. (i) Denote the group on the right-hand-side by Δ . The inclusion $\Delta < \Lambda_{(\bar{\Theta}_\varepsilon)}$ follows from Lemma 2.10 (iii) and (v). Take an element $h \in \Lambda_{(\bar{\Theta}_\varepsilon)}$. Proposition 2.2 (iv) implies that $h \in H_{-\varepsilon}$. If we assume $h = h(\sigma)$, then $\sigma \in \Gamma_{-\varepsilon}$, and therefore $h(\sigma) \in \Delta$. If h is not of the form $h(\sigma)$, Lemma 2.9 (ii) can be applied to $h^{-1} \in H_{-\varepsilon}$. It follows that

$$h = \prod_{k=1}^m h(i_1^k, \varepsilon_{k,1}, \dots, i_{n_k}^k, \varepsilon_{k,n_k}; \sigma_k) \cdot h(\sigma_{-\varepsilon}),$$

where $m \geq 0$, $\sigma_k \in \Gamma_{\varepsilon_{k,n_k}}$, $n_1 \geq n_2 \geq \dots \geq n_m \geq 1$, and $\sigma_{-\varepsilon} \in \Gamma_{-\varepsilon}$. Assume $h(i_1^l, \varepsilon_{l,1}, \dots, i_{n_l}^l, \varepsilon_{l,n_l}; \sigma_l) \notin \Delta$ for some $1 \leq l \leq m$ and that l is the biggest number with this property. We will derive a contradiction below. Then it is clear that $i_1^l = \iota_{-\varepsilon}$ and $\varepsilon_{l,1} = \varepsilon$. Also, $\sigma_l \in \Gamma_{\varepsilon_{l,n_l}}$ is not the identity, so there exist two different elements $\kappa, \rho \in I_{-1} \sqcup I_1$, such that $\sigma_l(\kappa) = \rho$. Let h act on

$$v = v(i_1^l, \varepsilon_{l,1}, \dots, i_{n_l}^l, \varepsilon_{l,n_l}, \kappa, \varepsilon_{l,n_l}, \alpha_1, e_1, \dots, \alpha_{n_1}, e_{n_1}),$$

where α 's and e 's are arbitrary and allowed. The terms $h(\sigma_{-\varepsilon})$ and $\prod_{k=l+1}^m h(i_1^k, \varepsilon_{k,1}, \dots, i_{n_k}^k, \varepsilon_{k,n_k}; \sigma_k)$ leave v fixed by the choice of l . From the final condition of Lemma 2.9 (ii) and from Lemma 2.10 (iv), it follows that the terms with length equal to n_l also leave v fixed. Finally, from Lemma 2.10 (i), it follows that the remaining terms act on v by eventually changing only the α 's. Therefore we conclude that

$$\begin{aligned} hv(i_1^l, \varepsilon_{l,1}, \dots, i_{n_l}^l, \varepsilon_{l,n_l}, \kappa, \varepsilon_{l,n_l}, \alpha_1, e_1, \dots, \alpha_{n_1}, e_{n_1}) \\ = v(i_1^l, \varepsilon_{l,1}, \dots, i_{n_l}^l, \varepsilon_{l,n_l}, \rho, \varepsilon_{l,n_l}, \beta_1, e_1, \dots, \beta_{n_1}, e_{n_1}) \end{aligned}$$

for some β 's. This shows that $h \notin \Lambda_{(\bar{\Theta}_\varepsilon)}$, a contradiction that proves (i).

(ii) From Proposition 2.2 (iii), it follows that

$$K_\varepsilon = \tau^{-\varepsilon} K_\varepsilon (\tau^{-\varepsilon}) \tau^\varepsilon = \tau^{-\varepsilon} \Lambda_{(\bar{\Theta}_\varepsilon)} \tau^\varepsilon = \theta^\varepsilon (\Lambda_{(\bar{\Theta}_\varepsilon)}).$$

The assertion follows from relation (R7) and Lemma 2.9 (i).

(iii) is obvious. □

Now, we want to explore the structure of the quasi-kernels of $\Lambda = \Lambda[I_{-1}, I_1, \iota_{-1}, \iota_1; \Sigma_{-1}, \Sigma_1]$, in particular, that of $\Lambda_{(\bar{\Theta}_\varepsilon)}$.

First, we note that Proposition 2.11 (ii) and relation (R6) imply that for $i \in I_\varepsilon$,

$$\lambda_\varepsilon^i \tau^{-\varepsilon} \Lambda_{(\bar{\Theta}_\varepsilon)} \tau^\varepsilon (\lambda_\varepsilon^i)^{-1} = \lambda_\varepsilon^i K_\varepsilon (\lambda_\varepsilon^i)^{-1}$$

$$= \langle \{h(i, -\varepsilon, i_1, \varepsilon_1, \dots, i_m, \varepsilon_m; \sigma_m) \mid m \geq 0, h(i, -\varepsilon, i_1, \varepsilon_1, \dots, i_m, \varepsilon_m; \sigma_m) \in \underline{H}\} \rangle.$$

It is clear that

$$\Lambda_{(\bar{\Theta}_\varepsilon)} = \langle \{h(\sigma_{-\varepsilon}) \mid \sigma_{-\varepsilon} \in \Gamma_{-\varepsilon}\} \cup \bigcup_{i \in I_\varepsilon} \lambda_\varepsilon^i \tau^{-\varepsilon} \Lambda_{(\bar{\Theta}_\varepsilon)} \tau^\varepsilon (\lambda_\varepsilon^i)^{-1} \cup \bigcup_{i \in I_{-\varepsilon}} \lambda_{-\varepsilon}^i \tau^\varepsilon \Lambda_{(\bar{\Theta}_\varepsilon)} \tau^{-\varepsilon} (\lambda_{-\varepsilon}^i)^{-1} \rangle$$

$$= \langle \{h(\sigma_{-\varepsilon}) \mid \sigma_{-\varepsilon} \in \Gamma_{-\varepsilon}\} \cup \mathcal{K}(0, -\varepsilon) \rangle.$$

In other words,

$$\Lambda_{(\bar{\Theta}_\varepsilon)} \cong \mathcal{K}(0, -\varepsilon) \rtimes \Gamma_{-\varepsilon}.$$

This can be written "recursively" as

$$K_\varepsilon \cong [\bigoplus_{\#(S'_{-\varepsilon})} K_{-\varepsilon} \oplus \bigoplus_{\#(S_\varepsilon)} K_\varepsilon] \rtimes \Gamma_{-\varepsilon}. \quad (5)$$

This is in a sense a "wreath product" representation.

Let's denote

$$\mathcal{H}_\varepsilon(0) = \langle \{h(\sigma_{-\varepsilon}) \mid \sigma_{-\varepsilon} \in \Gamma_{-\varepsilon}\} \rangle.$$

For $n \geq 1$, let

$$\mathcal{H}_\varepsilon(n) = \langle \{h(i_1, \varepsilon_1, \dots, i_n, \varepsilon_n; \sigma_n) \mid h(i_1, \varepsilon_1, \dots, i_n, \varepsilon_n; \sigma_n) \in H_{-\varepsilon} \text{ and } (i_1, \varepsilon_1) \neq (\iota_{-\varepsilon}, \varepsilon)\} \rangle.$$

Note that, each $\mathcal{H}_\varepsilon(n)$ is isomorphic to a direct sum of copies of Γ_1 and Γ_{-1} . Let us also denote

$$\mathcal{H}_\varepsilon[n] = \langle \mathcal{H}_\varepsilon(0) \cup \mathcal{H}_\varepsilon(1) \cup \dots \cup \mathcal{H}_\varepsilon(n) \rangle.$$

Relation (R3) implies that $\mathcal{H}_\varepsilon(n) \triangleleft \mathcal{H}_\varepsilon[n]$ and that there is an extension

$$\{1\} \longrightarrow \mathcal{H}_\varepsilon(n) \longrightarrow \mathcal{H}_\varepsilon[n] \longrightarrow \mathcal{H}_\varepsilon[n-1] \longrightarrow \{1\}. \quad (6)$$

The natural embeddings $\mathcal{H}_\varepsilon[m] \hookrightarrow \mathcal{H}_\varepsilon[n]$ give a representation of $\Lambda_{(\bar{\Theta}_\varepsilon)}$ as a direct limit of groups

$$\Lambda_{(\bar{\Theta}_\varepsilon)} = \varinjlim_n \mathcal{H}_\varepsilon[n]. \quad (7)$$

Lemma 2.12. K_{-1} is amenable if and only if K_1 is amenable, if and only if Γ_{-1} and Γ_1 are both amenable, and if and only if Σ_{-1} and Σ_1 are both amenable.

Proof. Assume that Γ_ε is not amenable for some $\varepsilon = \pm 1$. Then, by equation (5), it follows that $K_{-\varepsilon}$ is not amenable, so equation (5), applied once more, gives the nonamenability of K_ε .

Conversely, assume that Γ_{-1} and Γ_1 are both amenable. Then $\mathcal{H}_\varepsilon(n)$ is amenable as a direct sum of copies of Γ_{-1} and Γ_1 . Also, $\mathcal{H}_\varepsilon[0] = \mathcal{H}_\varepsilon(0) \cong \Gamma_{-\varepsilon}$ is amenable for $\varepsilon = \pm 1$. Therefore an easy induction based on the extension (6), gives the amenability of $\mathcal{H}_\varepsilon[n]$ for each $\varepsilon = \pm 1$ and each $n \geq 0$. Finally, the direct limit representation (7) of $\Lambda_{(\bar{\Theta}_\varepsilon)}$ implies the amenability of $\Lambda_{(\bar{\Theta}_\varepsilon)}$ for and therefore that of $K_\varepsilon = \tau^{-\varepsilon} \Lambda_{(\bar{\Theta}_\varepsilon)} \tau^\varepsilon$ for $\varepsilon = \pm 1$. \square

2.4. GROUP-THEORETIC STRUCTURE

We give a result about the structure of our groups.

Theorem 2.13. Take $\Lambda = \Lambda[I_{-1}, I_1, \iota_{-1}, \iota_1; \Sigma_{-1}, \Sigma_1]$. Let's assume that:

- (i) Σ_{-1} and Σ_1 are 2-transitive, that is, all stabilizers $(\Sigma_\varepsilon)_{i_\varepsilon}$ are transitive on the sets $I_\varepsilon \setminus \{i_\varepsilon\}$ for all $i_\varepsilon \in I_\varepsilon$ and $\varepsilon = \pm 1$;
- (ii) For each $\varepsilon = \pm 1$, either $\Sigma_\varepsilon = \langle (\Sigma_\varepsilon)_{i_\varepsilon} \mid i_\varepsilon \in I_\varepsilon \rangle$ or $\Sigma_\varepsilon = \text{Sym}(2)$.

Then Λ has a simple normal subgroup Ξ for which there is a group extension

$$1 \longrightarrow \Xi \longrightarrow \Lambda \xrightarrow{\eta} (\Gamma/[\Gamma, \Gamma]) \wr_{\mathbb{Z}} \mathbb{Z} \longrightarrow 1,$$

where η is defined on the generators by

$$\eta(h(\sigma)) = ((\dots, 0, \dots, 0, ([\sigma], 0), 0, \dots, 0, \dots), 0), \quad \eta(\tau) = ((\dots, 0, \dots), 1), \quad \text{and}$$

$$\eta(h(i_1, \varepsilon_1, \dots, i_n, \varepsilon_n; \sigma_n)) = ((\dots, 0, \dots, 0, ([\sigma_n], \varepsilon_1 + \dots + \varepsilon_n), 0, \dots, 0, \dots), 0).$$

Here $[\sigma]$ denotes the image of the permutation $\sigma \in \Gamma$ in $\Gamma/[\Gamma, \Gamma]$.

Proof. It follows from relations (R7), (R8), and (R9) that the action of θ on an element $h(i_1, \varepsilon_1, \dots, i_n, \varepsilon_n; \sigma_n)$ is consistent with the definition of η and the multiplication in the wreath product, that is,

$$\begin{aligned} \eta(\theta(h(i_1, \varepsilon_1, \dots, i_n, \varepsilon_n; \sigma_n))) &= \eta(\tau^{-1} h(i_1, \varepsilon_1, \dots, i_n, \varepsilon_n; \sigma_n) \tau) \\ &= ((\dots, 0, \dots, 0, ([\sigma_n], \varepsilon_1 + \dots + \varepsilon_n - 1), 0, \dots, 0, \dots), 0). \end{aligned}$$

It is easy to see that, since the commutant is in the kernel, the homomorphism $\eta : G \rightarrow (\Gamma/[\Gamma, \Gamma]) \wr_{\mathbb{Z}} \mathbb{Z}$ is well defined by

$$\eta(g) = ((\dots, (\prod_{\varepsilon_1 + \dots + \varepsilon_n = m} [\sigma_n], m), \dots), 0),$$

where the products are taken over all the factors $h(i_1, \varepsilon_1, \dots, i_n, \varepsilon_n; \sigma_n)$ of g . These two observations together with the universal property of the HNN-extensions (Remark 1.1) enable us to extend η to the entire group Λ .

Now, notice that if $\lambda = g_1 \tau^{\varepsilon_1} g_2 \tau^{\varepsilon_2} g_3 \tau^{\varepsilon_3} \dots g_n \tau^{\varepsilon_n} g_{n+1} \in \Xi$, then $\varepsilon_1 + \dots + \varepsilon_n = 0$. Thus

$$\lambda = g_1 (\tau^{\varepsilon_1} g_2 \tau^{-\varepsilon_1}) (\tau^{\varepsilon_1 + \varepsilon_2} g_3 \tau^{-\varepsilon_1 - \varepsilon_2}) \dots (\tau^{\varepsilon_1 + \varepsilon_2 + \dots + \varepsilon_{n-1}} g_n \tau^{-\varepsilon_1 - \varepsilon_2 - \dots - \varepsilon_{n-1}}) g_{n+1}$$

can be represented as products of τ -conjugates of elements from G .

Using Lemma 2.9 (ii), we see that every $\lambda = \tau^n g \tau^{-n}$ can be written as a product of elements of the form $\tau^n h(\sigma) \tau^{-n}$ and $\tau^n h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m; \sigma_m) \tau^{-n}$. The second element equals either $\tau^{n-m} h(\sigma_m) \tau^{m-n}$ or $h(j_1, \varepsilon'_1, \dots, j_k, \varepsilon'_k; \sigma_m)$ for some j_p 's and ε'_p 's. Therefore, it is easy to see that Ξ is generated by the following set

$$\begin{aligned} & \{h(i_1, \varepsilon_1, \dots, i_n, \varepsilon_n; \sigma_n) h(i'_1, \varepsilon_1, \dots, i'_n, \varepsilon_n; \sigma_n^{-1}) \mid \varepsilon_k = \pm 1, i_k, i'_k \in I_{-\varepsilon_k}, \forall k; n \geq 2, \sigma_n \in \Gamma_{\varepsilon_n}\} \\ & \cup \{i, \varepsilon, i_0, -\varepsilon, i_1, \varepsilon, i_2, \varepsilon_2, \dots, i_n, \varepsilon_n; \sigma_n\} h(\bar{i}, \varepsilon, i'_2, \varepsilon_2, \dots, i_n, \varepsilon_n; \sigma_n^{-1}) \mid \\ & \quad n \geq 2, i_0 \in I_\varepsilon, i'_2 \in I_{-\varepsilon_2}, i, \bar{i} \in I_{-\varepsilon}; i_k \in I_{-\varepsilon_k}, \varepsilon, \varepsilon_k = \pm 1, \forall k\} \\ & \cup \{h(\sigma_\varepsilon) h(i_\varepsilon, -\varepsilon, i_{-\varepsilon}, \varepsilon; \sigma_\varepsilon^{-1}) \mid \sigma_\varepsilon \in \Gamma_\varepsilon, i_{-\varepsilon} \in I'_\varepsilon, i_\varepsilon \in I_{-\varepsilon}, \varepsilon = \pm 1\} \\ & \cup \{h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m, i, \varepsilon, j, -\varepsilon, j_1, \varepsilon'_1, \dots, j_n, \varepsilon'_n; \sigma) \\ & \quad \cdot h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m, j', -\varepsilon, i', \varepsilon, j_1, \varepsilon'_1, \dots, j_n, \varepsilon'_n; \sigma^{-1}) \mid \\ & \quad m, n \in \mathbb{N}_0, i, i' \in I_{-\varepsilon}, j, j' \in I_\varepsilon, \sigma \in \Gamma_{\varepsilon'_n}; \varepsilon, \varepsilon_k, \varepsilon'_k = \pm 1, i_k \in I_{-\varepsilon_k}, j_k \in I_{-\varepsilon'_k}, \forall k\} \\ & \cup \{\tau^{\varepsilon n} h(\sigma_{-\varepsilon}) \tau^{-\varepsilon n} \underbrace{h(\underbrace{\iota_{-\varepsilon}, \varepsilon, \dots, \iota_{-\varepsilon}, \varepsilon}_{n \text{ times}}; \sigma_{-\varepsilon}^{-1})}_{n \text{ times}} \mid \sigma_{-\varepsilon} \in \Gamma_{-\varepsilon}, \varepsilon = \pm 1, n \in \mathbb{N}\} \\ & \cup \{\tau^{\varepsilon n} h(\sigma_{-\varepsilon}) \tau^{-\varepsilon n} \mid n \in \mathbb{N}, \sigma_{-\varepsilon} \in \Gamma_{-\varepsilon} \cap [\Gamma, \Gamma], \varepsilon = \pm 1\} \cup \{h(\sigma) \mid \sigma \in [\Gamma, \Gamma]\}. \end{aligned} \tag{8}$$

Take any element $a \in \Xi \setminus \{1\}$. It remains to show that $\langle\langle a \rangle\rangle_\Xi = \Xi$. Relations (R3), (R8), and (R9) and Lemma 2.9 (iii) imply that we can find a big enough n and i_k 's so that the element $h(i_1, \varepsilon_1, \dots, i_n, \varepsilon_n; \sigma_n)$ does not commute with a and does not modify a . Moreover, if we take

$$v \equiv h(i_1, \varepsilon_1, \dots, i_n, \varepsilon_n; \sigma_n) h(i'_1, \varepsilon_1, \dots, i'_n, \varepsilon_n; \sigma_n^{-1}) \in \Xi \setminus \{1\},$$

for any i'_k 's (not all equal to i_k 's), we will have

$$\begin{aligned} \langle\langle a \rangle\rangle_\Xi \ni b & \equiv a v a^{-1} v \\ & = h(p_1, l_1, \dots, p_m, l_m; \sigma_n) h(p'_1, l'_1, \dots, p'_d, l'_d; \sigma_n^{-1}) \\ & \quad \cdot h(i_1, \varepsilon_1, \dots, i_n, \varepsilon_n; \sigma_n) h(i'_1, \varepsilon_1, \dots, i'_n, \varepsilon_n; \sigma_n^{-1}) \end{aligned}$$

for some m, d, p_k 's, p'_k 's, l_k 's, and l'_k 's.

Now, it is clear that we can find big enough s and appropriate e_k 's, e''_k 's, j_k 's, and j''_k 's, so that $h(j''_1, e''_1, \dots, j''_s, e''_s; \sigma^{-1})$ commutes with b and $h(j_1, e_1, \dots, j_s, e_s; \sigma)$

does not. Then,

$$\begin{aligned} \langle\langle a \rangle\rangle_{\Xi} \ni b' &\equiv bh(j_1, e_1, \dots, j_s, e_s; \sigma)h(j_1'', e_1'', \dots, j_s'', e_s''; \sigma^{-1})b^{-1} \\ &\quad h(j_1'', e_1'', \dots, j_s'', e_s''; \sigma)h(j_1, e_1, \dots, j_s, e_s; \sigma^{-1}) \\ &= h(j_1', e_1, \dots, j_s', e_s; \sigma)h(j_1, e_1, \dots, j_s, e_s; \sigma^{-1}) \neq 1 \end{aligned}$$

for some j_k' 's, from relation (R3). We can take s to be big enough and adjust the 'tail' of $(j_1, e_1, \dots, j_s, e_s)$ so that $e_1 + \dots + e_n = 0$. Since the tuples $(j_1, e_1, \dots, j_s, e_s)$ and $(j_1', e_1, \dots, j_s', e_s)$ are different, it follows from Lemma 2.9 (i) and from the assumption $\varepsilon_1 + \dots + \varepsilon_n = 0$ that

$$\beta b' \beta^{-1} = h(p_1'', e_1''', \dots, p_k'', e_k''', p'', e_s; \sigma)h(\sigma^{-1}) \in \langle\langle a \rangle\rangle_{\Xi}$$

for some $k \in \mathbb{N}$, p_t'' 's, and e_t''' 's, where

$$\begin{aligned} \Xi \ni \beta &= \tau^{-e_s}(\lambda_{-e_s}^{j_s})^{-1} \dots \tau^{-e_1}(\lambda_{-e_s}^{j_1})^{-1} \\ &\cdot \prod_{e_k=-1} h(\rho_1^k, w_1^k, \dots, \rho_{t_k}^k, w_{t_k}^k, w, 1; \mu_{-e_k}^{j_k}) \cdot \prod_{e_k=1} h(\bar{\rho}_1^k, \bar{w}_1^k, \dots, \bar{\rho}_{t_k}^k, \bar{w}_{t_k}^k, \bar{w}, -1; \mu_{-e_k}^{j_k}), \end{aligned}$$

and where the last two factors are chosen appropriately to bring β into Ξ . This argument does not depend on the 'tail' of $(p_1, e_1, \dots, p_s, e_s)$, therefore we can take e_s to be either 1 or -1 .

We conclude that the following are elements of $\langle\langle a \rangle\rangle_{\Xi}$:

$$\begin{aligned} c &= h(\sigma_1)h(\iota_1, -1, p_1, e_1, \dots, p_k, e_k, p, 1; \sigma_1^{-1}) \text{ and} \\ d &= h(\sigma_{-1})h(\iota_{-1}, 1, q_1, l_1, \dots, q_k, l_k, q, -1; \sigma_{-1}^{-1}) \end{aligned}$$

for any big enough even number k , for any $\sigma_1 \in \Gamma_1$ and $\sigma_{-1} \in \Gamma_{-1}$, and for some p_m 's, q_m 's, e_m 's, and l_m 's.

We claim that, in the tuples $(\iota_1, -1, p_1, e_1, \dots, p_k, e_k, p, 1)$ and $(\iota_{-1}, 1, q_1, l_1, \dots, q_k, l_k, q, -1)$, the indices p , q , p_t 's, and q_t 's can be chosen arbitrary. To see this, consider

$$\Xi \ni f = h(\iota_1, -1, p_1, e_1, \dots, p_t, e_t; \omega_t)h(q_0, -1, q_1, o_1, \dots, q_r, o_r, q, e_t; \omega_t^{-1}),$$

where $q_0 \neq \iota_1$ and where the second factor is chosen appropriately. Then by relation (R3),

$$fcf^{-1} = h(\sigma_1)h(\iota_1, -1, p_1, e_1, \dots, \omega_t(p_{t+1}), \dots, p_k, e_k, p, 1; \sigma_1^{-1}) \in \langle\langle a \rangle\rangle_{\Xi}.$$

Because of the transitivity and 2-transitivity of Σ_{-1} and Σ_1 , the claim is proven. The element d can be manipulated similarly.

Now, consider

$$\Xi \ni s = h(\iota_{-1}, 1, i_2, \varepsilon_2, \dots, i_t, \varepsilon_t; \omega_t)h(\iota_1, -1, q_1', o_1', \dots, q_r', o_r', q', e_t; \omega_t^{-1})$$

for an appropriate choice of q'_i 's and p_i 's so it commutes with $h(\iota_1, -1, p_1, e_1, \dots, p_k, e_k, p, 1; \sigma_1^{-1})$. Therefore

$$scs^{-1}c^{-1} = h(\iota_{-1}, 1, i_2, \varepsilon_2, \dots, i_t, \varepsilon_t; \omega_t)h(\sigma_1(\iota_{-1}), 1, i_2, \varepsilon_2, \dots, i_t, \varepsilon_t; \omega_t^{-1}) \in \langle\langle a \rangle\rangle_{\Xi},$$

so by the transitivity of the group Σ_{-1} , we see that every element of the form

$$h(\iota_{-1}, 1, i_2, \varepsilon_2, \dots, i_t, \varepsilon_t; \omega_t)h(i_1, 1, i_2, \varepsilon_2, \dots, i_t, \varepsilon_t; \omega_t^{-1})$$

belongs to $\langle\langle a \rangle\rangle_{\Xi}$. Products of such elements yield

$$h(i'_1, 1, i_2, \varepsilon_2, \dots, i_t, \varepsilon_t; \omega_t)h(i_1, 1, i_2, \varepsilon_2, \dots, i_t, \varepsilon_t; \omega_t^{-1}) \in \langle\langle a \rangle\rangle_{\Xi}$$

for any $i_1, i'_1 \in I_{-1}$. By making the same argument that uses transitivity and 2-transitivity, we see that we can change the i_l indices of the first factor, so we infer that the first set of (8) belongs to $\langle\langle a \rangle\rangle_{\Xi}$.

Consider an integer $n \geq 2$, an even number $k \geq 2$, and an appropriate $h(j_1, \varepsilon'_1, \dots, j_k, \varepsilon'_k; \sigma)$ that commutes with $h(i_1, \varepsilon_1, i_2, \varepsilon_2, \dots, i_n, \varepsilon_n; \sigma_n)$ and with $h(\iota_{-\varepsilon_1}, \varepsilon_1, i_2, \varepsilon_2, \dots, i_n, \varepsilon_n; \sigma_n^{-1})$ and has the property that

$$\delta' \equiv \tau^{\varepsilon_1}h(\sigma)\tau^{-\varepsilon_1}h(j_1, \varepsilon'_1, \dots, j_k, \varepsilon'_k; \sigma^{-1})$$

belongs to Ξ . Then

$$\begin{aligned} & \delta'h(i_1, \varepsilon_1, i_2, \varepsilon_2, \dots, i_n, \varepsilon_n; \sigma_n)h(\iota_{-\varepsilon_1}, \varepsilon_1, i_2, \varepsilon_2, \dots, i_n, \varepsilon_n; \sigma_n^{-1})(\delta')^{-1} \\ &= h(\iota_{-\varepsilon_1}, \varepsilon_1, \sigma(\iota_{\varepsilon_1}), -\varepsilon_1, i_1, \varepsilon_1, i_2, \varepsilon_2, \dots, i_n, \varepsilon_n; \sigma_n) \\ & \quad h(\iota_{-\varepsilon_1}, \varepsilon_1, \sigma(i_2), \varepsilon_2, \dots, i_n, \varepsilon_n; \sigma_n^{-1}) \in \langle\langle a \rangle\rangle_{\Xi}. \end{aligned}$$

Products of those elements with elements from the first set give all the elements from the second set of (8), so it is included in $\langle\langle a \rangle\rangle_{\Xi}$.

The third set of (8) belongs to $\langle\langle a \rangle\rangle_{\Xi}$ since its elements are products of the elements c and d above with elements from the second set.

A generic element of the fourth set of (8) can be written as

$$\begin{aligned} & h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m, i, \varepsilon, j, -\varepsilon, \bar{i}, \varepsilon, j_2, \varepsilon'_2, \dots, j_n, \varepsilon'_n; \sigma) \cdot \\ & \quad h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m, j', -\varepsilon, i', \varepsilon, \bar{i}, \varepsilon, j_2, \varepsilon'_2, \dots, j_n, \varepsilon'_n; \sigma^{-1}), \quad (9) \end{aligned}$$

where we have written $\varepsilon'_1 = \varepsilon$. We must show that this element belongs to $\langle\langle a \rangle\rangle_{\Xi}$.

First, we start with the following element from the first set of (8)

$$\begin{aligned} \langle\langle a \rangle\rangle_{\Xi} \ni z &= h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m, i, \varepsilon, \iota_{-\varepsilon}, \varepsilon, q, -\varepsilon, j, -\varepsilon, \bar{i}, \varepsilon, j_2, \varepsilon'_2, \dots, j_n, \varepsilon'_n; \sigma) \cdot \\ & \quad h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m, i, \varepsilon, \iota_{-\varepsilon}, \varepsilon, q, -\varepsilon, \iota_{\varepsilon}, -\varepsilon, \bar{i}, \varepsilon, j_2, \varepsilon'_2, \dots, j_n, \varepsilon'_n; \sigma^{-1}), \end{aligned}$$

where $q \in I'_{\varepsilon}$.

Next, using Lemma 2.9 (i) and adopting the notations thereof, we define

$$\Xi \ni \gamma = \lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \cdots \lambda_{-\varepsilon_m}^{i_m} \tau^{\varepsilon_m} \lambda_{-\varepsilon}^i \tau^{2\varepsilon} (\lambda_{\varepsilon}^q)^{-1} \tau^{-2\varepsilon} (\lambda_{-\varepsilon}^i)^{-1} \\ \cdot \tau^{-\varepsilon_m} (\lambda_{-\varepsilon_m}^{i_m})^{-1} \cdots \tau^{-\varepsilon_1} (\lambda_{-\varepsilon_1}^{i_1})^{-1} \cdot h(r_1, e_1, \dots, r_{2l-1}, e_{2l-1}, \bar{r}_{-\varepsilon}, \varepsilon; \mu_{\varepsilon}^q)$$

for appropriate r_k 's and e_k 's satisfying $e_1 + \cdots + e_{2l-1} + \varepsilon = 0$ and for which the last factor commutes with everything in the next expressions. Then

$$\gamma z \gamma^{-1} = h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m, i, \varepsilon, j, -\varepsilon, \bar{i}, \varepsilon, j_2, \varepsilon'_2, \dots, j_n, \varepsilon'_n; \sigma) \cdot \bar{h},$$

where

$$\bar{h} \equiv \gamma h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m, i, \varepsilon, \iota_{-\varepsilon}, \varepsilon, q, -\varepsilon, \iota_{\varepsilon}, -\varepsilon, \bar{i}, \varepsilon, j_2, \varepsilon'_2, \dots, j_n, \varepsilon'_n; \sigma^{-1}) \gamma^{-1} \\ = \lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \cdots \lambda_{-\varepsilon_m}^{i_m} \tau^{\varepsilon_m} \lambda_{-\varepsilon}^i \lambda_{-\varepsilon}^{\bar{i}} \tau^{\varepsilon} \cdots \lambda_{-\varepsilon'_n}^{j_n} \tau^{\varepsilon'_n} h(\sigma^{-1}) \\ \cdot \tau^{-\varepsilon'_n} (\lambda_{-\varepsilon'_n}^{j_n})^{-1} \cdots \tau^{-\varepsilon} (\lambda_{-\varepsilon}^{\bar{i}})^{-1} (\lambda_{-\varepsilon}^i)^{-1} \tau^{-\varepsilon_m} (\lambda_{-\varepsilon_m}^{i_m})^{-1} \cdots \tau^{-\varepsilon_1} (\lambda_{-\varepsilon_1}^{i_1})^{-1} \\ = \lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \cdots \lambda_{-\varepsilon_m}^{i_m} \tau^{\varepsilon_m} \lambda_{-\varepsilon}^i h(\bar{i}, \varepsilon, j_2, \varepsilon'_2, \dots, j_n, \varepsilon'_n; \sigma^{-1}) (\lambda_{-\varepsilon}^i)^{-1} \tau^{-\varepsilon_m} (\lambda_{-\varepsilon_m}^{i_m})^{-1} \cdots \tau^{-\varepsilon_1} (\lambda_{-\varepsilon_1}^{i_1})^{-1} \\ = \lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \cdots \lambda_{-\varepsilon_m}^{i_m} \tau^{\varepsilon_m} h(\mu_{-\varepsilon}^i(\bar{i}), \varepsilon, j_2, \varepsilon'_2, \dots, j_n, \varepsilon'_n; \sigma^{-1}) \tau^{-\varepsilon_m} (\lambda_{-\varepsilon_m}^{i_m})^{-1} \cdots \tau^{-\varepsilon_1} (\lambda_{-\varepsilon_1}^{i_1})^{-1}.$$

Likewise, we consider the following element from the first set of (8)

$$\langle\langle a \rangle\rangle_{\Xi} \ni z' = h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m, j', -\varepsilon, \iota_{\varepsilon}, -\varepsilon, p, \varepsilon, \iota_{-\varepsilon}, \varepsilon, \mu_{-\varepsilon}^i(\bar{i}), \varepsilon, j_2, \varepsilon'_2, \dots, j_n, \varepsilon'_n; \sigma) \\ \cdot h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m, j', -\varepsilon, \iota_{\varepsilon}, -\varepsilon, p, \varepsilon, i', \varepsilon, \mu_{-\varepsilon}^i(\bar{i}), \varepsilon, j_2, \varepsilon'_2, \dots, j_n, \varepsilon'_n; \sigma^{-1}),$$

where $p \in I'_{-\varepsilon}$ and define

$$\Xi \ni \gamma' = \lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \cdots \lambda_{-\varepsilon_m}^{i_m} \tau^{\varepsilon_m} \lambda_{\varepsilon}^{j'} \tau^{-2\varepsilon} (\lambda_{-\varepsilon}^p)^{-1} \tau^{2\varepsilon} (\lambda_{\varepsilon}^{j'})^{-1} \\ \cdot \tau^{-\varepsilon_m} (\lambda_{-\varepsilon_m}^{i_m})^{-1} \cdots \tau^{-\varepsilon_1} (\lambda_{-\varepsilon_1}^{i_1})^{-1} \cdot h(r'_1, e_1, \dots, r'_{2l-1}, e_{2l-1}, \bar{r}_{-\varepsilon}, \varepsilon; \mu_{-\varepsilon}^p)$$

for appropriate r'_k 's. Then,

$$\gamma' z' (\gamma')^{-1} = \bar{h} \cdot h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m, j', -\varepsilon, i', \varepsilon, \mu_{-\varepsilon}^i(\bar{i}), \varepsilon, j_2, \varepsilon'_2, \dots, j_n, \varepsilon'_n; \sigma^{-1}),$$

where

$$\bar{h} \equiv \gamma' h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m, j', -\varepsilon, \iota_{\varepsilon}, -\varepsilon, p, \varepsilon, \iota_{-\varepsilon}, \varepsilon, \mu_{-\varepsilon}^i(\bar{i}), \varepsilon, j_2, \varepsilon'_2, \dots, j_n, \varepsilon'_n; \sigma) (\gamma')^{-1} \\ = \lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \cdots \lambda_{-\varepsilon_m}^{i_m} \tau^{\varepsilon_m} \lambda_{\varepsilon}^{j'} h(\mu_{-\varepsilon}^i(\bar{i}), \varepsilon, j_2, \varepsilon'_2, \dots, j_n, \varepsilon'_n; \sigma) (\lambda_{\varepsilon}^{j'})^{-1} \tau^{-\varepsilon_m} (\lambda_{-\varepsilon_m}^{i_m})^{-1} \cdots \tau^{-\varepsilon_1} (\lambda_{-\varepsilon_1}^{i_1})^{-1} \\ = \lambda_{-\varepsilon_1}^{i_1} \tau^{\varepsilon_1} \cdots \lambda_{-\varepsilon_m}^{i_m} \tau^{\varepsilon_m} h(\mu_{\varepsilon}^{j'}(\mu_{-\varepsilon}^i(\bar{i})), \varepsilon, j_2, \varepsilon'_2, \dots, j_n, \varepsilon'_n; \sigma) \tau^{-\varepsilon_m} (\lambda_{-\varepsilon_m}^{i_m})^{-1} \cdots \tau^{-\varepsilon_1} (\lambda_{-\varepsilon_1}^{i_1})^{-1} \\ = (\bar{h})^{-1},$$

since $\mu_{\varepsilon}^{j'}(\mu_{-\varepsilon}^i(\bar{i})) = \mu_{-\varepsilon}^i(\bar{i})$, due to relation (R6) and $\mu_{-\varepsilon}^i(\bar{i}) \in I_{-\varepsilon}$. Finally,

$$\langle\langle a \rangle\rangle_{\Xi} \ni \gamma z \gamma^{-1} \cdot \gamma' z' (\gamma')^{-1} = h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m, i, \varepsilon, j, -\varepsilon, \bar{i}, \varepsilon, j_2, \varepsilon'_2, \dots, j_n, \varepsilon'_n; \sigma) \\ \cdot h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m, j', -\varepsilon, i', \varepsilon, \mu_{-\varepsilon}^i(\bar{i}), \varepsilon, j_2, \varepsilon'_2, \dots, j_n, \varepsilon'_n; \sigma^{-1}),$$

and after a multiplication with an element from the first set of (8), we get the element (9).

Therefore the fourth set of (8) is in $\langle\langle a \rangle\rangle_{\Xi}$.

Repeating almost verbatim the corresponding part of the proof of Theorem [8, Theorem 3.16] gives us that the seventh set of (8) belongs to $\langle\langle a \rangle\rangle_{\Xi}$. Note that if $\Sigma_{\varepsilon} = \text{Sym}(2)$, then $[\Sigma_{\varepsilon}, \Sigma_{\varepsilon}]$ is the trivial group.

Next, we take numbers $m > n$ and

$$\gamma'' = \tau^{\varepsilon m} h(\sigma'_{-\varepsilon}) \tau^{-\varepsilon m} h(j_1, \varepsilon, \dots, j_{m+1}, \varepsilon, j, -\varepsilon; (\sigma'_{-\varepsilon})^{-1}) \in \Xi,$$

where $\sigma'_{-\varepsilon} \in \Gamma_{-\varepsilon}$, $j_k \in I'_{-\varepsilon}$, $\forall k$, and $j \in I'_{\varepsilon}$, with the relation $(\sigma'_{-\varepsilon})^{-1}(\iota_{\varepsilon}) = q$ for some $q \in I'_{\varepsilon}$.

After that, we take the following element of $\langle\langle a \rangle\rangle_{\Xi}$ (it is a product of elements from the second and fourth set)

$$\begin{aligned} x \equiv & \underbrace{h(\iota_{-\varepsilon}, \varepsilon, \dots, \iota_{-\varepsilon}, \varepsilon, q, -\varepsilon, \iota_{\varepsilon}, -\varepsilon, \dots, \iota_{\varepsilon}, -\varepsilon; \sigma_{-\varepsilon})}_{m \text{ times}} \cdot \underbrace{h(\iota_{\varepsilon}, -\varepsilon, \dots, \iota_{\varepsilon}, -\varepsilon; \sigma_{-\varepsilon})}_{m-n-1 \text{ times}} \\ & \cdot \underbrace{h(\iota_{-\varepsilon}, \varepsilon, \dots, \iota_{-\varepsilon}, \varepsilon, q, -\varepsilon, \iota_{\varepsilon}, -\varepsilon, \dots, \iota_{\varepsilon}, -\varepsilon, p, \varepsilon, \iota_{-\varepsilon}, \varepsilon, \dots, \iota_{-\varepsilon}, \varepsilon; \sigma_{-\varepsilon})}_{m \text{ times}} \cdot \underbrace{h(\iota_{-\varepsilon}, \varepsilon, \dots, \iota_{-\varepsilon}, \varepsilon; \sigma_{-\varepsilon})}_{n-1 \text{ times}}, \end{aligned}$$

where $p \in I'_{-\varepsilon}$. Then

$$\gamma'' x (\gamma'')^{-1} = \tau^{\varepsilon n} h(\sigma_{-\varepsilon}) \tau^{-\varepsilon n} \cdot h(p, \varepsilon, \underbrace{\iota_{-\varepsilon}, \varepsilon, \dots, \iota_{-\varepsilon}, \varepsilon}_{n-1 \text{ times}}; \sigma_{-\varepsilon}) \in \langle\langle a \rangle\rangle_{\Xi}.$$

Therefore upon a multiplication by an element from the first set of (8), we infer that the fifth set of (8) belongs to $\langle\langle a \rangle\rangle_{\Xi}$.

Finally, the argument from Theorem [8, Theorem 3.16] can be used for the sixth set of (8) the same way it was used for the seventh set.

This completes the proof. \square

Remark 2.14. *The example introduced in [3, Section 5] corresponds to the case $\Sigma_{-1} \cong \Sigma_1 \cong \text{Sym}(2)$. Theorem 2.13 corresponds to [3, Proposition 5.11] in this particular case.*

2.5. ANALYTIC STRUCTURE

In this section, we use some results from [8, Section 2].

Lemma 2.15. *The group $\Lambda = \Lambda[I_{-1}, I_1, \iota_{-1}, \iota_1; \Sigma_{-1}, \Sigma_1]$ is a non-ascending HNN-extension and its action on its Bass-Serre tree is minimal and of general type.*

Proof. Since the action is transitive, it is minimal. Since $H \neq G \neq \theta(H)$, then Λ is nondegenerate and non-ascending. The result now follows from [7, Proposition 20]. \square

Theorem 2.16. *The HNN-extension $\Lambda = \Lambda[I_{-1}, I_1, \iota_{-1}, \iota_1; \Sigma_{-1}, \Sigma_1]$ has a unique trace. It is C^* -simple if and only if either one of the groups Σ_{-1} and Σ_1 is non-amenable.*

Proof. Lemma 2.15 enables us to apply [3, Theorem 4.19] to conclude that Λ has the unique trace property since $\ker \Lambda$ is trivial. It also enables us to apply [3, Theorem 4.20] to conclude that Λ is C^* -simple if and only if K_{-1} and K_1 are non-amenable, which, by Lemma 2.12, is equivalent to the requirement that some of the groups Σ_{-1} and Σ_1 is non-amenable. \square

Finally, we prove

Theorem 2.17. *The HNN-extension $\Lambda = \Lambda[\Sigma_{-1}, \Sigma_1]$ is not inner amenable.*

Proof. Lemma 2.15 allows us to apply [8, Proposition 2.3], so we need to show that the action of $\Lambda = \Lambda[I_{-1}, I_1, \iota_{-1}, \iota_1; \Sigma_{-1}, \Sigma_1]$ on its Bass-Serre is finitely fledged.

For this, take any elliptic element $g \in \Lambda \setminus \{1\}$. Since g fixes some vertex, it is a conjugate of an element of G . The finite fledgedness property is conjugation invariant, so we can assume $g \in G \setminus \{1\}$.

From Lemma 2.9 (ii), we can write $g = h(\sigma)h_{-1}h_1$, where $\sigma \in \Gamma$,

$$h_{-1} = \prod_{k=1}^m h(i_1^k, -1, i_2^k, \varepsilon_{k,2}, \dots, i_{n_k}^k, \varepsilon_{k,n_k}; \sigma_k),$$

$$h_1 = \prod_{l=m+1}^r h(i_1^l, 1, i_2^l, \varepsilon_{l,2}, \dots, i_{n_l}^l, \varepsilon_{l,n_l}; \theta_l),$$

$r \geq m \geq 0$, $\sigma_k \in \Gamma_{\varepsilon_{k,n_k}}$, $\theta_l \in \Gamma_{\varepsilon_{l,n_l}}$, and $i_z^p \in I'_{\varepsilon_p, z}$. We also require $0 \leq n_1 \leq \dots \leq n_m$ and $0 \leq n_{m+1} \leq \dots \leq n_r$.

Let us assume that g fixes a vertex $v = v(i_1, \varepsilon_1, \dots, i_n, \varepsilon_n)$, where $n \geq \max\{n_m, n_r\} + 1$, and let's take $w = v(i_1, \varepsilon_1, \dots, i_n, \varepsilon_n, \dots, i_{n+d}, \varepsilon_{n+d})$ for any $d \geq 1$. We note that $h_{-\varepsilon_1}$ fixes w and $h(\sigma)h_{\varepsilon_1}$ modifies only indices with numbers no greater than $\{n_m, n_r\} + 1 \leq n$. Therefore

$$h(\sigma)h_{\varepsilon_1}v = v(i'_1, \varepsilon_1, \dots, i'_n, \varepsilon_n) \text{ and}$$

$$h(\sigma)h_{\varepsilon_1}w = v(i'_1, \varepsilon_1, \dots, i'_n, \varepsilon_n, i_{n+1}, \varepsilon_{n+1}, \dots, i_{n+d}, \varepsilon_{n+d})$$

for some $i'_k \in I'_{-\varepsilon_k}$. By our assumption, it follows that

$$v = gv = h(\sigma)h_{\varepsilon_1}v = v(i'_1, \varepsilon_1, \dots, i'_n, \varepsilon_n).$$

Thus $i'_k = i_k$ for all $1 \leq k \leq n$, and therefore $gw = w$.

This concludes the proof. \square

Corollary 2.18. *Theorems 2.16 and 2.13 imply:*

If either Σ_{-1} or Σ_1 is non-amenable, then the amenablsh radical of Λ is trivial.

If Σ_{-1} and Σ_1 are both amenable, then Λ is amenablsh.

Proof. If we show that the centralizer $C_\Lambda(\Xi)$ is trivial, [2, Theorem 4.1] will imply that Λ is C^* -simple if and only if Ξ is C^* -simple. Since Ξ is simple, if it is not C^* -simple, then it is amenablsh, and therefore Λ is also amenablsh because $(\Gamma/[\Gamma, \Gamma]) \wr_{\mathbb{Z}} \mathbb{Z}$ is amenable. If Ξ is C^* -simple, then so is Λ , thus both of their amenablsh radicals are trivial.

To illustrate that $C_\Lambda(\Xi)$ is trivial, assume that there is a nontrivial $g \in C_\Lambda(\Xi)$. Then g can be written as in Lemma 2.9 (iii), and using relations (R3), (R7), and (R8), we can find a non-trivial element of Ξ

$$h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m, j_1, \varepsilon'_1, \dots, j_n, \varepsilon'_n; \sigma) \cdot h(i_1, \varepsilon_1, \dots, i_m, \varepsilon_m, j'_1, \varepsilon''_1, \dots, j'_n, \varepsilon''_n; \sigma^{-1})$$

that does not commute with g , a contradiction. \square

3. REFERENCES

- [1] Baumslag, G.: *Topics in Combinatorial Group Theory*, Birkhäuser, 1993.
- [2] Breuillard, E., Kalantar, M., Kennedy, M., and Ozawa, N.: C^* -simplicity and the unique trace property for discrete groups. *Publications mathématiques de l'IHÉS* **126** (2017), 35–71.
- [3] Bryder, R. S., Ivanov, N. A., and T. Omland, T.: C^* -simplicity of HNN-extensions and groups acting on trees. *Annales de L'Institut Fourier* **70**, no 4, (2020), 1497–1543.
- [4] Cohen, D.: *Combinatorial Group Theory: a Topological Approach*, Cambridge University Press, 1989.
- [5] E. Effros, E.: Property Γ and inner amenability. *Proc. Amer. Math. Soc.* **47**, no 2 (1975), 483–486.
- [6] de la Harpe, P.: On simplicity of reduced C^* -algebras of groups. *Bull. Lond. Math. Soc.* **39**, no. 1 (2007), 1–26.
- [7] de la Harpe, P., and Préaux, J.-P.: C^* -simple groups: amalgamated free products, HNN-extensions, and fundamental groups of 3-manifolds. *J. Topol. Anal.* **3**, no. 4 (2011), 451–489.
- [8] Ivanov, N. A.: Examples of group amalgamations with nontrivial quasi-kernels. *Serdica Math. J.l* **46**, no. 4 (2020), 357–396.
- [9] Ivanov, N. A. and Omland, T.: C^* -simplicity of free products with amalgamation and radical classes of groups. *J. Funct. Anal.* **272**, no. 9 (2017), 3712–3741.
- [10] Kalantar, M. and Kennedy, M.: Boundaries of reduced C^* -algebras of discrete groups. *J. Reine Angew. Math.* **2017**, Issue 727, 247–267.
- [11] Le Boudec, A.: C^* -simplicity and the amenable radical. *Invent. Math.* **209**, no. 1 (2017), 159–174.

- [12] F. J. Murray, F. J. and von Neumann, J.: On rings of operators. IV. *Ann. Math. (2)* **44** (1943), 716–808.
- [13] Jean-Pierre Serre, J.-P.: *Trees* (translation of "Arbres, Amalgames, SL_2 "), Springer, 2003.
- [14] Vaes, S.: An inner amenable group whose von Neumann algebra does not have property Gamma. *Acta Math.* **208**, no. 2 (2012), 389–394.

Received on April 14, 2021

NIKOLAY A. IVANOV
Faculty of Mathematics and Informatics
Sofia University St Kliment Ohridski
5 James Bourchier Blvd.
1164 Sofia
BULGARIA
E-mail: nivanov@fmi.uni-sofia.bg

ГОДИШНИК НА СОФИЙСКИЯ УНИВЕРСИТЕТ „СВ. КЛИМЕНТ ОХРИДСКИ“

ФАКУЛТЕТ ПО МАТЕМАТИКА И ИНФОРМАТИКА

Том 106

ANNUAL OF SOFIA UNIVERSITY „ST. KLIMENT OHRIDSKI“

FACULTY OF MATHEMATICS AND INFORMATICS

Volume 107

ON THE REGULARITY OF CERTAIN THREE-ROW ALMOST HERMITIAN INCIDENCE MATRICES

GENO NIKOLOV, BORISLAVA PETROVA

In this note we present a new proof of the regularity of a class of three-row almost Hermitian matrices, based on some properties of Legendre polynomials.

Keywords: Incidence matrices, Birkhoff interpolation, Legendre polynomials, Gegenbauer polynomials.

2000 Math. Subject Classification: 41A55, 65D30, 65D32.

1. INTRODUCTION AND STATEMENT OF THE RESULT

Throughout this paper, the notation π_n will stand for the set of algebraic polynomials of degree not exceeding n . In 1906 G. D. Birkhoff [2] formulated a general problem on interpolation by algebraic polynomials, which includes as particular cases the Lagrange and Hermite interpolation problems. Before formulating the Birkhoff interpolation problem (BIP), we need the following:

Definition 1. An *incidence matrix* $E = \{e_{ij}\}_{i=1, j=0}^n$ is a matrix with elements $e_{ij} \in \{0, 1\}$. The number of 1-entries in E is denoted by $|E|$, and we shall assume always that E is a *normal incidence matrix*, i.e., $|E| = r + 1$.

The Birkhoff interpolation problem (BIP). Given an incidence matrix $E = \{e_{ij}\}_{i=1, j=0}^n$, a vector of interpolation nodes $\mathbf{X} = (x_1, x_2, \dots, x_n) \in \mathbb{R}^n$,

$x_1 < x_2 < \dots < x_n$, and a data set $\{\gamma_{ij} \in \mathbb{C} : e_{ij} = 1\}$, find a polynomial $p \in \pi_{|E|-1}$ such that

$$p^{(j)}(x_i) = \gamma_{ij}, \quad \{i, j\} : e_{ij} = 1. \quad (1.1)$$

It should be pointed out that, unlike the Lagrange and Hermite interpolation problems, which are known to have a unique solution, the general BIP is not always solvable.

Definition 2. An incidence matrix $E = \{e_{ij}\}_{i=1, j=0}^{n, r}$ is said to be (*order*) *regular*, if for every vector of interpolation nodes $\mathbf{X} = (x_1, x_2, \dots, x_n) \in \mathbb{R}^n$, $x_1 < x_2 < \dots < x_n$, and a data set $\{\gamma_{ij} \in \mathbb{C} : e_{ij} = 1\}$, the BIP (1.1) has a unique solution.

Surprisingly enough, despite the efforts of many mathematicians, the problem of complete characterization of the regular incidence matrices remains open. A simple necessary condition for regularity was found by Pólya.

Pólya condition. A necessary condition for $E = \{e_{ij}\}_{i=1, j=0}^{n, r}$ to be regular is

$$\sum_{i=1}^n \sum_{j=0}^k e_{ij} \geq k + 1, \quad k = 0, \dots, |E| - 1. \quad (1.2)$$

In 1969 Atkinson and Sharma [1] found a simple sufficient condition for regularity. We need another definition before formulating their result.

Definition 3. A *block* is called any maximal sequence of 1-entries in a row of E . A block $e_{ij} = e_{i, j+1} = \dots = e_{i, j+\ell-1} = 1$ is called *even*, resp. *odd*, if its length ℓ is even, resp. odd number. The smallest column index j of 1-entry in a block defines its *level*. *Hermitian block* is a block with level 0.

A row $\mathbf{e}_i = (e_{i,0}, e_{i,1}, \dots, e_{i,r})$ of E is called *Hermitian row of length k* if it contains a single block which is Hermitian with length k .

A block $e_{ij} = e_{i, j+1} = \dots = e_{i, j+\ell-1} = 1$ in an interior row \mathbf{e}_i , $1 < i < n$, is called *supported*, if there are 1-entries in rows i_1 and i_2 , $i_1 < i < i_2$ with column indices $j_1, j_2 < j$.

Atkinson–Sharma Theorem. *Every incidence matrix $E = \{e_{ij}\}_{i=1, j=0}^{n, r}$ which satisfies the Pólya condition (1.2) and does not contain supported odd blocks is regular.*

Note that the incidence matrices corresponding to Lagrange’s and Hermite’s interpolation problems fulfill the assumptions of the Atkinson–Sharma Theorem. Indeed, their incidence matrices contain only Hermitian rows (with length one in the Lagrange case), therefore obviously satisfy the Pólya condition and, as their rows contain only blocks with level 0, these blocks are not be supported.

Atkinson and Sharma also conjectured that all matrices that contain odd supported blocks are not regular. However, Lorentz and Zeller [11] found a counterexample to this conjecture, showing that the three-row incidence matrix

$$E = \begin{pmatrix} 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 \end{pmatrix} \quad (1.3)$$

is regular, despite having two odd supported blocks.

Since the problem of characterizing the regularity of general incidence matrices turns out to be a very difficult one, some authors [3, 4, 5, 6, 9, 10] have studied the special class of almost Hermitian matrices, which are incidence matrices which have only one (interior) non-Hermitian row. Special attention has been paid to the three-row almost Hermitian matrices. Particular reason for the interest in three-row matrices is that, by applying technique of splitting (de-coalescence) of rows, singularity of such matrices can imply singularity of incidence matrices with more rows, see e.g. [8].

Definition 4. A three-row almost Hermitian incidence matrix $E(p, q; k_1, k_2)$ is an incidence matrix with its first and third row Hermitian of length p and q , respectively (with $p \leq q$), and single 1-entries (blocks of length one) in the middle row in positions k_1 and k_2 , where $1 \leq k_1 < k_2 - 1$ (the case $k_1 = k_2 - 1$ is handled by the Atkinson-Sharma theorem).

It follows from the results in [9, 10] that $E(p, q; k_1, k_2)$ is not regular unless one of the following conditions is satisfied (see [13, Theorem 8.5]):

$$p \leq k_1 < k_2 - 1 \leq q, \quad (1.4)$$

$$q + 1 < k_2 \quad \text{and} \quad k_1 + k_2 = p + q + 1. \quad (1.5)$$

Only in the second case (called in [13, p. 104] as *the symmetric exterior case*) the regularity is completely characterized. Precisely, in this case $E(p, q; k_1, k_2)$ is regular if and only if $p = q$ (for more details, see [13, Theorem 8.15]). In the present note we present a short proof of the “if part” (the sufficiency). More precisely, we prove the following

Theorem 1. *The almost Hermitian matrix $E(m, m; k, 2m + 1 - k)$ is regular for every $k \in \mathbb{N}$, $1 \leq k < m$.*

Notice that the matrix in (1.3) corresponds to the case $m = 2, k = 1$.

Our proof of Theorem 1 makes use of some properties of the Gegenbauer polynomials, in particular of the Legendre polynomials.

2. PROOF OF THEOREM 1

The claim of Theorem 1 is equivalent to the following statement:

Proposition 1. *Let $m, k \in \mathbb{N}$, $1 \leq k < m$. Then for every $x \in (-1, 1)$ and data set $\{(a_j, b_j), j = 0, 1, \dots, m-1; c, d\}$ there exist a unique algebraic polynomial $Q(x)$ of degree not exceeding $2m + 1$ satisfying the interpolation conditions*

$$\begin{aligned} Q^{(j)}(-1) &= a_j, & j = 0, 1, \dots, m-1, \\ Q^{(j)}(1) &= b_j, & j = 0, 1, \dots, m-1, \\ Q^{(k)}(x) &= c, \\ Q^{(2m+1-k)}(x) &= d. \end{aligned} \tag{2.1}$$

The linear system for the coefficients of Q has a unique solution if and only if the corresponding homogeneous system has only trivial solution. The polynomial Q which satisfy the homogeneous system has zeros of multiplicity m at ± 1 , therefore is of the form

$$Q(t) = \omega(t) [A(t - x) + B], \quad \omega(t) = (x^2 - 1)^m$$

with constants A and B determined by $Q^{(k)}(x) = Q^{(2m+1-k)}(x) = 0$, i.e., by the linear system

$$\begin{cases} B\omega^{(k)}(x) + Ak\omega^{(k-1)}(x) = 0 \\ B\omega^{(2m+1-k)}(x) + A(2m+1-k)\omega^{(2m-k)}(x) = 0. \end{cases}$$

To prove Proposition 1, and thereby Theorem 1, we need to show that the unique solution of this last system is $A = B = 0$, which is equivalent to showing that $\Delta(x) \neq 0$ for every $x \in (-1, 1)$, where

$$\Delta(x) = k\omega^{(k-1)}(x)\omega^{(2m+1-k)}(x) - (2m+1-k)\omega^{(k)}(x)\omega^{(2m-k)}(x). \tag{2.2}$$

For the proof of (2.2) we shall use some properties of the Legendre polynomials, the orthogonal polynomials in $[-1, 1]$ with respect to the constant weight function. Recall that the n -th Legendre polynomial P_n is defined by

$$P_n(x) = \frac{1}{2^n n!} \left(\frac{d}{dx} \right)^n \{(x^2 - 1)^n\}.$$

For $j = 1, 2, \dots, m$, we define recursively the j -fold anti-derivative $S_j(x)$ of P_m by

$$S_j(x) = \int_{-1}^x S_{j-1}(t) dt, \quad S_0(x) = P_m(x).$$

In view of the definition of Legendre polynomials, we have

$$S_j(x) = \frac{1}{2^m m!} \omega^{(m-j)}(x), \quad j = 0, 1, \dots, m. \tag{2.3}$$

For the proof of (2.2) we shall need the following lemma.

Lemma 1. For $j = 1, 2, \dots, m$, there holds

$$S_j(x) = \frac{(m-j)!}{(m+j)!} (x^2-1)^j \left(\frac{d}{dx}\right)^j \{P_m(x)\}. \quad (2.4)$$

Proof. We apply backward induction on j . Since $\left(\frac{d}{dx}\right)^m \{P_m(x)\} = \frac{(2m)!}{2^m m!}$, (2.3) shows that equality (2.4) is true for $j = m$. Assuming that (2.4) is true for some j , $1 \leq j \leq m$, we obtain

$$\begin{aligned} S_{j-1}(x) &= S'_j(x) = \frac{(m-j)!}{(m+j)!} \frac{d}{dx} \left\{ (x^2-1)^j \left(\frac{d}{dx}\right)^j \{P_m(x)\} \right\} \\ &= \frac{(m-j)!}{(m+j)!} (x^2-1)^{j-1} \left\{ (x^2-1) \left(\frac{d}{dx}\right)^{j+1} \{P_m(x)\} + 2jx \left(\frac{d}{dx}\right)^j \{P_m(x)\} \right\} \\ &= \frac{(m-j)!}{(m+j)!} (x^2-1)^{j-1} \left\{ (x^2-1)z'' + 2jxz' \right\}, \end{aligned} \quad (2.5)$$

where

$$z(x) = \left(\frac{d}{dx}\right)^{j-1} \{P_m(x)\}. \quad (2.6)$$

At this point we exploit some well-known properties of the Gegenbauer polynomials. The Gegenbauer polynomial C_n^λ is the n -th orthogonal polynomial in $[-1, 1]$ with respect to the weight function $w_\lambda(x) = (1-x^2)^{\lambda-1/2}$ (and the n -th Legendre polynomials P_n equals $C_n^{1/2}$). The Gegenbauer polynomials satisfy the ordinary differential equation

$$(1-x^2)y'' - (2\lambda+1)xy' + n(n+2\lambda)y = 0, \quad y = C_n^\lambda(x) \quad (2.7)$$

and their derivatives satisfy $\frac{d}{dx} \{C_n^\lambda(x)\} = 2\lambda C_{n-1}^{\lambda+1}(x)$ (see [14, eqns. (4.7.5) and (4.7.14)]). From this last property we observe that, apart from a constant factor, the polynomial $z(x)$ in (2.6) is equal to $C_{m-j+1}^{j-1/2}(x)$. Then, according to (2.7),

$$(x^2-1)z'' + 2jxz' = (m-j+1)(m+j)z,$$

and substituting this expression in (2.5) we obtain

$$S_{j-1}(x) = \frac{(m-j+1)!}{(m+j-1)!} (x^2-1)^{j-1} z(x).$$

With this the induction step from j to $j-1$ is done. Lemma 1 is proved. \square

We proceed with the proof of (2.2). From (2.3) and

$$\left(\frac{d}{dx}\right)^j \{P_m(x)\} = \frac{1}{2^m m!} \omega^{(m+j)}(x)$$

we observe that Lemma 1 is equivalent to the identity

$$\frac{\omega^{(m-j)}(x)}{(m-j)!} = (x^2 - 1)^j \frac{\omega^{(m+j)}(x)}{(m+j)!}, \quad j = 1, \dots, m. \quad (2.8)$$

With $j = m - k + 1$ and $j = m - k$ this yields

$$\frac{\omega^{(k-1)}(x)}{(k-1)!} = (x^2 - 1)^{m-k+1} \frac{\omega^{(2m+1-k)}(x)}{(2m+1-k)!},$$

$$\frac{\omega^{(k)}(x)}{(k)!} = (x^2 - 1)^{m-k} \frac{\omega^{(2m-k)}(x)}{(2m-k)!}.$$

By expressing $\omega^{(2m+1-k)}$ and $\omega^{(2m-k)}$ and substitution in (2.2) we find that

$$\begin{aligned} \Delta(x) &= k!(2m+1-k)! \left\{ \frac{\omega^{(k-1)}(x)}{(k-1)!} \frac{\omega^{(2m+1-k)}(x)}{(2m+1-k)!} - \frac{\omega^{(k)}(x)}{k!} \frac{\omega^{(2m-k)}(x)}{(2m-k)!} \right\} \\ &= \frac{(2m+1-k)!}{k!} (x^2 - 1)^{k-m-1} \left\{ [k\omega^{(k-1)}(x)]^2 + (1-x^2)[\omega^{(k)}(x)]^2 \right\}. \end{aligned}$$

Since the zeros of $\omega^{(k-1)}$ and $\omega^{(k)}$ interlace, the sum in the last curl brackets is positive for $x \in (-1, 1)$, and consequently $\Delta(x) \neq 0$ for $x \in (-1, 1)$. With this the proof of Proposition 1 is complete.

ACKNOWLEDGEMENT. The authors are supported by the Sofia University Research Fund through Contract No. 80-10-20/22.03.2021.

3. REFERENCES

- [1] Atkinson, K. and A. Sharma, A.: A partial characterization of poised Hermite-Birkhoff interpolation problems. *SIAM J. Numer. Anal.* **6** (1969), 230–235.
- [2] Birkhoff, G.D.: General mean value theorem and remainder theorems with application to mechanical differentiation and quadrature. *Trans. Amer. Math. Soc.* **7** (1906), 107–136.
- [3] DeVore, R.A., Meir, A., and Sharma, A.: Strongly and weakly non-poised H–B interpolation problems. *Canad. J. Math.* **25** (1973), 1040–1050.
- [4] Drols, W.: On a problem of DeVore, Meir and Sharma. In: *Approximation Theory*, vol. III (E. W. Cheney, ed.), Academic Press, New York, 1980, pp. 361–366.
- [5] Drols, W.: Zur Hermite–Birkhoff Interpolation: DMS–Matrizen. *Math. Z.* **172** (1980), 179–194.
- [6] Drols, W.: Fasthermitesche Imzidenzmatrizen. *Z. Angew. Math. Mech.* **61** (1981), T275–276.
- [7] Dyn, N., Lorentz, G.G., and Riemenschneider, Sh.D. : Continuity of the Birkhoff interpolation. *SIAM J. Numer. Anal* **19** (1982), 507–509.

- [8] Karlin, S. and Karon, J. M.: On Hermite–Birkhoff interpolation. *J. Approx. Theory* **6** (1972), 90–114.
- [9] Lorentz, G. G.: The Birkhoff interpolation problem: New methods and results. In: *Linear Operators and Approximation*, vol. 2 (P. L. Butzer and B. Sz. Nagy, eds.), Birkhauser Verlag, Basel, 1975, pp. 481–501.
- [10] Lorentz, G. G., Stangler, S. S., and Zeller, K.: Regularity of some special Birkhoff matrices. In: *Approximation Theory*, vol. II (G. G. Lorentz et al., eds.), Academic Press, New York, 1976, pp. 423–436.
- [11] Lorentz, G. G. and Zeller, K.: Birkhoff interpolation. *SIAM J. Numer. Math.* **8** (1971), 43–48.
- [12] Lorentz, G. G. and Zeller, K.: Birkhoff interpolation problem: Coalescence of rows. *Arch. Math.* **26** (1975), 189–192.
- [13] Lorentz, G. G., Jetter, K., and Riemenschneider, Sh. D.: “Birkhoff interpolation”. *Encyclopedia in Mathematics and its Applications*, Vol. 19, Reading, Mass., 1983.
- [14] Szegő, G.: “Orthogonal Polynomials”, 4th ed., Amer. Math. Soc. Coll. Publ., Vol. 23, Providence, RI, 1975.

Received on July 9, 2021

GENO NIKOLOV, BORISLAVA PETROVA
 Faculty of Mathematics and Informatics
 Sofia University St Kliment Ohridski
 5, J. Bourchier blvd., BG-1164 Sofia
 BULGARIA
 E-mails: geno@fmi.uni-sofia.bg
 borislava.petrova1997@gmail.com

Submission of manuscripts. The *Annual* is published once a year. No deadline exists. Once received by the editors, the manuscript will be subjected to rapid, but thorough review process. If accepted, it is immediately scheduled for the nearest forthcoming issue. No page charge is made. The author(s) will be provided with a free of charge printable pdf file of their published paper.

The submission of a paper implies that it has not been published, or is not under consideration for publication elsewhere. In case it is accepted, it implies as well that the author(s) transfers the copyright to the Faculty of Mathematics and Informatics at the “St. Kliment Ohridski” University of Sofia, including the right to adapt the article for use in conjunction with computer systems and programs and also reproduction or publication in machine-readable form and incorporation in retrieval systems.

Instructions to Contributors. Preferences will be given to papers, not longer than 25 pages, written in English and typeset by means of a T_EX system. A simple specimen file, exposing in detail the instruction for preparation of the manuscripts, is available upon request from the electronic address of the Editorial Board.

Manuscripts should be submitted for editorial consideration in pdf-format by e-mail to `annuaire@fmi.uni-sofia.bg`. Upon acceptance of the paper, the authors will be asked to send the text of the papers in .`tex` format and the appropriate graphic files (preferably in .`eps` format).

The manuscripts should be prepared in accordance with the instructions, given below.

The first page of manuscripts must contain a title, name(s) of the author(s), a short abstract, a list of keywords and the appropriate 2010 MSC codes (primary and secondary, if necessary). The affiliation(s), including the electronic address, should be given at the end of the manuscripts.

Figures have to be inserted in the text near their first reference. If the author cannot supply and/or incorporate the graphic files, drawings (in black ink and on a good quality paper) should be enclosed separately. If photographs are to be used, only black and white ones are acceptable.

Tables should be inserted in the text as close to the point of reference as possible. Some space should be left above and below the table.

Footnotes, which should be kept to a minimum and should be brief, must be numbered consecutively.

References must be cited in the text in square brackets, like [3], or [5, 7], or [11, p. 123], or [16, Ch. 2.12]. They have to be numbered either in the order they appear in the text or alphabetically. Examples (please note order, style and punctuation):

For books: Obreshkoff, N.: *Higher Algebra*. Nauka i Izkustvo, Second edition, Sofia, 1963 (in Bulgarian).

For journal articles: Frisch, H. L.: Statistics of random media. *Trans. Soc. Rheology*, **9**, 1965, 293–312.

For articles in edited volumes or proceedings: Friedman, H. Axiomatic recursive function theory. In: *Logic Colloquium 95*, (R. Gandy and F. Yates, eds.), North-Holland, 1971, 188–195.

