

<https://doi.org/10.15388/vu.thesis.723>  
<https://orcid.org/0000-0001-5107-530X>

VILNIUS UNIVERSITY

Ugnė Čižikovienė

# Investigations of Binomial Group Testing Models

DOCTORAL DISSERTATION

Natural Sciences,  
Mathematics (N 001)

VILNIUS 2025

This dissertation was written between 2020 and 2024 at Vilnius University.

**Academic supervisor:**

**Assoc. Prof. Dr. Viktor Skorniakov** (Vilnius University, Natural Sciences, Mathematics – N 001).

**Academic consultant:**

**Prof. Habil. Dr. Remigijus Leipus** (Vilnius University, Natural Sciences, Mathematics – N 001).

Dissertation Defence Panel:

Chair – **Prof. Dr. Jurgita Markevičiūtė** (Vilnius University, Natural Sciences, Mathematics – N 001).

Members:

**Prof. Dr. Jonas Šiaulys** (Vilnius University, Natural Sciences, Mathematics – N 001).

**Prof. Habil. Dr. Vydas Čekanavičius** (Vilnius University, Natural Sciences, Mathematics – N 001).

**Assoc. Prof. Dr. Rūta Levulienė** (Vilnius University, Natural Sciences, Mathematics – N 001).

**Prof. Dr. Yaakov Malinovsky** (University of Maryland, Baltimore County, Natural Sciences, Mathematics – N 001).

The dissertation shall be defended at a public meeting of the Dissertation Defence Panel at 27th of January, 2025 at the Institute of Mathematics of Vilnius University. Address: Naugarduko str. 24, Vilnius, Lithuania, tel. +370 5 219 3050; e-mail: mif@mif.vu.lt

The text of this dissertation can be accessed at the Library of Vilnius University and on the website of Vilnius University:

[www.vu.lt/lt/naujienos/ivykiu-kalendorius](http://www.vu.lt/lt/naujienos/ivykiu-kalendorius).

<https://doi.org/10.15388/vu.thesis.723>  
<https://orcid.org/0000-0001-5107-530X>

VILNIAUS UNIVERSITETAS

Ugnė Čižikovienė

# Binominių Grupinio Testavimo Modelių Tyrimai

**DAKTARO DISERTACIJA**

Gamtos mokslai,  
matematika (N 001)

VILNIUS 2025

Disertacija rengta 2020 – 2024 metais Vilniaus universitete.

**Mokslinis vadovas:**

**doc. dr. Viktor Skorniakov** (Vilniaus universitetas, gamtos mokslai, matematika – N 001).

**Mokslinis konsultantas:**

**prof. habil. dr. Remigijus Leipus** (Vilniaus universitetas, gamtos mokslai, matematika – N 001).

Gynimo taryba:

Pirmininkė – **prof. dr. Jurgita Markevičiūtė** (Vilniaus universitetas, gamtos mokslai, matematika – N 001).

Nariai:

**prof. dr. Jonas Šiaulys** (Vilniaus universitetas, gamtos mokslai, matematika – N 001).

**prof. habil. dr. Vydas Čekanavičius** (Vilniaus universitetas, gamtos mokslai, matematika – N 001.)

**doc. dr. Rūta Levulienė** (Vilniaus universitetas, gamtos mokslai, matematika - N 001).

**prof. dr. Yaakov Malinovsky** (University of Maryland, Baltimore County, gamtos mokslai, matematika – N 001).

Disertacija ginama viešame Gynimo tarybos posėdyje 2025 m. sausio 27 d. Vilniaus universiteto Matematikos ir informatikos fakulteto Matematikos instituto 201 auditorijoje. Adresas: Naugarduko g. 24, Vilnius, Lietuva, tel. +370 5 219 3050; el. paštas: [mif@mif.vu.lt](mailto:mif@mif.vu.lt).

Disertaciją galima peržiūrėti Vilniaus universiteto bibliotekoje ir Vilniaus universiteto interneto svetainėje adresu: [www.vu.lt/lt/naujienos/ivykiu-kalendorius](http://www.vu.lt/lt/naujienos/ivykiu-kalendorius).

# Table of Contents

<b>LIST OF ABBREVIATIONS</b>	<b>7</b>
<b>Notation</b>	<b>8</b>
<b>1 Introduction</b>	<b>9</b>
1.1 Topic and early history of Group Testing . . . . .	9
1.2 Examples of applications . . . . .	10
1.3 Classification of Group Testing methods . . . . .	12
1.4 Structure of the Thesis . . . . .	13
<b>2 Background</b>	<b>15</b>
2.1 Binomial Testing Assumptions . . . . .	15
2.2 Several Group Testing procedures . . . . .	16
2.3 Typical Group Testing tasks . . . . .	24
<b>3 Results</b>	<b>26</b>
3.1 Optimal configurations for Modified Dorfman, Sterrett, and Square Array procedures . . . . .	26
3.1.1 Theoretical results . . . . .	26
3.1.2 Examples . . . . .	29
3.2 Algorithm for finding optimal cut-point . . . . .	33
3.2.1 Theoretical results . . . . .	33
3.2.2 Examples . . . . .	36
3.3 Probabilistic analysis of the Pairwise testing procedure . . . . .	43
<b>4 Discussion and conclusions</b>	<b>46</b>
4.1 About results of Section 3.1 . . . . .	46
4.2 About results of Section 3.2 . . . . .	47
4.3 About results of Section 3.3 . . . . .	49
4.4 Several concluding remarks . . . . .	50
<b>5 Proofs</b>	<b>52</b>
5.1 Proofs of results stated in Section 3.1 . . . . .	52

5.2 Proofs of results stated in Section 3.2 . . . . .	67
5.3 Proofs of results stated in Section 3.3 . . . . .	69
<b>Bibliography</b>	<b>78</b>
<b>Santrauka (Summary in Lithuanian)</b>	<b>87</b>
Tyrimų sritis . . . . .	87
Tyrimų objektas ir uždaviniai . . . . .	88
Savokos, žymenys ir prielaidos . . . . .	88
Nagrinėti modeliai . . . . .	89
Disertacijoje spręsti Grupinio testavimo uždaviniai . . . . .	92
Disertacijos struktūra . . . . .	94
Rezultatų apžvalga . . . . .	94
Optimalios Modifikuotos Dorfman, Sterrett ir Kvadra-tinės matricos procedūrų konfigūracijos . . . . .	94
Optimalaus tikimybinio slenksčio paieškos algoritmas . . . . .	97
Porinės testavimo procedūros analizė . . . . .	100
Išvados . . . . .	103
Rezultatų naujumas . . . . .	103
Aprobacija . . . . .	104
Publikacijos . . . . .	104
Trumpos žinios apie autore . . . . .	105
<b>Declarations</b>	<b>108</b>
<b>Publications by the Author</b>	<b>109</b>

# LIST OF ABBREVIATIONS

Abbreviation	Explanation
a.s.	Almost surely
A2	Square Array procedure
BTA	Binomial testing assumptions
BGT	Binomial group testing
$Be(p)$	Bernoulli distribution with parameter $p$
CGT	Combinatorial group testing
CLT	Central limit theorem
COCP	Continuous scale optimal cut-point
DOCP	Discrete scale optimal cut-point
D	Dorman procedure
GT	Group testing
H	Halving procedure
LDP	Large deviation principle
LLN	Law of large numbers
lhs	Left-hand side
MD	Modified Dorman procedure
MGF	Moment generating function
OCP	Optimal cut-point
PGT	Probalistic group testing
PT	Pairwise testing procedure
QC	Quality control
rhs	Right-hand side
ST	Sterrett procedure
s.t.	Such that
UCP	Universal cut-point
w.r.t	With respect to

# Notation

Notation	Explanation
$G_X$	Average testing savings per item (aka gain)
$L^X$	Loss function
$N$	Cohort size
$N_{opt}^X$	Optimal value of the tested cohort size when applying procedure $X$ for the case of a known prevalence $p$
$N_*^X$	Optimal value of the tested cohort size when applying procedure $X$ for the case of an unknown prevalence $p$
$\dot{N}$	Derivative of a latent function $N$
$T_X$	Random number of tests required to identify all defectives when applying procedure $X$
$\theta_X$	Average number of tests required to identify all defectives when applying procedure $X$
$t_X$	Average number of tests per single item in a cohort when applying procedure $X$
$p_{X,c}$	Optimal continuous scale cut-point for procedure $X$
$p_c^X$	Optimal (discrete scale) cut-point for procedure $X$

# Chapter 1

## Introduction

### 1.1 Topic and early history of Group Testing

In many real-life settings, the task of testing is frequent. For example, physicians routinely apply tests to screen for diseases, manufacturers test goods to identify defective items, statisticians utilize tests to validate hypotheses, etc. Group Testing (GT), also known as Pool Testing, refers to a set of methods devoted to this task. The distinguishing feature of any GT method is an attempt to save testing efforts by replacing single items' tests with tests of groups of items. Hence the name.

Today, more than 80 years have passed since the recorded historical beginning of GT. It is usually associated with a seminal paper [19] published by Robert Dorfman in 1943. In that paper, Dorfman addressed the following problem. During World War II, the US Army conscripts had to be tested for syphilis. Syphilis required an accurate blood test, known nowadays as the Wassermann test, which was carried out after a sample of each soldier's blood had been taken. Under usual circumstances, syphilis was a rare disease, and the vast majority of the tests carried out came back negative. Therefore, Dorfman pointed out that the total number of tests for syphilis testing could be significantly reduced by pooling samples. This meant that blood samples from a large number of soldiers had to be taken and mixed into a single pooled sample so that the initial syphilis test was not performed on each soldier's blood sample. If the result was negative for the pooled sample, all soldiers whose blood samples were included in the pooled sample were healthy and free of syphilis. If the test was positive, one knew that at least one of the soldiers in the group had syphilis, and the whole cohort needed

to be retested one by one. This first GT model demonstrated how to save by testing groups. Unfortunately, as reported in [20], it was not put into practice then, yet it inspired further investigations in this direction. These investigations were not immediate since World War II ended, the need for GT for massive testing disappeared, and the scientific community put it aside for quite a long time. In 1957, Sterret published a paper [73] describing a new testing scheme<sup>1</sup> which has some advantages compared with the original Dorfman's, and it took two years more for the seminal work by Sobel and Groll [72], the two Bell Laboratories Scientists, to make a break-through. This long (74 page) paper published in 1959 spawned a sequence of important GT-related papers to appear in 1960's [24, 25, 70, 71, 79] and to set in motion GT field afresh.

## 1.2 Examples of applications

Although the concept of group testing was first formulated in the context of medical testing of patients, it soon became apparent that use cases are not limited to this field alone. Below, we provide a sample of various applications illustrating its spread. The list is far from exhaustive yet sufficient to give an impression of breadth.

**Medicine and biology** It is not surprising that GT counts down a vast number of applications devoted to screening for a particular infectious disease like HIV [58, 61, 65], hepatitis B [15, 23, 28, 60], and most recently COVID-19 [13, 47, 49, 52, 64]. There are, however, many others.

In DNA testing, one looks for true genomic sub-sequences in relatively short fragments of DNA. The possibility of mixing samples taken from patients means that one can apply group testing, significantly reducing the number of tests and, hence, costs. The topic was so important that Du and Hwang (two famous investigators in the GT field) have published a dedicated monograph [41].

Another application in biology is the counting of infected/exposed objects. Not all cases require the identification of every item that is infected. Often, only the proportion of infected matters. For example, when monitoring the spread of diseases while preserving the confidentiality of individuals [68], or inspecting the spread of diseases among a wide variety of insects [21].

---

<sup>1</sup>it will be explicated in the sequel

**Communications and networking** Multi-access channels are channels that can be accessed, communicated, and messaged by multiple users. The basic problem here is identifying the users having information to transmit (active users). Hayes [35] was among the first who discovered that the GT approach is applicable in this context and can be successfully applied for the identification of active users. Several subsequent works were discussed by Wolf [82], who related the communications community's efforts to employ GT to the general GT theory developed so far. Among many investigations after that, paper [50] by Goodrich and Hirschberg is worth mentioning. Similar to predecessors, they studied GT algorithms for resolving broadcast conflicts on multi-access channels and for identification of the dead sensors in a mobile and wireless network. Their approach enriched the standard GT model, allowing the result of each test to be non-binary and indicate the number of defective items contained in the tested subset.

**Cybersecurity, database systems, data compression** An important cybersecurity problem is to efficiently determine which files on a computer/distributed storage have changed. In the general GT context, the changed files correspond to the defective items. M. T. Goodrich, M. J. Atallah, R. Tamassia [31] and T. Madej [51] proposed GT procedures devoted to solving problems of this kind.

Another cybersecurity problem is the detection of various virtual attacks. Khattab et al. [44], Xuan et al. [87] and Gurani et al. [32] described how GT could be used to detect denial-of-service attacks. They divided a server into several virtual servers and monitored which received high traffic. This way, the highest traffic users are identified.

For effective database management, it may be helpful to classify items as having high demand. Cormode and Muthukrishnan [26] suggested how to use GT for achieving such item classification.

Hong and Ladner [36] described an adaptive algorithm for image data compression, whereas Shavit et al. [67] did this for video data compression.

**Relation to Theoretical Computer Science** There are many works demonstrating a clear relationship between GT and Theoretical Computer Science. For example, Hwang and coauthors in a series of papers [27, 40, 85] employ the fact that GT procedures for finding a single defective can be cast in terms of optimal binary search trees (aka Hu-Tucker

trees). Chen et al. [11] propose GT procedures based on the Shannon-entropy criteria, whereas Hsu [37] develops GT procedures based on the Huffman lower bound and Shannon-entropy criteria. Triesch [78] and later Allemann [6] employ hypergraphs to construct very effective GT procedures. Among many others, Aldridge's paper [4] delineates a relationship with coding theory, whereas his recent review with coauthors [5] provides a systematic view of GT from the informational perspective and contains many examples of applications.

**Quality control** The field of Quality Control (QC) is rich in GT applications, pretty much like that of medicine and biology. We do not list here specific ones yet point out the monograph [42] and review [75]. Inspecting the references therein, the reader can find many various GT procedures related to the field of QC.

### 1.3 Classification of Group Testing methods

There are several ways to classify GT procedures. The first and most common is to distinguish between combinatorial GT (CGT) and probabilistic GT (PGT). In CGT, one assumes that, given a set of  $N \in \mathbb{N}$  objects, there can be at most or exactly  $1 \leq d \leq N$  defectives<sup>2</sup>. In order to identify these defective items and reduce the number of tests utilized, combinatorial methods are applied. In CGT, the worst case scenario (i.e., the one requiring the largest possible number of tests) is of primary interest. In PGT, one assumes that each item to be classified has a probability of defectiveness and also specifies dependence between items. This way, the data-generating mechanism is introduced, and because of randomness, the focus usually lies in reducing an average number of tests by taking into account the data-generating mechanism.

The second way to classify GT procedures is to distinguish between adaptive ones and non-adaptive ones. In adaptive procedures, testing is done sequentially, and one can utilize the results of the tests carried out so far in the subsequent tests. In non-adaptive procedures, all tests usually run in parallel, and there is no opportunity to make use of accrued information. Hybrid procedures are also met in the literature.

---

<sup>2</sup>or else contaminated, infected, etc.; in all the Thesis, we use these adjectives as well as their counterparts (non-defective, pure, non-infected, etc.) interchangeably; in all cases, they mean the same: an object or a group of objects having (respectively, not having) property of interest

There are several other criteria employed in the classification of GT procedures. Some of them apply only to specific sub-classes (e.g., only to adaptive procedures). Below, we list a few.

*Test kit quality.* In many practical applications, it is unreasonable to assume that the test kit at hand obeys 100 percent sensitivity and/or specificity. Therefore, one can split the GT procedures into those assuming perfect testing and those assuming testing with errors.

*Cohort size.* Often there are natural constraints preventing from testing arbitrarily large groups. E.g., in medical screening, it is common to observe the so-called dilution effect: pooling of too many specimens inflates the test's sensitivity/specificity drastically; thus, only procedures with upper bounded cohort size are applicable. It is, therefore, possible to classify GT procedures into size-constrained/unconstrained.

*Nesting.* One says that an adaptive GT procedure belongs to the class of nested procedures provided the following holds: at each testing stage (except the first one), the cohort to be tested next is a subset of a contaminated items' set, i.e., the set known to contain defectives. The class of nested procedures is very large and common in applications.

*Test output.* The procedures can also be classified by their output. Most common scenario is that of binary procedures: given the set to test, such procedure outputs 1 provided at least one item in the tested set is contaminated and it outputs 0 provided all items are pure; there are no further indications which items (if any) are contaminated and how much items of this kind there are in the tested set. However, some procedures can output the (estimated) number of contaminated items.

In this Thesis, we investigate only size unconstrained binary PGT procedures corresponding to perfect testing and satisfying Binomial Testing Assumptions described in Subsection 2.1. Almost all of them are nested.

## 1.4 Structure of the Thesis

The remaining part of the thesis contains three chapters. In Chapter 2, we give notation and provide all necessary background. It includes

description of:

- (a) general Binomial Testing Assumptions under which all subsequent analysis is done;
- (b) several GT procedures, which are repeatedly used in the sequel;
- (c) general GT results required for understanding of our ones;
- (d) typical GT analysis tasks.

Chapter 3 is devoted to the statement of our results, some examples of applications, and a related literature review. Corresponding proofs are given in Chapter 5. Chapter 4 is devoted to discussion. Finally, there is a summary in Lithuanian placed at the very end.

# Chapter 2

## Background

### 2.1 Binomial Testing Assumptions

As mentioned in the Introduction 1, we consider only PGT models. However, even here, the most general model is intractable analytically and concrete GT procedures are formulated under additional assumptions. Before proceeding to their statement, we remind that, in the Thesis, we consider only binary tests: when applied to the group at hand, such test outputs one if the group contains at least one contaminated item, and it outputs zero provided the group is pure. Generally,  $N$  stands for the size of the initial cohort of the items to be tested. It can be further parameterized by an integer  $n \geq 1$ , and such parametrization (including a trivial one  $N(n) = n$ ) depends on the algorithm of the procedure under consideration (see Subsection 2.2 for concrete examples).

In what follows, we always assume that:

- (BTA1) initially, each item in the cohort to be tested can be contaminated with the same constant probability  $p \in (0, 1)$ ;
- (BTA2) items are independent;
- (BTA3) test kit under consideration is perfect (i.e., sensitivity = specificity = 100%) and does not depend on the size of the tested group.

Assumptions (BTA1)–(BTA3) are called Binomial Testing Assumptions (BTA). Imposing them, we can formalize our setup by the following probabilistic model:

- Let  $Y_1, \dots, Y_N$  be i.i.d.  $\sim Be(p)$  random variables (r.vs.) representing the cohort of interest;
- $Y_i = 1 \iff$  item  $i$  is contaminated (then  $Y_i = 0 \iff$  item  $i$  is pure);
- For any testing procedure  $X$  and any non-empty  $A \subset \{Y_1, \dots, Y_N\}$ ,  $X(A) = \mathbb{1}\{\sum_{Y_i \in A} Y_i > 0\}$ .

Given cohort  $C = \{Y_1, \dots, Y_N\}$  and test procedure  $X$ , let  $T_X$  denote the random number of tests required to identify all contaminated in  $C$  and let  $\theta_X(N, p) = E T_X$  be its average. Note that  $\theta_X$  is a function of two arguments:  $N \in \mathbb{N}, p \in (0, 1)$ . For convenience,  $q$  stands for  $1 - p$  in the entire Thesis.

## 2.2 Several Group Testing procedures

In this subsection, we describe several binomial GT procedures serving as examples and objects of investigation in the sequel. The graphic schemes of the procedures are presented in figures 2.1–2.5.

**Dorfman procedure D** This procedure was already described in the Introduction 1. It spans two steps:

Step1: test initial pooled sample  $IP$ ;

Step2: if  $IP$  tests negative, finish; otherwise retest each item individually.

One can see that  $T_D = 1 + N\mathbb{1}\{Y_1 + \dots + Y_N > 0\}$ . Therefore,

$$\theta_D(N, p) = 1 + N E \mathbb{1}\{Y_1 + \dots + Y_N > 0\} = 1 + N(1 - q^N). \quad (2.1)$$

**Modified Dorfman procedure MD** Sobel and Groll [72] observed that the D procedure is inconsistent and suggested a fix: if the initial pool tested positively and having retested  $N - 1$  items there were no defectives, there is no need to retest the last one. In what follows, we denote this procedure MD. By description,

$$T_{MD} = 1 + (N - 1)\mathbb{1}\{Y_1 + \dots + Y_N > 0\} + \mathbb{1}\{Y_1 + \dots + Y_{N-1} > 0\}.$$

Therefore,

$$\theta_{MD}(N, p) = 1 + (N - 1)(1 - q^N) + 1 - q^{N-1} = 1 - pq^{N-1} + N(1 - q^N). \quad (2.2)$$

**Sterrett procedure ST** Sterrett [73] proposed another modification of the D procedure. His testing algorithm is as follows:

Step1: test initial pooled sample  $IP$ ;

Step2: if  $IP$  tests negative, finish; otherwise go to Step3;

Step3: retest each item one by one until the first contaminated is identified; if the whole cohort is already tested, finish; otherwise, treating the set of the remaining untested items as a new initial cohort, go to Step1 and start over.

In [72] it was shown that

$$\theta_{ST}(N, p) = 2q - p^{-1}(1 - q^{N+1}) + (2 - q)N. \quad (2.3)$$

**Pairwise testing algorithm PT** Pairwise testing algorithm was investigated by Yao and Hwang in [85]. Its description is as follows:

Step1: if there is only one item in the cohort  $C$ , test it and finish; if there are no items at all, finish; otherwise proceed to Step2;

Step2: having a cohort of  $N \geq 2$ , choose a pair of items, form a pool, and test it;

Step3:

- if the pool tests negative, classify tested items as pure, set  $N = N - 2, C = C \setminus \{\text{tested pair}\}$ ;
- if the pool tests positive, retest one item; in case that item tests negative, classify the remaining item as contaminated, set  $N = N - 2, C = C \setminus \{\text{tested pair}\}$ ; if that item tests positive, set  $N = N - 1, C = C \setminus \{\text{retested item}\}$ .

Step4: go to Step1.

**Remark 2.2.1.** There is no difference which pair of items to choose in Step 2. In what follows, we always choose the last two.  $\square$

In [85] it was demonstrated that

$$\theta_{PT}(N, p) = N \frac{2 - q^2}{1 + q} + \frac{q^2 + q - 1}{(1 + q)^2} (1 - (-q)^N). \quad (2.4)$$

**Square Array procedure A2** This procedure was introduced by Phatarfod and Sudbury [63] and later generalized by Berger, Mandell and Subrahmanyam [9]. To apply it, one has to have  $N = n^2, n \in \mathbb{N}$ , items in total. Assuming this holds true, the algorithm reads as follows:

Step1: put samples of all items on a square  $n \times n$  matrix;

Step2: make  $n$  pools corresponding to rows and  $n$  pools corresponding to columns; test them;

Step3: if all pools test negatively, finish; otherwise retest items  $I_{ij}$  satisfying condition "row  $i$  and column  $j$  tested positively".

An average number of tests was computed by Phatarfod and Sudbury [63]:

$$\begin{aligned} \theta_{A2}(N, p) &= 2n + n^2 (1 - 2q^n + q^{2n-1}) = \\ &2\sqrt{N} + N \left(1 - 2q^{\sqrt{N}} + q^{2\sqrt{N}-1}\right). \end{aligned} \quad (2.5)$$

**Halving procedure H** It is difficult to trace back who has introduced this procedure first as it has tight relationships to binary search problem. For the purposes of GT in quality control, it was discussed by Johnson et. al. [42]. To apply it, one has to have  $N = 2^n, n \geq 1$  items in total. As the name suggests, the corresponding algorithm works by testing halves. Below comes a precise description.

Step1: Test initial pool. If it tests negative, finish; otherwise, go to Step 2.

Step2: Split the initial cohort into two equal subsets. Apply the procedure recursively (i.e., starting from Step 1) to each half.

An average number of tests required by this procedure (see, e.g., [69]) equals

$$\theta_H(N, p) = 1 + 2N \sum_{k=1}^{\log_2 N} \frac{1 - q^{2^k}}{2^k}. \quad (2.6)$$

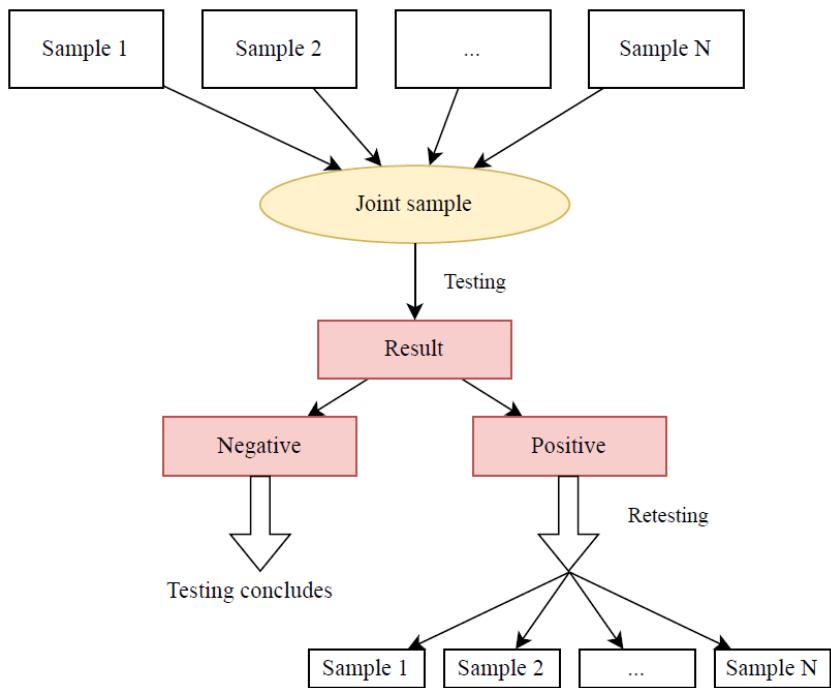


Figure 2.1: Dorfman testing scheme [19].

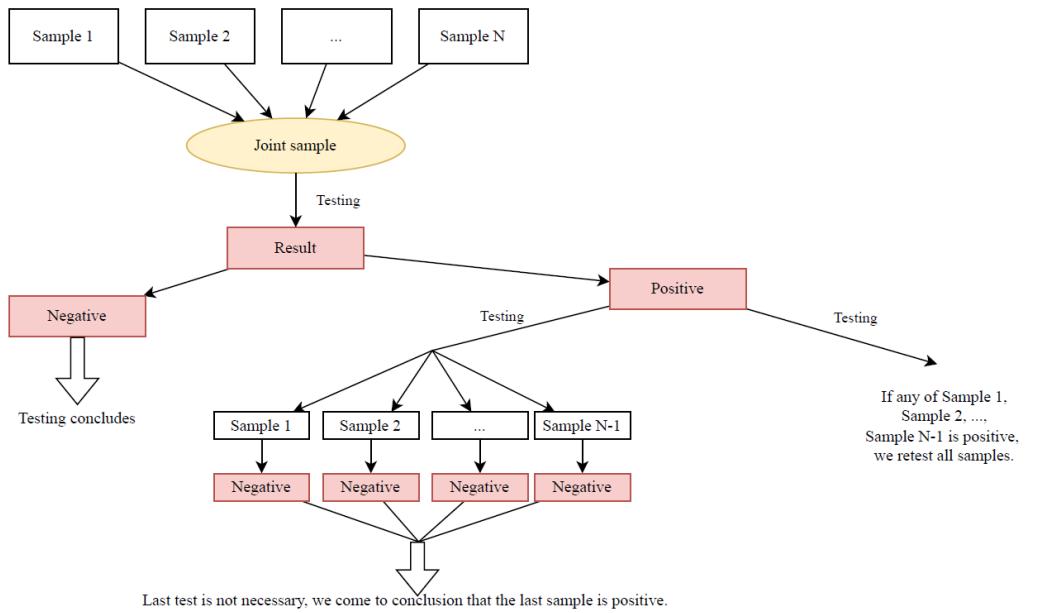


Figure 2.2: Modified Dorfman testing scheme [72].

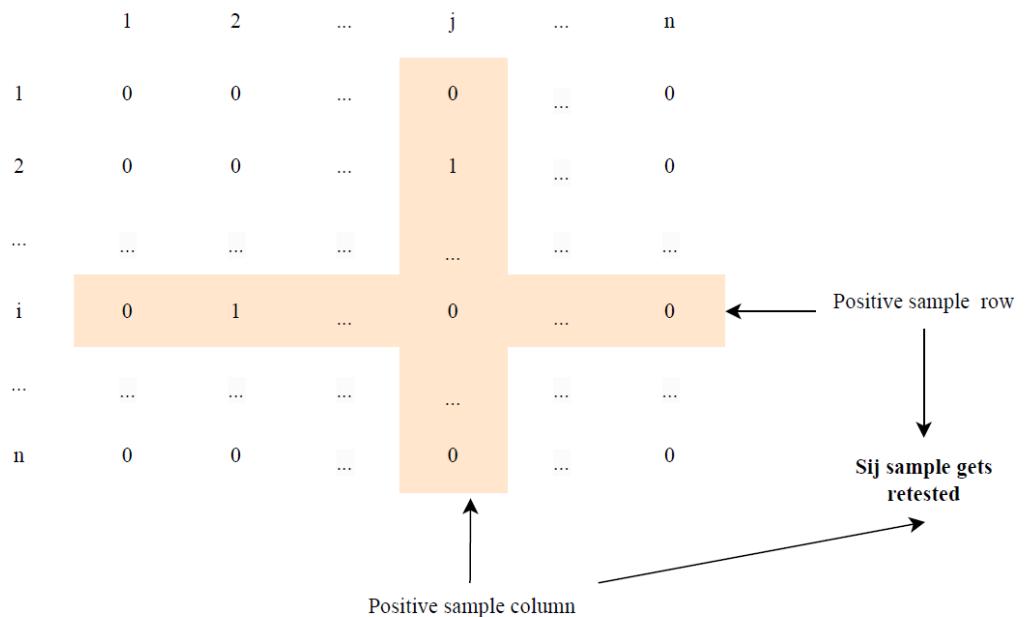


Figure 2.3: Square Array testing scheme [63].

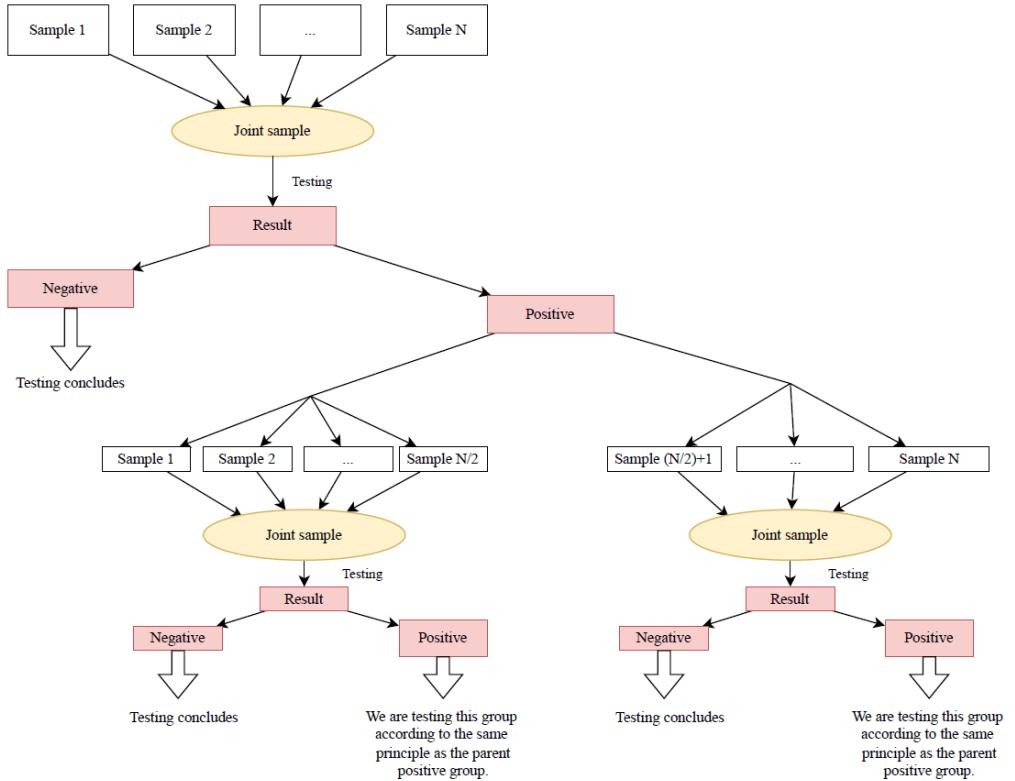


Figure 2.4: Halving testing scheme.

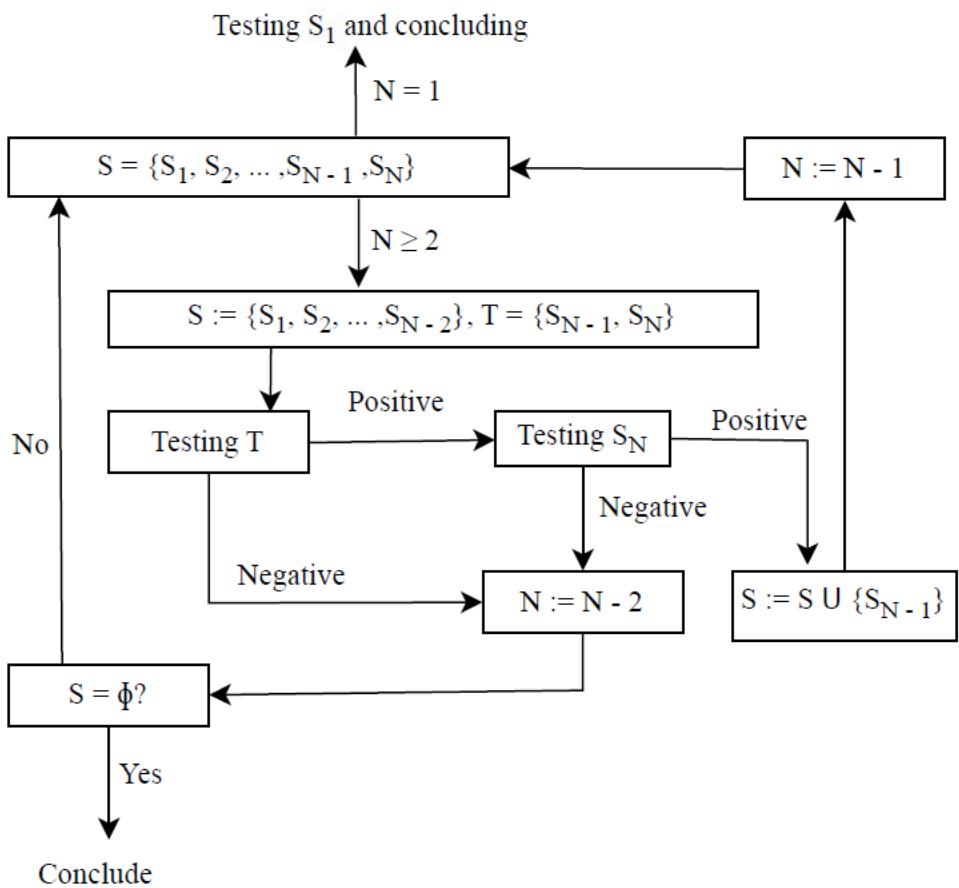


Figure 2.5: Pairwise testing scheme [85].

## 2.3 Typical Group Testing tasks

Let  $X$  be a procedure of interest. It is reasonable to assume that it differs from other procedures not only w.r.t. the testing scheme but also w.r.t. an average number of tests  $\theta_X(N, p)$ . Since the primary aim of GT is to minimize this quantity and the contamination probability  $p$  may be regarded as fixed by nature, one of the typical GT tasks is to find the optimal value of the tested cohort size  $N_{opt}^X$ . To state this optimization problem formally, one usually considers the function

$$\mathbb{N} \times (0, 1) \ni (N, p) \mapsto t_X(N, p) := \frac{\theta_X(N, p)}{N} \quad (2.7)$$

whose value is an average number of tests per single item in a cohort of size  $N$ . Then, any minimizer of this function w.r.t.  $N$  is called optimal configuration and denoted by  $N_{opt}^X$ . Note that we assume fixed  $p$ , i.e.  $N_{opt}^X \in \arg \min_{N \in \mathbb{N}} t_X(N, p)$ . Therefore,  $N_{opt}^X$  depends on  $p$ . This has several important consequences discussed below.

The function  $(0, 1) \ni p \mapsto N_{opt}^X(p)$  shows how large should be the sizes of the tested cohorts to have average testing savings per item (aka gain)

$$G_X(p) := 1 - t_X(N_{opt}^X(p), p) \quad (2.8)$$

maximal (in the long run). It is intuitively clear that, for any reasonable  $X$ ,  $(0, 1) \ni p \mapsto N_{opt}^X(p)$  should be non-increasing, and in all the remaining parts of the Thesis, we assume this by default. Then values of  $G_X(p)$  close to 1 should be observed for small  $p$ 's. In practice, testing pools of arbitrarily large sizes may be infeasible because of the dilution effect: large pool sizes change test kit characteristics (usually, sensitivity drops down drastically). Knowing the maximal tolerable pool size  $N_{\max}$  of the test kit at hand allows one to find the range  $R_X$  of  $p$ 's for which  $N_{opt}^X(p) \leq N_{\max}$ . For this, explicit functional form of  $(0, 1) \ni p \mapsto N_{opt}^X(p)$  is desirable.

An explicit form of the latter function yields an ability to address another typical GT problem: a comparison of several competing procedures. Say,  $X_i, i = 1, \dots, k$ , are procedures applicable to the problem at hand. One can investigate both gains  $G_{X_i}(p)$  (see Eq. (2.8)) and sizes  $N_{opt}^{X_i}(p)$  for different values of  $p$  and, taking into account various scenarios, choose preferable procedures.

Finally, there is one more typical GT task tightly related to  $N_{opt}^X$ : finding optimal cut-point (OCP)  $p_c^X \in (0, 1)$ . As noted earlier, intuition

suggests that group testing should be profitable only for sufficiently small prevalence  $p$ . To be more specific, recall Dorfman procedure D: if  $p$  is small, then testing will oftentimes end up after a single initial pool test; however, if  $p$  is close to 1, almost every group of size  $N \geq 2$  will be contaminated, and in addition to single pool test there will be  $N$  individual tests; hence, instead of having only  $N$  individual tests, one will have  $N + 1$  tests. The above OCP problem therefore translates to finding value  $p_c^X \in (0, 1)$  s.t.  $\forall p \in (p_c^X, 1) N_{opt}^X(p) = 1$ . In the general BTA context, this problem was addressed by Ungar. In [79], he obtained the following fundamental GT result.

**Theorem 2.3.1.** *Binomial GT makes sense if and only if  $p \in \left(0, \frac{3-\sqrt{5}}{2}\right)$ : for  $p \in \left(\frac{3-\sqrt{5}}{2}, 1\right)$ , there does not exist GT procedure performing better than individual testing; for  $p \in \left(0, \frac{3-\sqrt{5}}{2}\right)$  there always exists at least one procedure  $X$  s.t.  $t_X(N_{opt}^X(p), p) < 1$ .*

In the binomial GT, the value  $\frac{3-\sqrt{5}}{2}$  is known under the name of Universal Cut-Point (UCP). For some GT procedures, it is achievable. That is,

$$G_X(p) = 0 \iff p \geq \frac{3 - \sqrt{5}}{2}.$$

However, there are procedures for which

$$G_X(p) = 0 \iff p \geq p_c^X \text{ with } p_c^X < \frac{3 - \sqrt{5}}{2}.$$

Therefore, when investigating characteristics of a particular procedure  $X$ , it is first of all important to find  $p_c^X$ : if the problem at hand is s.t.  $p > p_c^X$ , one needs to resort to individual testing or look for another GT procedure if  $p_c^X < \frac{3-\sqrt{5}}{2}$ . Because of Theorem 2.3.1, in all the Thesis, we restrict the range of  $p$  to  $\left(0, \frac{3-\sqrt{5}}{2}\right)$ .

In the end, we mention that there are other binomial GT problems not discussed here. Some of them are difficult and open for quite a long period of time. In our opinion, those described above are the most typical. Nonetheless, dealing with a particular procedure still poses a challenge and requires extensive analysis.

# Chapter 3

## Results

### 3.1 Optimal configurations for Modified Dorfman, Sterrett, and Square Array procedures

#### 3.1.1 Theoretical results

In Chapter 2, we have mentioned that finding  $N_{opt}^X$  is one of the typical yet very important GT tasks because knowledge of  $N_{opt}^X$  allows to utilize  $X$  most efficiently. As we know, there are no a lot of results providing explicit expressions of  $N_{opt}^X$ .

For the case of the Dorfman procedure D, the optimal configuration was derived by Samuels [66] quite long ago. He has demonstrated that, for<sup>1</sup>  $p \in (0, 1 - (1/3)^{1/3}) \approx (0, 0.31)$ ,  $N_{opt}^D(p) \in \{\lfloor \sqrt{p^{-1}} \rfloor + 1, \lfloor \sqrt{p^{-1}} \rfloor + 2\}$  whereas, for  $p \geq 1 - (1/3)^{1/3}$ ,  $N_{opt}^D(p) = 1$  (the latter result also means that optimal cut-point  $p_c^D = 1 - (1/3)^{1/3}$  is strictly less than  $UCP = \frac{3-\sqrt{5}}{2}$  discussed in Subsection 2.3). However, at the start of our investigations, we discovered that for the modified Dorfman procedure MD, the Sterrett procedure ST, and the square array procedure A2, the analytical expressions of  $N_{opt}^X$  were not completely made explicit. To be more precise, for MD and ST procedures, Malinovsky and Albert [55] numerically derived explicit analytical expressions of optimal configurations for a wide range of  $p$ 's and conjectured that the same expressions are valid

---

<sup>1</sup>here and in the sequel,  $\lfloor x \rfloor$  denotes an integer part of  $x \in \mathbb{R}$ ;  $\lceil x \rceil$  equals  $x$  if  $x \in \mathbb{Z}$  and it equals  $\lfloor x \rfloor + 1$  otherwise

for all  $p \in (0, UCP)$ . In our paper [88], we justified the correctness of their conjectures and proved the following statements.

**Theorem 3.1.1.** *Let*

$$g_0(p) := \frac{1}{q} \left( \frac{1 - 2pq}{q \left( 1 - \ln q \sqrt{\frac{2}{p}} \right)} \right)^{\sqrt{\frac{p}{2}}} \quad \text{for } p \in \left( 0, \frac{3 - \sqrt{5}}{2} \right). \quad (3.1)$$

*Then equation  $g_0(p) = 1$  admits a unique solution  $p_* \approx 0.1711$  and the following relations hold:*

$$\begin{aligned} N_{opt}^S(p) &\in \left\{ \lfloor \sqrt{2p^{-1}} \rfloor, \lfloor \sqrt{2p^{-1}} \rfloor + 1 \right\} \text{ for } p \in \left( p_*, \frac{3 - \sqrt{5}}{2} \right); \\ N_{opt}^S(p) &\in \left\{ \lfloor \sqrt{2p^{-1}} \rfloor, \lfloor \sqrt{2p^{-1}} \rfloor + 1, \lfloor \sqrt{2p^{-1}} \rfloor + 2 \right\} \text{ for } p \in (0, p_*]. \end{aligned}$$

**Theorem 3.1.2.** *For all  $p \in \left( 0, \frac{3 - \sqrt{5}}{2} \right)$ ,*

$$N_{opt}^{MD}(p) \in \left\{ \lfloor \sqrt{p^{-1}} \rfloor, \lfloor \sqrt{p^{-1}} \rfloor + 1 \right\}. \quad (3.2)$$

The A2 procedure was investigated by Hudgens and Kim [38]. They derived quite tight lower and upper bounds yet did not succeed in providing exact expression of optimal A2 configuration and left explicit formulae undiscovered. They also addressed the problem of finding OCP (see Section 2.3) and actually found optimal value  $p_c^{A2}$ . In [89], we obtained expression for  $N_{opt}^{A2}(p)$  and along the way produced additional insights about  $p_c^{A2}$  (see clarifying Remark 3.1.5 below). Before stating our results, we remind that, in the case of A2, cohort size  $N$  is parameterized by  $n \in \mathbb{N}$  as follows<sup>2</sup>:  $N(n) = n^2$ . Also, in the setting of A2, it was more convenient for us to treat  $t_{A2}$  as a function  $[2, \infty) \times (UCP, 1) \ni (n, q) \mapsto t_{A2}(n, q)$  with  $n \in [2, \infty)$  ranging continuously rather than a function of  $(N, p)$  with a domain  $\{n^2 : n \in \mathbb{N}\} \times (0, UCP)$ . Keeping this in view, our results read as follows.

**Theorem 3.1.3.** *Let  $g(q, n) = \frac{2}{n} - 2q^n + q^{2n-1} = t_{A2}(n, q) - 1$ . The following statements hold true:*

---

<sup>2</sup>see Section 2.2 for the description of A2

(i) For any  $(q, n) \in (1/2, 1) \times (2, \infty)$ , system of equations

$$\begin{cases} 1 = nq^n \left(1 - \frac{q^{n-1}}{2}\right) \\ n \ln q = -\frac{\left(1 - \frac{q^{n-1}}{2}\right)}{(1-q^{n-1})} \end{cases} \quad (3.3)$$

has a unique solution  $(q_*, n_*) \approx (0.748416, 4.453524)$ .

- (ii) For any fixed  $q \in (q_*, 1)$  and with respect to  $n$ , equation  $g(q, n) = 0$  admits two solutions  $n_L, n_U : 2 < n_L < n_* < n_U < \infty$ . On  $(n_L, n_U)$ ,  $n \mapsto g(q, n)$  attains values in  $(-\infty, 0)$  whereas on  $(2, \infty) \setminus [n_L, n_U]$  it attains values in  $(0, \infty)$ .
- (iii) For any fixed  $q \in (q_*, 1)$ , the region  $(n_L, n_U)$  is the one where A2 is efficient, i.e.,  $t_{A2}(q, n) < 1$  for  $n \in (n_L, n_U)$ . In that region, there exists a unique (and, therefore, global) minimizer  $n_{min}$  of  $(2, \infty) \ni n \mapsto t_{A2}(q, n)$ . For  $q \in [0.755, 1)$ , it is given by

$$n_{min} = \frac{1}{p^{\frac{2}{3}}} + \frac{1}{2p^{\frac{1}{3}}} + 0.2 + 3p^2 + t_* \quad (3.4)$$

for some  $t_* \in [0, 1]$ .

- (iv)  $(2, \infty) \ni n \mapsto t_{A2}(n, q)$  also has a unique (and, therefore, global) maximizer located in the region  $(n_U, \infty)$ . For any fixed  $q \in (0, q_*)$ , A2 is never optimal, i.e.,  $(2, \infty) \ni n \mapsto t_{A2}(n, q)$  attains values in  $(1, \infty)$ .

**Corollary 3.1.4.** Let  $g(q, n)$  be as in Theorem 3.1.3. Then  $g(q, 5) = 0$  has a unique solution  $q_5 \approx 0.750209961$ . For all  $q \in (q_5, 1)$ ,  $n_{opt}(q)$  belongs to the set

$$\left\{ \left\lfloor \frac{1}{p^{\frac{2}{3}}} + \frac{1}{2p^{\frac{1}{3}}} + 3p^2 + 0.2 \right\rfloor + i : i = 0, 1, 2 \right\}. \quad (3.5)$$

**Remark 3.1.5.** Inspecting the statements above, it might be tempting to conclude that our results are not exhaustive: it is unclear what happens in the region  $(q_*, q_5)$  and what is an expression for  $n_{opt}(q)$  there. Applying Theorem 3.1.3 for a fixed  $q \in (q_*, q_5)$ , we have that  $t_{A2}(n_{min}, q) < 1$  with  $n_{min} = n_{min}(q)$  of (iii). However, this  $n_{min}(q) \in (4, 5)$  and  $\min(t_{A2}(4, q), t_{A2}(5, q)) > 1$  when  $q \in (q_*, q_5)$ . We do not provide separate proof of this fact and only mention that technical details can be filled in after inspection of the proofs presented in Chapter 5. Hudgens and Kim [38] operated on the discrete scale and deduced that region  $(q_5, 1)$

is the one where  $\{5, 6, \dots\} \ni n \mapsto t_{A2}(n, q)$  is efficient. Having proved that  $n = 2, 3, 4$  are never optimal, they actually verified that  $(q_5, 1)$  is the region where practical application of A2 makes sense. Therefore, in this direction, our input adds the missing part to the theoretical characterization of A2. However, it is an honest deal to say that theoretical findings about OCP, though interesting, have no additional practical value.  $\square$

### 3.1.2 Examples

In this subsection, we provide two examples illustrating the usefulness of our results. The first of our examples is most interesting in the asymptotic regime as  $p \rightarrow 0+$ . Since in this case Dorfman procedure D and modified Dorfman procedure MD are equivalent, we exclude MD. Also, in both examples, we switch back to our usual convention and treat  $t_X$  as a function of  $(N, p)$ .

The figures and calculations were produced by making use of an open-source computer algebra system SymPy [59] and the Python packages lying in the core kit for scientific numerical programming with Python: NumPy [34], SciPy [80], Matplotlib [39] and Pandas [81], [77]. The same software was employed in the proofs provided in Chapter 5.

**Example 1: Comparison of procedures** To illustrate performance of different procedures, we compare A2, D, S and H in terms of magnitude of optimal configuration  $N_{opt}^X$  and gain (defined by Eq. (2.8)) across the range  $p \in (0, 0.249790)$  where application of all of them seems reasonable<sup>3</sup>. In what follows,  $N_*^X$  stands for the unrounded optimal configuration, i.e. the minimizer of the continuous argument function  $[1, \infty) \ni N \mapsto t_X(N, p)$ . We have chosen to operate on the continuous scale and use  $t_X(N_*^X(p), p)$  instead of  $t_X(N_{opt}^X(p), p)$  since it is much easier to interpret information visually in comparison to the discrete case. Recall that, in the case of A2,  $N = n^2$ , where  $n$  is the number of rows (columns) in the square array used in the definition of A2. In this example (as well as in the second one), this reparametrization remains in force: writing  $t_{A2}(N, p)$ , we actually mean  $t_{A2}(n_{min}(q), q)$  of Theorem 3.1.3 with  $n_{min}$  given by (3.4). From results of [66], [69] and those above

---

<sup>3</sup>the restriction of the range is due to A2 which does not make sense for larger  $p$ 's

it follows that

$$N_*^{A2}(p) = n_{min}^2(p) = p^{-\frac{4}{3}}(1 + o(1)), \quad N_*^D(p) = \sqrt{p^{-1}}(1 + o(1)),$$

$$N_*^S(p) = \sqrt{2p^{-1}}(1 + o(1)), \quad N_*^H(p) = -\frac{1}{2 \log_2 q}(1 + o(1))$$

and that

$$t_{A2}(N_*^{A2}(p), p) = 3p^{\frac{2}{3}}(1 + o(1)), \quad t_D(N_*^D(p), p) = 2\sqrt{p}(1 + o(1)),$$

$$t_{ST}(N_*^S(p), p) = \sqrt{2p}(1 + o(1)),$$

$$t_H(N_*^H(p), p) = -(2p \log_2 p)(1 + o(1))$$

as  $p \rightarrow 0+$ .

Figure 3.1 shows the behavior of  $N_*^X$  and gain  $G_X$  of the four procedures considered. Due to the raise to the square of  $n_{min}$ ,  $N_*^{A2}$  grows to infinity much faster than the counterparts of the remaining schemes, and, because of this, in the top left sub-figure, the range of  $p$  starts quite far from the origin and the accompanying bottom left sub-figure on the log scale is given. The latter clearly depicts the relationships

$$\ln(N_*^D(p)), \ln(N_*^S(p)) \sim -\frac{1}{2} \ln(p),$$

$$\ln(N_*^H(p)) \sim -\ln(p), \quad \ln(N_*^{A2}(p)) \sim -\frac{4}{3} \ln(p), \quad (3.6)$$

following from the formulae given above and clearly showing that the asymptotic slope of  $\ln(N_*^{A2}(p))$  on the  $-\ln p$  scale is the largest one.

Talking about gains depicted in the right top sub-figure, one can see that the D procedure always performs worse than other competitors. However, each of A2, S and H have their own regions where they perform best. The bottom right sub-figure illustrates that, for  $p$  tending to zero, the H procedure's gain growth rate is the biggest one.

**Example 2: Optimal configuration when the prevalence is unknown**  
In reference [54], the authors looked for the pool size leading to optimal testing by making use of procedure D when the prevalence is unknown. They employed two approaches. Both (approaches) were based on the following loss function. Given procedure  $X$ , define

$$L^X(N, p) = t_X(N, p) - t_X(N_{opt}^X(p), p), \quad (3.7)$$

where  $N_{opt}^X = N_{opt}^X(p) \in \arg \min_{N \in \{1, 2, \dots\}} t_X(N, p)$  is the optimal configuration when the prevalence  $p$  is known. It is clear that  $L^X(N, p) \geq$

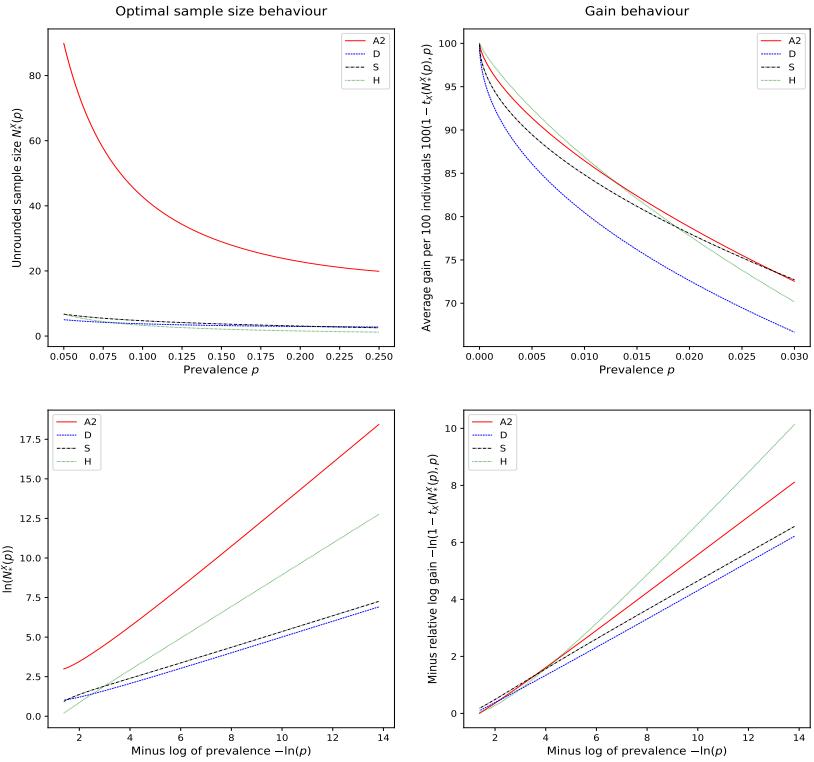


Figure 3.1: Graph showing the behavior of optimal pool sizes and gains on original and log-log scales.

0  $\forall (N, p) \in \mathbb{N} \times (0, 1)$  and, for a given  $p$ ,  $L^X(N, p) = 0$  precisely when  $N \in \arg \min_{N \in \{1, 2, \dots\}} t_X(N, p)$ . In what follows, to distinguish between  $N_{opt}^X(p)$  and optimal configuration suitable for unknown  $p$ 's, the latter configuration is denoted by  $N_\star^X$ .

The first approach in [54] was to make use of mini–max strategy and take  $N_\star^X$  as a minimizer of

$$\{1, 2, \dots\} \ni N \mapsto \sup_p L^X(N, p). \quad (3.8)$$

The second approach was to make use of the Bayesian paradigm and,

after putting the prior  $\pi$  on  $p$ , to take  $N_\star^X$  as a minimizer of

$$\{1, 2, \dots\} \ni N \mapsto \mathbb{E}_\pi(L^X(N, p)) = \mathbb{E}_\pi(t_X(N, p)) - c(\pi), \\ c(\pi) = \mathbb{E}_\pi(t_X(N_{opt}^X(p), p)). \quad (3.9)$$

In this example, we have adopted both approaches to the case of A2. When using the Bayesian one,  $\pi$  was taken uniform over  $(0, 0.249790)$ . Thereby, we have modeled a situation when the only prior information is that the application of A2 makes sense (see Corollary 3.1.4). Also, we have modified (3.9) and used

$$\{1, 2, \dots\} \ni N \mapsto \mathbb{E}_\pi(L^X(N, p))^2 = \mathbb{E}_\pi(t_X(N, p) - t_X(N_{opt}^X(p), p))^2 \quad (3.10)$$

instead. To justify our choice, note that, in (3.9),  $c(\pi)$  does not depend on  $N$ . Therefore, minimization of the target function amounts to minimization of  $\{1, 2, \dots\} \ni N \mapsto \mathbb{E}_\pi(t_X(N, p))$  and the corresponding minimizer depends only on the prior  $\pi$ . This way important information carrying function  $p \mapsto t_X(N_{opt}^X(p), p)$  remains unutilized. It seems, however, more reasonable to look for an estimate minimizing the distance to optimal value function and depending on this function.

Figures 3.2–3.3 show graphs of (3.8) and (3.10) for the case of  $X = A2$  and the previously mentioned prior  $\pi$ . Numerical estimation yielded the following values:

- $N_\star^{A2}(p) = 12^2$  for the case of mini–max approach;
- $N_\star^{A2}(p) = 7^2$  for the case of Bayesian approach.

We did not make any attempt to rigorously prove that these values are the only global minimizers of the target loss.

Finishing the example, it is important to note that, though the strategy discussed above leads to sub-optimal testing in a stable environment when the prevalence is close to constant and its reliable estimation is possible, it appears to be a reasonable strategy when the prevalence is varying rapidly and is difficult to capture by data at hand. Therefore, at least in the initial stage, it can be considered as a good alternative for optimal testing during pandemics like COVID-19. Of course, under such circumstances, one can (and should) use various priors motivated by expert knowledge and/or domain-specific factors.

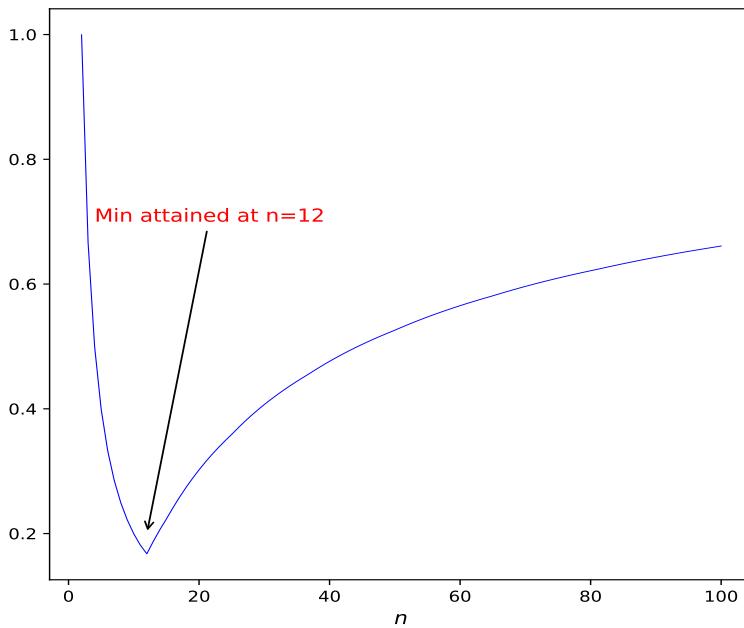


Figure 3.2: Graph of  $\{2, 3, \dots\} \ni n \mapsto \sup_q (t_{A2}(n, q) - t_{A2}(n_{opt}(q), q)$

## 3.2 Algorithm for finding optimal cut-point

### 3.2.1 Theoretical results

In Subsection 2.3, we have discussed the importance of finding OCP for a given procedure  $X$ . However, to our best knowledge, Ungar's Theorem 2.3.1 is the only result of the general nature addressing this problem. All other works were tied to investigations of particular procedures. In our work [90], we proposed an algorithm suitable for finding an approximate value of OCP for a class of BTA satisfying procedures and allowing to recover exact OCP in many cases. Along the way, we have discovered an interesting connection of independent interest between GT and Bifurcation Theory. Our method applies to the class of binomial GT procedures satisfying the following constraints.

(M0)  $\exists c \geq 2$  s.t.  $X$  is a-priori known to be useless for  $N \in [1, c)$ .

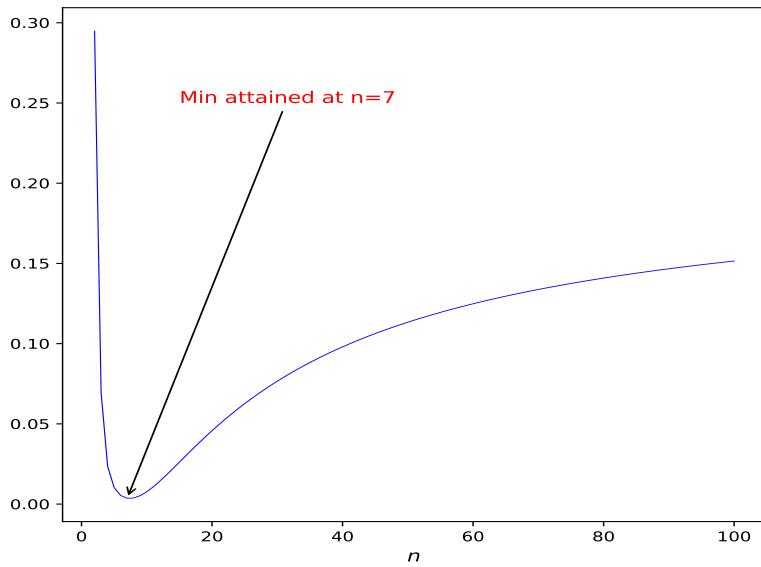


Figure 3.3: Graph of  $\{2, 3, \dots\} \ni n \mapsto E_\pi(t_{A2}(n, q) - t_{A2}(n_{opt}^{A2}(q), q))^2$

(M1) Function  $\mathbb{N} \times (0, UCP] \ni (N, p) \mapsto \theta_X(N, p)$  can be treated as a continuous function on  $[c, \infty) \times (0, UCP]$  differentiable in the whole interior of its domain.

(M2)  $\forall N \in (c, \infty)$  function  $(0, UCP] \ni p \mapsto \theta_X(N, p)$  is strictly increasing.

(M3)  $\forall N \in (c, \infty) t_X(N, UCP) > 1$ .

(M4)  $\forall N \in (c, \infty) \exists p \in (0, UCP) : t_X(N, p) < 1$ .

Our results are given in two propositions stated below. The first one characterizes the properties of the OCP.

**Proposition 3.2.1.** *Assume (M0)–(M4). Let  $p_{X,c} = \sup\{p \in (0, UCP) \mid \exists N \in (c, \infty) : t_X(N, p) < 1\}$ . Then  $\forall p \in (0, p_{X,c})$  procedure X makes sense on the continuous scale;  $\forall p \in (p_{X,c}, UCP]$  it makes no sense at all, that is,*

$$(N, p) \in (c, \infty) \times (p_{X,c}, UCP] \Rightarrow t_X(N, p) > 1. \quad (3.11)$$

The second proposition demonstrates that under (M0)–(M4), there exists a generic procedure for finding  $p_{X,c}$ , and it can be naturally cast in terms of the Bifurcation Theory as follows. Treating  $p \in (0, UCP]$  as a control parameter and  $N \in (c, \infty)$  as a function of some latent continuous argument, consider an autonomous dynamical system

$$\dot{N} = t_X(N, p) - 1. \quad (3.12)$$

**Proposition 3.2.2.** *Assume (M0)–(M4).  $p_{X,c}$  is a bifurcation point of the system (3.12) and one can distinguish between three types of possible bifurcations.*

(b0)  $p_{X,c}$  is the only value of the control parameter for which (3.12) admits fixed points in  $(c, \infty)$ . In this case  $p_{X,c} < UCP$  and all  $N \in (c, \infty)$  solve  $t_X(N, p_{X,c}) = 1$ .

If there exists  $p_l \in (0, p_{X,c})$  for which (3.12) admits a fixed point  $N \in (c, \infty)$ , there are two possibilities:

(b1) (3.12) has fixed points in  $(c, \infty)$  for all  $p \in [p_l, p_{X,c})$  yet there are no fixed points corresponding to  $p_{X,c}$ ;

(b2) (3.12) has fixed points in  $(c, \infty)$  for all  $p \in [p_l, p_{X,c}]$  including  $p_{X,c}$  which then is necessary smaller than  $UCP$ .

In all cases, bifurcation curve induces a differentiable map  $(c, \infty) \ni N \mapsto p_N \in (0, p_{X,c}]$ . Therefore,  $p_{X,c}$  can be determined by finding its maximum. For bifurcations of types (b0) and (b2), this amounts to solving a system

$$\begin{cases} t_X(N, p) = 1, \\ \frac{\partial}{\partial N} t_X(N, p) = 0 \end{cases} \quad (3.13)$$

(with respect to both  $N$  and  $p$ ) and then picking up a largest  $p$  value from the set  $S = \{(N, p) \in (c, \infty) \times (0, UCP] \mid (N, p) \text{ solves (3.12)}\}$ . For the bifurcation of type (b1),  $p_{X,c} = \max(\lim_{N \rightarrow c+} p_N, \lim_{N \rightarrow \infty} p_N)$ . In particular, this holds true when (3.13) has no solution lying in  $(c, \infty) \times (0, UCP]$ .

In this subsection, we do not comment on the restrictiveness of conditions (M0–M4) and postpone this to the dedicated Section of Chapter 4. However, before proceeding to examples, we provide several clarifying remarks.

**Remark 3.2.3.** Prop. 3.2.2 establishes the procedure for finding OCP on the continuous scale (COCP). In practice, one operates on the discrete one since the number of tested items  $N$  is integer. As a rule, discrete scale OCP (DOCP), so far in the Thesis denoted as  $p_c^X$ , is lower than  $COCP = p_{X,c}$  of Prop. 3.2.1. However, the difference is usually small (see examples in Section 3.2.2) whereas the determination of DOCP is often times quite involved. Moreover, in the case of (b2),  $DOCP = p_c^X$  can often be recovered as follows:

- take  $N_c$  s.t.  $(N_c, p_{X,c})$  solves (3.13);
- set  $DOCP = \max(p_{\lfloor N_c \rfloor}, p_{\lceil N_c \rceil})$ .

In the case of (b1), DOCP is very likely to coincide with COCP.  $\square$

**Remark 3.2.4.** We have seen that, in some cases (like that of A2), cohort size  $N = N(n)$  is a function of  $n \in \mathbb{N}$  and it is more convenient to treat  $\theta_X, t_X$  and other related functions as functions of  $(n, p)$  rather than functions of  $(N, p)$ . Replacing  $N$  by  $n$  in (M0)–(M4) and then in all functions in Propositions 3.2.1–3.2.2 does not change conclusions of these Propositions provided  $(c, \infty) \ni n \mapsto N(n)$  is differentiable and strictly increasing.  $\square$

**Remark 3.2.5.** We are inclined to think that bifurcations (b1)–(b2) are the prevalent ones since we are unaware of practical examples of (b0) satisfying our assumptions. Yet Subsection 3.2.2.5 contains an example showing that a counterpart of (b0) may occur on the discrete scale. We were unable to exclude this type theoretically. Thus, appealing to the mentioned example, we are inclined to think that (b0) is not a redundant case but an exceptional one, corresponding to optimal procedures (see the discussion in Chapter 4).  $\square$

## 3.2.2 Examples

In this section, we provide several examples demonstrating applications of Prop. 3.2.2. We also provide two examples of the procedures violating our conditions. Figures appearing in this subsection were produced by making use of Desmos Graphing Calculator [17].

### 3.2.2.1 Dorfman procedure D

Recall that

$$\begin{aligned}\theta_D(N, p) &= 1 \cdot q^N + (N+1)(1-q^N) = N + 1 - Nq^N, \\ t_D(N, p) &= \frac{\theta_D(N, p)}{N} = 1 + \frac{1}{N} - q^N.\end{aligned}\quad (3.14)$$

We put  $c = 2$  in (M0) and have it (testing one item does not require D procedure). (M1) obviously holds. Since  $\frac{\partial}{\partial p} t_D(N, p) = Nq^{N-1} > 0$  for all  $N \in (2, \infty)$ , (M2) holds as well. As for (M3), note that

$$\begin{aligned}\frac{d}{dN} t_D(N, UCP) &= \frac{d}{dN} \left( 1 + \frac{1}{N} - \left( \frac{\sqrt{5}-1}{2} \right)^N \right) = \\ &\quad - \frac{1}{N^2} - \left( \frac{\sqrt{5}-1}{2} \right)^N \ln \left( \frac{\sqrt{5}-1}{2} \right).\end{aligned}$$

Equating this to zero (and solving numerically) one finds out that this function has a unique minimum at  $N_{\min} \approx 2.888$  and  $\min_{N>2} t_D(N, UCP) \approx t_D(2.888, UCP) = 1.097$ . Moreover, it has a maximum at  $N_{\max} \approx 5.75$  and then decreases to  $\lim_{N \rightarrow \infty} t_D(N, UCP) = 1$ . Since  $t_D(2, UCP) = \frac{\sqrt{5}}{2} > 1$ , (M3) holds. Finally, from (3.14) it follows that

$$\forall N \in (2, \infty) \lim_{p \rightarrow 0+} t_D(N, p) = \frac{1}{N}.$$

Therefore, (M4) holds as well.

Figure 3.4 shows a plot of the inverted bifurcation map  $N \mapsto p_N$  described in Prop. 3.2.2. In this case, it admits analytical expression:  $p_N = 1 - \left(\frac{1}{N}\right)^{\frac{1}{N}}$ . System (3.13) is given by

$$\begin{cases} \frac{1}{N} = q^N, \\ -\frac{1}{N^2} = q^N \ln q, \end{cases} \quad (3.15)$$

and can be solved analytically too. Its solution is  $(N_*, p_{D,c}) = (\text{e}, 1 - e^{-e^{-1}})$ . Recall that Samuels' [66] analysis led to  $DOCP = 1 - 3^{-3^{-1}}$ . Our COCP is quite close. Moreover, we can recover DOCP by applying the method described in the Remark 3.2.3. In this particular case the method works and affirms Samuels' result.

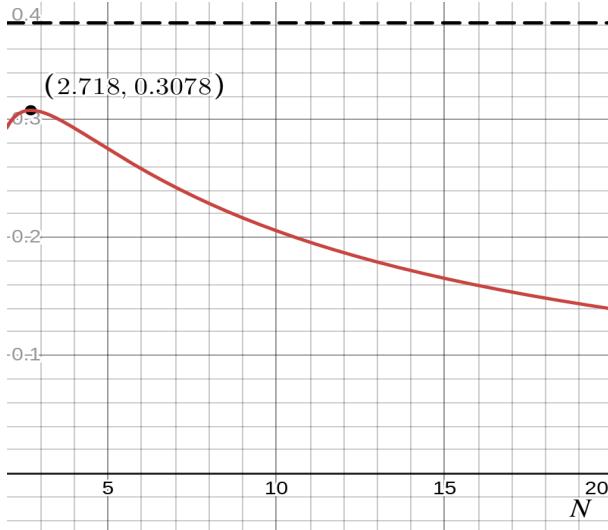


Figure 3.4: plot of  $N \mapsto p_N$  (solid line) for the procedure D; dashed line shows constant line  $p = UCP = \frac{3-\sqrt{5}}{2}$ ; maximal value yields  $COPC = 1 - e^{-e^{-1}} \approx 0.3078$ ; value  $p_{\lceil e \rceil} = 1 - 3^{-3^{-1}}$  is equal to DOCP.

### 3.2.2.2 Square Array procedure A2

As we know, in the case of A2,

$$\theta_{A2}(N, p) = 2\sqrt{N} + N \left( 1 - 2q^{\sqrt{N}} + q^{2\sqrt{N}-1} \right) \quad (3.16)$$

for  $N(n) = n^2, n \in \mathbb{N}$ . Therefore, for the sake of convenience and in view of Remark 3.2.4, we make a change of variables and treat both  $t_{A2}$  and  $\theta_{A2}$  as functions of  $(n, p)$ . We carry this convention to all related functions as well. This yields

$$t_{A2}(n, p) = \frac{2}{n} + 1 - 2q^n + q^{2n-1} = \frac{2}{n} + (1 - q^n)^2 + pq^{2n-1}. \quad (3.17)$$

From the latter expression, it follows that A2 makes sense only for  $n > 2$ . Bearing in mind the practical aspect (i.e., the fact that cohort sizes are integers) we therefore set  $c = 3$  in<sup>4</sup> (M0). Note that, due to the design of the procedure, this truncation actually restricts sizes of the tested cohorts to start from  $9 = 3 \times 3$  (not 3). It is obvious that (M1) holds.

---

<sup>4</sup>relying on results of Subsection 3.1.1, we could even set  $c = 5$

Since  $\forall p \in (0, UCP]$

$$\begin{aligned}\frac{\partial}{\partial p} \theta_{A2}(n, p) &= n^2 (2nq^{n-1} - (2n-1)q^{2n-2}) = \\ &n^2 (q^{2n-2} + 2nq^{n-1}(1-q^{n-1})) > 0,\end{aligned}$$

(M2) holds as well.

Justification of (M3) can be done by accomplishing the following steps:

- check that  $\frac{\partial}{\partial n} t_{A2}(n, p) = -\frac{2}{n^2} - 2q^n \ln q(1-q^{n-1})$  and  $\frac{\partial^2}{\partial n^2} t_{A2}(n, p) = \frac{4}{n^3} - 2q^n \ln^2 q(1-2q^{n-1})$ ;
- numerically solve  $\frac{\partial^2}{\partial n^2} t_{A2}(n, p) \Big|_{p=UCP} = 0$  and obtain two roots:  $n_1 \approx 5.278, n_2 \approx 9.448$ ;
- verify that  $n_1$  corresponds to the maximum whereas  $n_2$  corresponds to the minimum of  $n \mapsto \frac{\partial}{\partial n} t_{A2}(n, UCP)$  and that

$$\frac{\partial}{\partial n} t_{A2}(n_1, UCP) < -0.0055 < 0;$$

- conclude that  $n \mapsto t_{A2}(n, UCP)$  is decreasing and (M3) holds since

$$\lim_{n \rightarrow \infty} t_{A2}(n, UCP) = 1.$$

Finally, from the last expression given in (3.17), it follows that

$$t_{A2}(n, p) \xrightarrow[p \rightarrow 0+]{ } \frac{2}{n} \leq \frac{2}{3}.$$

Hence (M4).

Figure 3.5 shows the inverted bifurcation curve  $(3, \infty) \ni n \mapsto p_n$  of Prop. 3.2.2. System (3.13) is given by

$$\begin{cases} \frac{2}{n} - 2q^n + q^{2n-1} = 0, \\ -\frac{1}{n^2} - q^n(1-q^{n-1}) \ln q = 0. \end{cases}$$

As can be seen from the curve, it has a unique solution  $(n_*, p_{A2,c}) \approx (4.454, 0.252)$ . Recall that Kim and Hudgens [38] analysis on the discrete scale yielded  $DOCP = 0.2498$ . This point precisely coincides with  $p_{\lceil n_* \rceil}$ . Thus, the method described in 3.2.3 again led to the recovery of the DOCP.

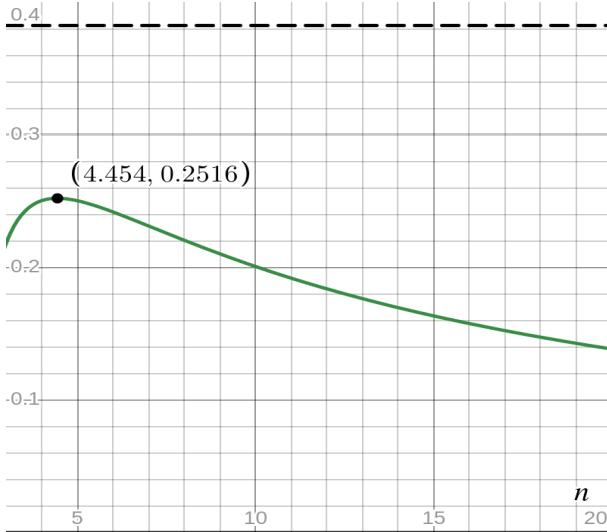


Figure 3.5: plot of  $n \mapsto p_n$  (solid line) for the procedure A2; dashed line shows constant line  $p = UCP = \frac{3-\sqrt{5}}{2}$ ; maximal value yields  $COPC \approx 0.2516$ ;  $p_5 \approx 0.2498$  is equal to DOCP.

### 3.2.2.3 Modified Dorfman procedure MD

For MD,

$$t_{MD}(N, p) = 1 - q^N + \frac{1 - pq^{N-1}}{N}. \quad (3.18)$$

We set  $c = 2$  and (M0) is satisfied (again, as in the case of procedure D, testing one item does not require procedure MD). As usually, (M1) is obvious. Since

$$\begin{aligned} \frac{\partial}{\partial p} t_{MD}(N, p) &= Nq^{N-1} - \frac{1}{N} (q^{N-1} - (N-1)pq^{N-2}) = \\ &= q^{N-1} \left( N - \frac{1}{N} \right) + \frac{N-1}{N} pq^{N-2} > 0 \end{aligned}$$

for any fixed  $N \in (2, \infty)$ , (M2) holds.  $\forall N \in (2, \infty) \lim_{p \rightarrow 0^+} t_{MD}(N, p) = \frac{1}{N} < 1$ . Hence (M4). Finally, verification of (M3) can be done in the same way as in the case of procedure A2. Since an exercise is quite lengthy and tedious, we omit the details as well as check that system (3.13) does not admit solution lying in  $(2, \infty) \times (0, UCP]$ . The latter means that we have a bifurcation of type (b1). Since  $t_{MD}(2, UCP) = 1$ , we conclude that

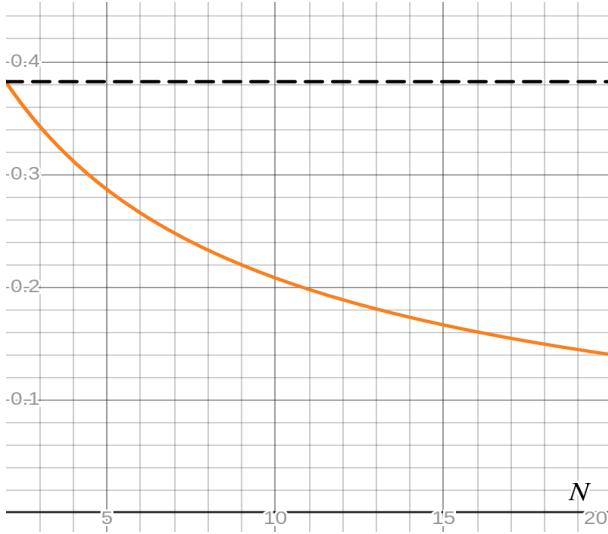


Figure 3.6: plot of  $N \mapsto p_N$  (solid line) for the procedure MD; dashed line shows constant line  $p = UCP = \frac{3-\sqrt{5}}{2}$ .

$p_{MD,c} = \lim_{N \rightarrow 2+} p_N = UCP$ . Figure 3.6 provides graphical illustration of the said.

As noted in the Remark 3.2.5, this time  $COPC = DOCP = UCP$ .

### 3.2.2.4 Sterrett procedure ST

As we know, for procedure ST,

$$t_{ST}(N, p) = 2 - q + \frac{2q - (1-q)^{-1}(1-q^{N+1})}{N}. \quad (3.19)$$

Setting  $c = 2$  in (M0), one has that (M0)–(M1) readily hold. A way to verify (M2) lies in showing that  $(1 - UCP, 1) \ni q \mapsto t_{ST}(N, p)$  is decreasing. Since

$$\frac{\partial}{\partial q} t_{ST}(N, p) = \left( \frac{2}{N} - 1 \right) - \frac{1}{N} \left( \frac{1 - q^{N+1}}{1 - q} \right)'_q,$$

one sees that the term  $2/N - 1 < 0$  for any  $N > 2$  and it suffices to note that

$$\left( \frac{1 - q^{N+1}}{1 - q} \right)'_q = (1 + q + \cdots + q^N)'_q > 0.$$

(M4) follows by noting that, for any fixed  $N > 2$ ,

$$\lim_{p \rightarrow 0+} t_{ST}(N, p) = \lim_{q \rightarrow 1-} t_{ST}(N, p) = 1 + \frac{2 - \lim_{q \rightarrow 1-} \frac{(1-q^{N+1})'_q}{(1-q)'_q}}{N} = \\ 1 + \frac{2 - \frac{N+1}{1}}{N} = \frac{1}{N} < 1.$$

As in the previous example, we omit verification of (M3). It amounts to a careful analysis of the derivative of  $(2, \infty) \ni N \mapsto t_{ST}(N, UCP)$ . One can also show that the system (3.13) does not admit solutions  $(N, p)$  lying in  $(2, \infty) \times (0, UCP]$ . Since  $UCP$  solves  $t_{ST}(2, p) = 1$  (w.r.t.  $p$ ), we again have that  $p_{ST,c} = \lim_{N \rightarrow 2+} p_N = UCP$  as in the previous example. Figure 3.7 demonstrates that the bifurcation curve qualitatively exhibits the same behavior too. Again, note that  $COPC = DOCP = UCP$ .

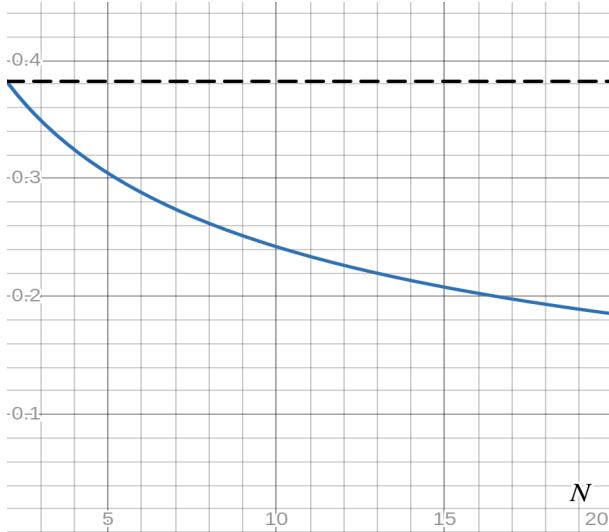


Figure 3.7: plot of  $N \mapsto p_N$  (solid line) for the procedure ST; dashed line shows constant line  $p = UCP = \frac{3-\sqrt{5}}{2}$ .

### 3.2.2.5 Examples violating our assumptions

The first example is the PT procedure. Since

$$\theta_{PT}(N, p) = N \frac{2 - q^2}{1 + q} + \frac{q^2 + q - 1}{(1 + q)^2} (1 - (-q)^N),$$

it is clear that one can not extend  $\theta_{PT}$  to the continuously differentiable function w.r.t.  $N$ . Hence, (M1) does not hold. (M3) does not hold as well. Indeed,  $UCP$  solves  $q^2 + q - 1 = 0$ ; also  $\left. \frac{2-q^2}{1+q} \right|_{p=UCP} = 1$  and

$$\frac{2-q^2}{1+q} < 1 \iff 0 < q^2 + q - 1 \iff p \in (0, UCP). \quad (3.20)$$

Another example of this kind is the procedure H. Recall that here the cohort size  $N(n) = 2^n, n \in \mathbb{N}$ , and

$$\theta_H(N, p) = 1 + 2^{n+1} \sum_{k=1}^n \frac{1 - q^{2^k}}{2^k}.$$

Thus, it again violates (M1).

### 3.3 Probabilistic analysis of the Pairwise testing procedure

In [85], it was proved that, for  $p \in \left[1 - \frac{1}{\sqrt{2}}, UCP\right]$ , the PT procedure is optimal in the class of nested procedures. That is, for any nested procedure  $X$  and for any  $N \in \mathbb{N}, \theta_{PT}(N, p) \leq \theta_X(N, p)$  uniformly over  $p \in \left[1 - \frac{1}{\sqrt{2}}, UCP\right]$ . Despite this fundamental property, the PT procedure did not receive considerable attention in the literature. After getting familiar with the PT and retrieving citing literature, we have discovered that out of fifteen citing references [1–3, 12, 22, 29, 30, 43, 48, 53, 55, 56, 76, 83, 84] retrieved by us<sup>5</sup> from Google Scholar, Malinovsky [53] was the only who investigated a problem having a direct relationship to the PT procedure. All other researchers touched the paper of Yao and Hwang [85] merely as a reference having a connection to GT with a mild relation to their own problem. These circumstances motivated us to give a broader probabilistic characterization of the PT procedure [91]. We succeeded in deriving an exact analytical expression of the moment-generating function (MGF) for the number of tests performed by the PT procedure. With the help of the MGF, it was possible to obtain common limiting theorems: strong law of large numbers (SLLN), central limit theorem (CLT), and large deviations principle (LDP). To state formal results, we need several notions.

For short, let  $\Theta_N \equiv T_{PT}$  denote the number of conducted tests required for an identification of all defectives in a given binomial set

---

<sup>5</sup>the list was generated on 28th of June, 2022; non-English references were excluded

having  $N$  items, and let  $Y_i, i = 1, \dots, N$ , be an indicator of an  $i$ th item status (1 stands for the defective one). Then  $Y_i \sim Be(p)$ . Also, let  $\bar{Y}_i := 1 - Y_i$  and

$$M_0 = \begin{pmatrix} 1 & 0 \\ 1 & 0 \end{pmatrix}, \quad M_1 = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}. \quad (3.21)$$

Our first result gives an explicit expression for  $\Theta_N$  in terms of the above quantities.

**Proposition 3.3.1.** *Let  $A = \left\{ \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}, \begin{pmatrix} 0 & 1 \\ 0 & 1 \end{pmatrix} \right\}$  and  $B_k = \begin{pmatrix} Y_k & \bar{Y}_k \\ 1 & 0 \end{pmatrix}$  for  $k = 1, \dots, N$ . Then  $\Theta_2 = 3Y_2 + \bar{Y}_2(1 + Y_1)$ ,  $\Theta_3 = 2 + \bar{Y}_3Y_2 + Y_3\Theta_2$ , and*

$$\begin{aligned} \Theta_N = 1 + Y_N(\bar{Y}_{N-1}Y_{N-2} + 2) + Y_{N-1} + \\ \sum_{j=3}^{N-1} (\bar{Y}_{j-1}Y_{j-2} + Y_{j-1} + 1) (Y_j + \bar{Y}_{j-1}\mathbb{1}\{B_N B_{N-1} \cdots B_{j+1} \in A\}) + \\ Y_2 + \bar{Y}_2\mathbb{1}\{B_N B_{N-1} \cdots B_3 \in A\} \text{ for } N \geq 4. \end{aligned} \quad (3.22)$$

The expression above provides insight into the structure of  $\Theta_N$ . Our next result provides the announced explicit formula of the MGF.

**Theorem 3.3.2.** *Let  $M_{\Theta_N}(\lambda)$  denote the moment generating function of  $\Theta_N$  at  $\lambda \in \mathbb{R}$ . Set*

$$\alpha_i = \alpha_i(\lambda) = \frac{1}{2} \left( p e^{2\lambda} + (-1)^i \sqrt{p^2 e^{4\lambda} + 4qe^\lambda(q + pe^\lambda)} \right), \quad i = 0, 1; \quad (3.23)$$

$$\kappa_N = \kappa_N(\lambda) = \frac{\alpha_0^N - \alpha_1^N}{\alpha_0 - \alpha_1} \text{ for } N \geq 0. \quad (3.24)$$

Then  $M_{\Theta_N}(\lambda) =$

$$\begin{aligned} e^{2\lambda} \left[ \left( (1-q)^2 e^{3\lambda} + q(1-q)^2 e^{2\lambda} + q(1-q^2)e^\lambda + q^2 \right) \kappa_{N-2} + \right. \\ \left. q \left( (1-q)^2 e^{3\lambda} + q(1-q)(2-q)e^{2\lambda} + 2q^2(1-q)e^\lambda + q^3 \right) \kappa_{N-3} \right] \end{aligned} \quad (3.25)$$

for  $N \geq 3$ .

The remaining results are the consequences of the previous one.

**Corollary 3.3.3.**  $E \Theta_N = N \frac{2-q^2}{1+q} + \frac{q^2+q-1}{(1+q)^2} (1 - (-q)^N)$ ,

$$\begin{aligned} \text{Var } \Theta_N &= N \frac{(1-q)}{(q+1)^3} \left( q (q^3 + 3q^2 + 5q + 4) + \right. \\ &\quad \left. (-q)^N (2q+4) (q^2 + q - 1) \right) + \\ &\quad \frac{\left(1 - (-q)^N\right)}{(q+1)^4} \left( q (5q^2 + 3q - 7) + (-q)^N (q^2 + q - 1)^2 \right), N \geq 3. \end{aligned} \tag{3.26}$$

**Corollary 3.3.4.** *The following asymptotic results apply to  $\Theta_N$  as  $N \rightarrow \infty$ .*

LLN:  $\frac{\Theta_N}{N} \xrightarrow{L_2} \frac{2-q^2}{1+q}$  and  $\frac{\Theta_N}{N} \xrightarrow{a.s.} \frac{2-q^2}{1+q}$ .

CLT:  $\sqrt{N} \left( \frac{\Theta_N}{N} - \frac{2-q^2}{1+q} \right) \xrightarrow{d} N(0, \sigma^2)$ ,  $\sigma^2 = \frac{q(1-q)(q^3+3q^2+5q+4)}{(q+1)^3}$ .

LDP:  $\frac{\Theta_N}{N}$  satisfies the Large Deviation Principle with a good rate function  $I$  equal to the Legendre transform of  $\mathbb{R} \ni \lambda \mapsto \ln \alpha_0(\lambda)$  with  $\alpha_0(\lambda)$  given by (3.23). That is, for any closed  $C \subset \mathbb{R}$  and any open  $O \subset \mathbb{R}$ ,

$$\limsup_{N \rightarrow \infty} \frac{1}{N} \ln P \left( \frac{\Theta_N}{N} \in C \right) \leq - \inf_{x \in C} I(x)$$

and

$$- \inf_{x \in O} I(x) \leq \liminf_{N \rightarrow \infty} \frac{1}{N} \ln P \left( \frac{\Theta_N}{N} \in O \right),$$

where  $I(x) = \sup_{\lambda \in \mathbb{R}} (x\lambda - \ln \alpha_0(\lambda))$ .

# Chapter 4

## Discussion and conclusions

In this chapter, we have collected several remarks about the results stated in Chapter 3. In our opinion, these remarks fall out of the scope of the related literature review yet are worth mentioning. For each section of Chapter 3, we devote a separate section. The last section is devoted to several concluding observations of a general nature.

### 4.1 About results of Section 3.1

It is always possible to obtain numerical solutions of optimal configurations. However, examples given in Subsection 3.1.2 clearly demonstrate that analytic expressions are very useful. Moreover, in an asymptotic regime, numerical solutions are of minimal applicability because of computational errors, and the absence of analytic expressions prevents the asymptotic analysis.

To get a quick example, consider the asymptotic regime when  $p \rightarrow 0+$ . From results of Section 3.1

$$t_{MD}(N_{opt}^{MD}(p), p) = \frac{2 + o(1)}{N_{opt}^{MD}(p)} \text{ and } t_{ST}(N_{opt}^{ST}(p), p) = \frac{2 + o(1)}{N_{opt}^{ST}(p)}.$$

Therefore, taking into account Theorems 3.1.1–3.1.2,

$$\lim_{p \rightarrow 0+} \frac{t_{MD}(N_{opt}^{MD}(p), p)}{t_{ST}(N_{opt}^{ST}(p), p)} = \lim_{p \rightarrow 0+} \frac{N_{opt}^{ST}(p)}{N_{opt}^{MD}(p)} = \sqrt{2}. \quad (4.1)$$

Hence, for small  $p$ 's, testing optimally by using the procedure MD results in an average number of tests per object, which is approximately 1.4

times larger than the average number of tests per object obtained by applying procedure ST at its optimal configuration. Analyses of a similar kind appear in the literature. E.g., in [57], the limit (4.1) appears with MD replaced by D and is the same. This means that asymptotically the original and the modified Dorfman procedures are identical. However, their behavior differs for  $p$ 's far from the origin (see [55]).

## 4.2 About results of Section 3.2

In this Section, we discuss conditions (M0)–(M4) of Section 3.2 with a focus on their meaning and restrictiveness.

(M0) may be viewed as a condition required to confine the range of the dynamical system we make use of when looking for the  $COPC = p_{X,c}$ . Together with (M3), it also defines the boundary value for the determination of  $p_{X,c}$  for the case  $p_{X,c} = UCP$ . In a usual case, one can set  $c = 2$  and have it since, in this context,  $c = 2$  means that taking one item, we do not need any GT procedure: to identify the defectiveness of a single item, one always needs one test<sup>1</sup>.

Constraint (M1) is the most restrictive since not all binomial GT procedures can be naturally extended to have a differentiable mean with respect to both arguments (though in the case of  $p$  this always holds true). Nonetheless, this particular assumption is the one we heavily rely on. It also enables us to draw the connection with the bifurcation theory.

Other constraints can be justified naturally and attributed to many binomial GT procedures in general.

(M2) states that an average number of tests per batch spanning  $N$  items should increase together with the rate of defectiveness. For justification, we mention another fundamental result due to Yao and Hwang [86] who have demonstrated that  $\forall N \in \mathbb{N}$ , function  $(0, UCP] \ni p \mapsto \inf_X \theta_X(N, p)$ , with an infimum being taken over all possible BGT procedures, is strictly increasing.

At first glance, it may seem that (M3) rules out procedures having  $p_{X,c} = UCP$ . As demonstrated by example, this is not the case. In fact, it restricts the subset of GT procedures to those having  $p_{X,c}$  on the boundary of the domain of the bifurcation curve. We are inclined to

---

<sup>1</sup>see examples in Subsection 3.2.2

think that this way, we rule out optimal procedures, i.e., those which are best performing in certain classes. However, we do not treat this as a drawback since, for optimal procedures, one generally expects  $p_{X,c} = UCP$ .

Finally, (M4) in technical terms states that we focus on the procedures applicable to any number of tested items, at least for some  $p$ 's in the range of their sensibility. This is very often the case since many procedures are suitable for large-scale testing when the rate of defectiveness is small.

In case (M0)–(M4) hold, our algorithm appears to be efficient. There is a word of caution: one has to choose  $c$  in (M0) carefully. The point is that the dynamical system defined by (3.12), when viewed on a wider domain, may exhibit more complicated bifurcations. Figure 4.1 provides a convincing graphical illustration.

Turning to the types of bifurcations, one sees that, in terms suggested by Strogatz [74], (b1) is usually the saddle point bifurcation, whereas for (b2) the system admits fixed points for all but boundary value of the control parameter  $p \in (0, UCP]$ . We are inclined to think that the dynamical system approach we took could be extended to GT procedures violating (M1) and successfully used to investigate other general properties of GT procedures, yet one needs to work on the discrete scale.

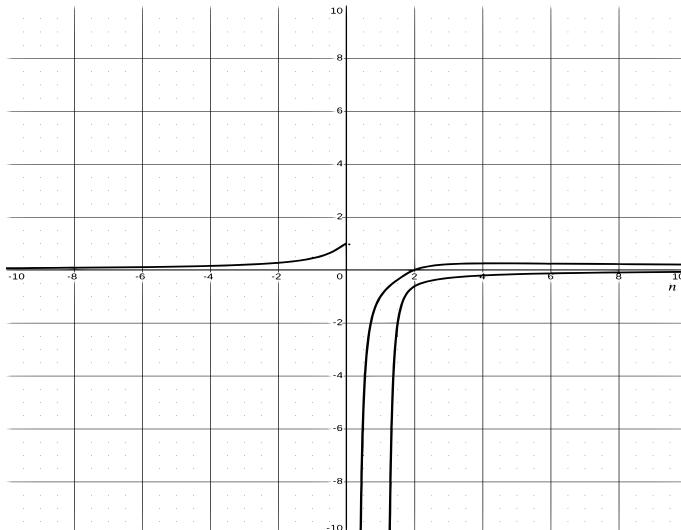


Figure 4.1: plot of bifurcation curve for procedure A2 when the domain of the dynamical system (3.12) is extended to the whole real line.

We end by noting that relationships (3.20) uncover an interesting fact:

$$\forall N \in (2, \infty) t_{PT}(N, UCP) = 1,$$

which means that the PT procedure is "almost" of type (b0). Yao and Hwang [85] proved that the PT is an optimal nested testing procedure if and only if  $p \in \left[1 - \frac{\sqrt{2}}{2}, \frac{3-\sqrt{5}}{2}\right]$ . This suggests that procedures of type (b0) are likely to be those which are optimal in some region.

### 4.3 About results of Section 3.3

We did not provide any specific examples of applications of results stated in Section 3.3 beyond moment calculations presented in Corollary 3.3.3. Despite that, one can list several reasons supporting the relevance of our analysis.

- Though the definition of an optimal procedure is usually tied to an average number of tests, when choosing between several procedures, it is desirable to evaluate their performance by taking into account multiple aspects. For example, procedure  $Pr1$  may perform slightly better than  $Pr2$  regarding an average number of tests. However,  $Pr1$  may have a considerably larger variance than  $Pr2$  and, therefore, the previously mentioned slight gain of  $Pr1$  could be gladly traded by the practitioner in favor of  $Pr2$ .
- We have already mentioned that the importance of the PT remained unrecognized in the literature, and there is more to say on that.
  - Many GT procedures described in the literature have limited applicability in certain areas due to the *dilution effect* described in 1.3. Recall that basically this means the following. Given procedure  $X$  and fixed  $p$ ,  $N_{opt}^X(p)$  may be too large and inflate the operating characteristics (sensitivity and/or specificity) of the test kit at hand making them unacceptably low (aka diluting) for that particular application<sup>2</sup>. With respect to this property, the PT procedure is a very favorable option: it requires only pools of size  $N = 2$ , and this holds true for all  $p$ 's in the region of its optimality  $\left[\frac{2-\sqrt{2}}{2}, \frac{3-\sqrt{5}}{2}\right]$ .

---

<sup>2</sup>in theory, (BTA3) stated in Section 2.1 prevents from this; however, in practice, it may be a serious obstacle

- The region  $\left[\frac{2-\sqrt{2}}{2}, \frac{3-\sqrt{5}}{2}\right]$  where the PT procedure performs optimally is bounded away from zero in contrast to many other GT procedures which do better for  $p$ 's close to zero. In certain applications, this property may be of significant importance. For example, when screening for a quite widespread infectious disease.
- In [85] it was conjectured that there exists such  $p_0 \in \left[\frac{2-\sqrt{2}}{2}, \frac{3-\sqrt{5}}{2}\right]$  that for  $p \in \left[p_0, \frac{3-\sqrt{5}}{2}\right]$  the PT procedure is optimal over all (not necessarily nested) procedures satisfying the BTA.
- Our Proposition 3.3.1 demonstrates that, despite apparently simple recurrence governing the evolution of  $\Theta_N$  (see Eq. (5.22)), the resulting dependence structure is not so simple. At least we were not able to analyze its behavior neither by making use of Markov chains theory, nor by making use of martingale theory. A well-developed apparatus of weakly dependent sequences also did not promise easy deduction of Corollary 3.3.4. More than that, even direct moment calculation exercise, though accomplishable for  $E\Theta_N$  at a reasonable price (see Lemma in Section 4 of [85]), becomes much more involved when it comes to  $\text{Var } \Theta_N$  and higher order moments. This way,  $(\Theta_N)_{N \geq 2}$  yields an example of a sequence of positive integer-valued random variables having an interesting probabilistic structure encountered in practical application and not designed artificially for learning or other purposes.

Finally, we believe that our results may be useful for the solution of several unresolved conjectures. Namely, the one stated in [85] and mentioned above, and the generalized PT optimality conjecture stated in [53].

## 4.4 Several concluding remarks

In the Thesis, we have focused on the PGT procedures satisfying the BTA. Such a framework may seem too simplified for the real setting. However, this is not generally true. For some applications (say, in quality control or computer science), these assumptions are justifiable. Moreover, though developed quite long ago, simple procedures like  $D$ ,  $MD$ , and  $ST$  are not outdated and still in use even in large scale projects. E.g., the American Red Cross makes use of Dorfman procedure for the screening of blood

donations for HIV and hepatitis [18] whereas, in Lithuania, Dorfman procedure was not long ago applied to test for COVID-19 employees of larger firms and pupils attending public schools [7]. There are several reasons explaining why the BTA-based procedures are still in use.

*Convenience.* Even though since the original work of Dorfman many procedures tied to particular needs of applications considered were developed, the organizational flow of the whole project may not afford to apply more elaborated procedures. In such cases, simple procedures appear to be a good alternative, resulting in cost savings.

*Specifics of application.* It appears that one of the reasons for the emergence of A2 in genetic applications was testing equipment: it was such that A2 was very handy option.

*Tolerable errors.* Though there exist a lot of generalizations allowing imperfection of the test (e.g., [46], [45], [33], [8], [10]), as noted by several authors [9], [38], [55], procedures assuming perfect tests can be quite accurate since modern tests exhibit very small errors.

Investigations of the BTA-based PGT procedures remain essential due to other reasons as well.

For example, binomial testing procedures may serve as a basement for more elaborated ones: procedures MD and ST were built on the top of D; in [46] and [33], A2 serves as a basis for extensions incorporating imperfection of the test and dilution effect; in [8] extensions addressing subject-specific risk characteristics and imperfect tests are proposed for the D procedure whereas in [10] the authors do the same focusing only on the heterogeneity of the population.

For another example of the usefulness of the BTA-based procedures, consider benchmarking. The BTA-based procedures, being more simple to treat analytically, provide theoretically justified benchmark thresholds for more elaborated procedures that assume the imperfection of the test and/or other specific conditions.

In view of the said, our findings seem to be useful input to the existing GT knowledge base.

# Chapter 5

## Proofs

### 5.1 Proofs of results stated in Section 3.1

Before proceeding to the proofs, we give several remarks. In [55] it was demonstrated that, for a fixed  $p \in (0, \frac{3-\sqrt{5}}{2})$ , function  $[1, \infty) \ni N \mapsto t_{ST}(N, p)$  admits a unique absolute minimum which is attained at the unique zero of  $[1, \infty) \ni N \mapsto \frac{\partial}{\partial N}t_{ST}(N, p)$ , say  $N_*^{ST}(p)$ . Therefore,  $N_{opt}^{ST}(p) \in \{[N_*^{ST}(p)], [N_*^{ST}(p)] + 1\}$ , and the choice between two possible values is made by evaluating whether  $t_{ST}([N_*^{ST}(p)], p) > t_{ST}([N_*^{ST}(p)] + 1, p)$  or  $t_{ST}([N_*^{ST}(p)], p) \leq t_{ST}([N_*^{ST}(p)] + 1, p)$  holds true. For the case of the procedure MD, Pfeifer and Enis [62] have obtained a quite similar result stated below.

**Theorem 5.1.1.** [Pfeifer and Enis [62], Lemma 2] For a fixed  $p \in \left(0, \frac{3-\sqrt{5}}{2}\right)$ , function  $[1, \infty) \ni N \mapsto t_{MD}(N, p)$  admits a unique absolute minimum attained at the smallest zero of  $[1, \infty) \ni N \mapsto \frac{\partial}{\partial N}t_{MD}(N, p)$ , say  $N_*^{MD}(p)$ . In the set

$$A^{MD} = \left\{ (p, N) \in \left(0, \frac{3-\sqrt{5}}{2}\right) \times [1, \infty) : t_{MD}(N, p) < 1 \right\},$$

$N_*^{MD}(p)$  is the only zero of  $N \mapsto \frac{\partial}{\partial N}t_{MD}(N, p)$ .

From (2.1) and (2.2) it follows that

$$t_D(N, p) = 1 - q^N + \frac{1}{N} \text{ and } t_{MD}(N, p) = t_D(N, p) - \frac{pq^{N-1}}{N}.$$

Therefore,  $t_D(N, p) - t_{MD}(N, p) = \frac{pq^{N-1}}{N} > 0$ . Hence, appealing to the result of Samuels [66] mentioned in Section 3.1, we thus conclude that looking for  $N_{opt}^{MD}(p)$  corresponding to  $p \in (0, 1 - (1/3)^{1/3})$ , one can apply the same algorithm as for  $N_*^{ST}(p)$  above: it is enough to find the unique  $N_*^{MD}(p)$  and then select  $N_{opt}^{MD}(p) \in \{\lfloor N_*^{MD}(p) \rfloor, \lfloor N_*^{MD}(p) \rfloor + 1\}$ . In the region  $\left[1 - (1/3)^{1/3}, \frac{3-\sqrt{5}}{2}\right)$  additional care is needed.

*Proof of Theorem 3.1.1.* Consider equation  $\frac{\partial}{\partial N} t_{ST}(N, p) = 0$ . Simple rearrangement shows that it is equivalent to equality

$$\frac{1}{\ln q} - \frac{1 - 2pq}{\ln q} \left(\frac{1}{q}\right)^{N+1} = N. \quad (5.1)$$

Denote the lhs by  $h(N)$ . Then (5.1) means that  $N_*^{ST}(p)$  is a fixed point of  $h : [0, \infty) \rightarrow \mathbb{R}$ . Since  $h$  is translated and scaled increasing exponential function with  $h(0) < 0$ , that fixed point is unique in agreement with the results discussed above. Moreover, in view of these results, it suffices to demonstrate that  $N_*^{ST}(p) \in \left[\sqrt{2p^{-1}} - 1, \sqrt{2p^{-1}} + 1\right]$  in order to deduce that  $N_{opt}^{ST}(p) \in \left\{\lfloor \sqrt{2p^{-1}} \rfloor + i : i \in \{-1, 0, 1, 2\}\right\}$ . Taking into account the exponential form of  $h$ , the latter will follow provided we show that

$$h\left(\sqrt{\frac{2}{p}} + 1\right) > \sqrt{\frac{2}{p}} + 1 \quad \text{and} \quad h\left(\sqrt{\frac{2}{p}} - 1\right) < \sqrt{\frac{2}{p}} - 1. \quad (5.2)$$

For each  $m \in \{-1, 0, 1\}$ , define a function  $\left(0, \frac{3-\sqrt{5}}{2}\right) \ni p \mapsto g_m(p)$  by

$$g_m(p) = \frac{1}{q} \left( \frac{1 - 2pq}{q^{1+m} \left(1 - \ln q \sqrt{\frac{2}{p}} \left(1 + m \sqrt{\frac{p}{2}}\right)\right)}\right)^{\sqrt{\frac{p}{2}}}. \quad (5.3)$$

By simple rearrangement, it follows that (5.2) is equivalent to

$$g_1(p) > 1 \quad \text{and} \quad g_{-1}(p) < 1. \quad (5.4)$$

Figure 5.1 shows the graphs of  $g_1, g_0$ , and  $g_{-1}$ . These suggest that relationships (5.4) do hold outside the zero neighborhood though, due to resolution issues, the true behavior close to the origin may be masked.

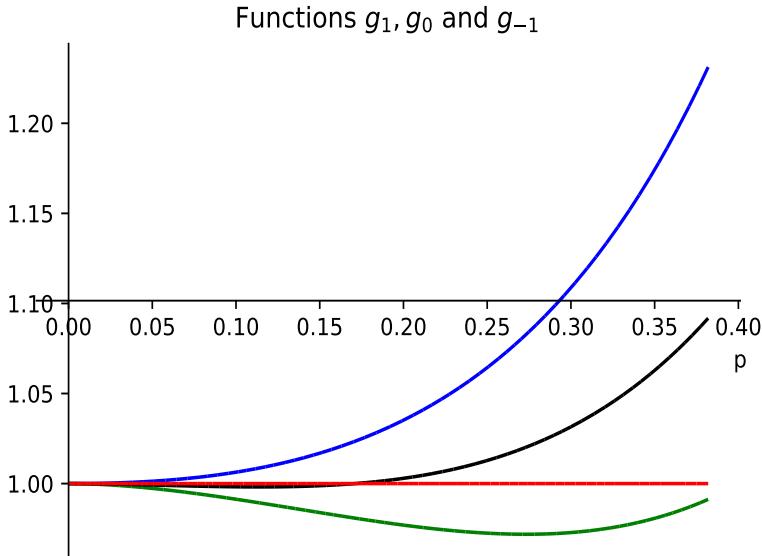


Figure 5.1: Functions  $g_m$  on  $\left(0, \frac{3-\sqrt{5}}{2}\right)$ :  $g_1$  plotted in blue,  $g_0$  plotted in black, and  $g_{-1}$  plotted in green.

To see that here  $g_1(p) > 1$  and  $g_{-1}(p) < 1$ , consider expansion

$$\begin{aligned}
 \ln g_m(p) = & \\
 & \sqrt{\frac{p}{2}} \left( \ln(1 - 2pq) - (1 + m) \ln q - \ln \left( 1 - \ln q \left( m + \sqrt{\frac{2}{p}} \right) \right) \right) \\
 & - \ln q = \sqrt{\frac{p}{2}} \left( p(1 + m - 2q) + p^2 \left( \frac{1+m}{2} - 2q^2 \right) + O(p^3) \right) + \\
 & m\sqrt{\frac{p}{2}} \ln q + \sqrt{\frac{1}{2p}} (\ln q)^2 \left( 1 + m\sqrt{\frac{p}{2}} \right)^2 + \frac{2}{3} \frac{(\ln q)^3}{p} \left( 1 + m\sqrt{\frac{p}{2}} \right)^3 + \\
 & \frac{1}{\sqrt{2}} \frac{(\ln q)^4}{p^{3/2}} \left( 1 + m\sqrt{\frac{p}{2}} \right)^4 + O(p^3). \quad (5.5)
 \end{aligned}$$

Plugging in  $m = \pm 1$  yields

$$\begin{aligned}\ln g_{-1}(p) &= \sqrt{\frac{p}{2}} (-2pq(1 + pq + O(p^2)) - \sqrt{\frac{p}{2}} \ln q \left(1 - \frac{\ln q}{p}\right) - \\ &\quad (\ln q)^2 + \frac{\sqrt{p}(\ln q)^2}{2\sqrt{2}} + \frac{2}{3} \frac{(\ln q)^3}{p} \left(1 - \sqrt{\frac{p}{2}}\right)^3 + \frac{(\ln q)^4}{\sqrt{2}p^{\frac{3}{2}}} \left(1 - \sqrt{\frac{p}{2}}\right)^4 + \\ &\quad O(p^3) = -\sqrt{\frac{p}{2}} (2p + \ln q (2 + O(p))) - \\ &\quad (\ln q)^2 \left(1 - \frac{2}{3} \frac{\ln q}{p} \left(1 - \sqrt{\frac{p}{2}}\right)^3\right) + O\left(p^{\frac{5}{2}}\right) = -\frac{5}{3}p^2 + O\left(p^{\frac{5}{2}}\right).\end{aligned}$$

and

$$\begin{aligned}\ln g_1(p) &= \sqrt{\frac{p}{2}} (3p^2 - 2(pq)^2 + O(p^3)) + \sqrt{\frac{p}{2}} \ln q \left(1 + \frac{\ln q}{p}\right) + \\ &\quad (\ln q)^2 + \frac{\sqrt{p}(\ln q)^2}{2\sqrt{2}} + \frac{2}{3} \frac{(\ln q)^3}{p} \left(1 + \sqrt{\frac{p}{2}}\right)^3 + \frac{(\ln q)^4}{\sqrt{2}p^{\frac{3}{2}}} \left(1 + \sqrt{\frac{p}{2}}\right)^4 + \\ &\quad O(p^3) = (\ln q)^2 \left(1 + \frac{2}{3} \frac{\ln q}{p} \left(1 + \sqrt{\frac{p}{2}}\right)^3\right) = \frac{1}{3}p^2 + O\left(p^{\frac{5}{2}}\right).\end{aligned}$$

Hence,  $g_1(0 + 0) = g_{-1}(0 + 0) = 1$  and  $g_1(p) > 1, g_{-1}(p) < 1$  for all  $p \in (0, \delta)$  provided  $\delta > 0$  is small enough. One can further show that  $g'_1$  is positive on  $(0, \frac{3-\sqrt{5}}{2})$ , whereas, in case of  $g'_{-1}$ , the following hold true: it has a unique zero  $x_0 \in (0, \frac{3-\sqrt{5}}{2})$ ; it is negative on  $(0, x_0)$  and positive on  $(x_0, \frac{3-\sqrt{5}}{2})$ ; finally,

$$\lim_{x \rightarrow \frac{3-\sqrt{5}}{2}^-} g_{-1}(x) \approx 0.9912.$$

Calculations being lengthy and tedious nonetheless require only standard calculus and we, therefore, omit the details. Putting all together, relationships (5.4) do hold. Taking into account all the said and then by the similar argument as above, it follows that

$$N_*^{ST}(p) \in \left[\sqrt{2p^{-1}} - 1, \sqrt{2p^{-1}}\right] \iff g_0(p) \geq 1$$

and

$$N_*^{ST}(p) \in \left[\sqrt{2p^{-1}}, \sqrt{2p^{-1}} + 1\right] \iff g_0(p) \leq 1.$$

Also, since  $\left(0, \frac{3-\sqrt{5}}{2}\right) \ni p \mapsto N_*^{ST}(p)$  is continuous and strictly decreasing<sup>1</sup>,  $g_0 - 1$  has a unique zero  $p_* \in (0, \frac{3-\sqrt{5}}{2})$  (see figure 5.2) and sets  $g_0^{-1}((-\infty, 1]), g_0^{-1}([1, \infty))$  are connected, i.e., intervals  $(0, g_0^{-1}(\{1\}))$ ,

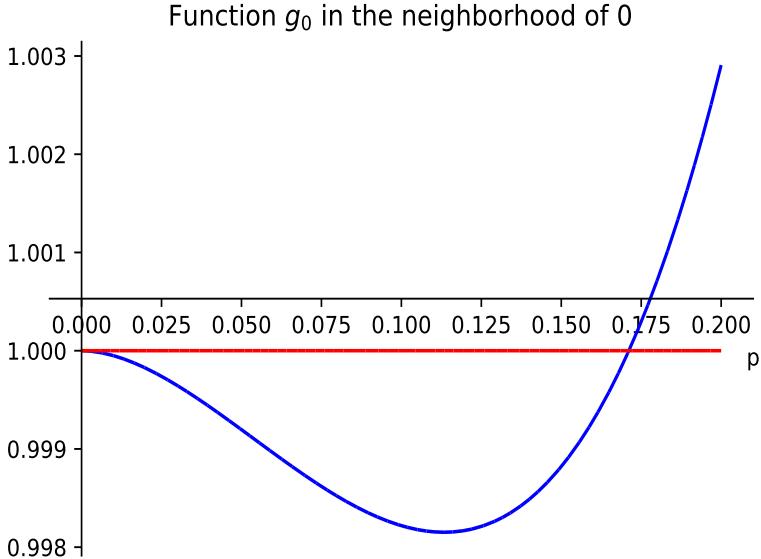


Figure 5.2: The behavior of  $g_0$  near the origin: the point of intersection with 1 is the cut-off point  $p_*$ . For  $p$ 's on the left from it,  $N_{opt}^{ST}(p) \in \{\lfloor \sqrt{2p^{-1}} \rfloor, \lfloor \sqrt{2p^{-1}} \rfloor + 1, \lfloor \sqrt{2p^{-1}} \rfloor + 2\}$ , whereas for  $p$ 's on the right,  $N_{opt}^{ST}(p) \in \{\lfloor \sqrt{2p^{-1}} \rfloor, \lfloor \sqrt{2p^{-1}} \rfloor + 1\}$

$\left[ g_0^{-1}(\{1\}), \frac{3-\sqrt{5}}{2} \right)$  respectively. To finish the proof, it remains to demonstrate that  $\left\lfloor \sqrt{\frac{2}{p}} \right\rfloor - 1$  is never optimal on  $\left( p_*, \frac{3-\sqrt{5}}{2} \right)$ . For this, consider a function

$$\begin{aligned} f(p) &= t_{ST} \left( \left\lfloor \sqrt{\frac{2}{p}} \right\rfloor - 1, p \right) - t_{ST} \left( \left\lfloor \sqrt{\frac{2}{p}} \right\rfloor, p \right) = \\ &\quad \mathbb{1} \left\{ p \in \left( p_*, \frac{2}{9} \right] \right\} (t_{ST}(1, p) - t_{ST}(2, p)) + \\ &\quad \mathbb{1} \left\{ p \in \left( \frac{2}{9}, \frac{3-\sqrt{5}}{2} \right) \right\} (t_{ST}(2, p) - t_{ST}(3, p)) \end{aligned}$$

on  $\left( p_*, \frac{3-\sqrt{5}}{2} \right)$ . It is left continuous and has a single point of discontinuity equal to  $\frac{2}{9}$ . Since  $f'$  is negative on  $(p_*, \frac{2}{9}) \cup \left( \frac{2}{9}, \frac{3-\sqrt{5}}{2} \right)$ , invoking left

---

<sup>1</sup>this needs some reasoning yet we omit the details

continuity, we infer that its minimal values are  $f(2/9)$  on  $(p_*, 2/9]$  and  $f\left(\frac{3-\sqrt{5}}{2}\right)$  on  $\left(\frac{2}{9}, \frac{3-\sqrt{5}}{2}\right)$  respectively. By direct substitution and because of left continuity,  $f\left(\frac{3-\sqrt{5}}{2}\right) = 0$ , whereas numerical estimation yields  $f(2/9) \approx 0.018976$ . The verification of the enumerated properties of  $f$  requires only lengthy standard calculus and we omit it. The graph of  $f$  illustrating its behavior is given in figure 5.3.  $\square$

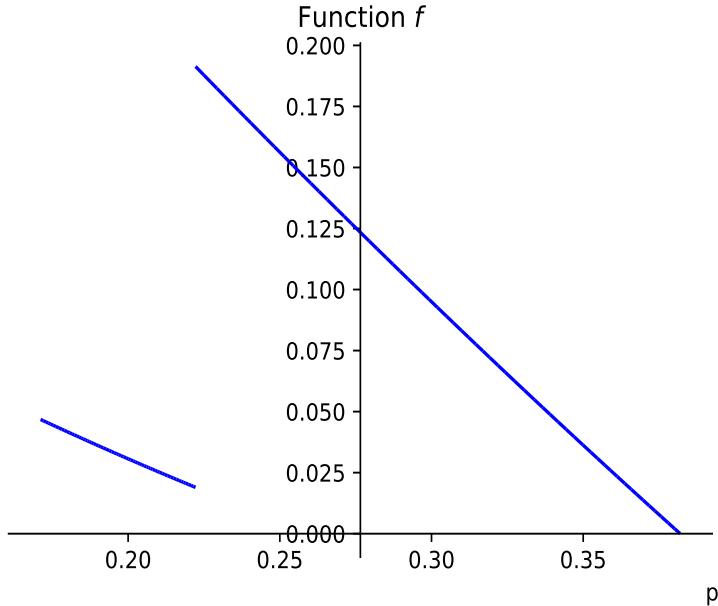


Figure 5.3: The behavior of function  $f(p) = t_{ST}\left(\left\lfloor \sqrt{\frac{2}{p}} \right\rfloor - 1, p\right) - t_{ST}\left(\left\lfloor \sqrt{\frac{2}{p}} \right\rfloor, p\right)$  on  $\left(p_*, \frac{3-\sqrt{5}}{2}\right)$ .

*Proof of Theorem 3.1.2.* For the sake of clarity, we split the proof into three steps. As in the proof of Theorem 3.1.1, some tedious details are omitted and only the sketch is given.

*Step 1: slicing  $A^{MD}$ .* In this step, we show that

$$\left\{ \left( p, \frac{1}{\sqrt{p}} + 1 - \frac{5}{2}p \right) : p \in \left( 0, \frac{3 - \sqrt{5}}{2} \right) \right\} \subseteq A^{MD}, \quad (5.6)$$

where  $A^{MD}$  is the same as in the statement of Theorem 5.1.1. To this end, note that, by elementary rearrangement,

$$t_{MD}(N, p) < 1 \iff 1 < q^{N-1}(Nq + p) = q^{N-1}((N-1)q + 1) = \\ [\text{denoting } N-1 = x] = q^x(1 + qx) \iff 0 < x \ln q + \ln(1 + qx).$$

Plugging in  $N = \sqrt{p^{-1}} + 1 - (5/2)p$ , we then obtain the condition we need to check:

$$\left(\frac{1}{\sqrt{p}} - \frac{5p}{2}\right) \ln q + \ln \left(1 + q \left(\frac{1}{\sqrt{p}} - \frac{5p}{2}\right)\right) > 0.$$

Putting  $y = \sqrt{p}$  this translates to checking that

$$f(y) = \left(\frac{1}{y} - \frac{5y^2}{2}\right) \ln(1 - y^2) + \ln \left(1 + (1 - y^2) \left(\frac{1}{y} - \frac{5y^2}{2}\right)\right)$$

is positive on its domain  $\left(0, \sqrt{\frac{3-\sqrt{5}}{2}}\right)$ . One can show that  $f$  is strictly convex on  $\left(0, \sqrt{\frac{3-\sqrt{5}}{2}}\right)$ . Consequently,  $f'$  is non-decreasing and upper bounded by  $\lim_{x \rightarrow \sqrt{\frac{3-\sqrt{5}}{2}}^-} f'(x) \approx -1.66$ . This, in turn, yields that  $f$  is decreasing and lower bounded by  $\lim_{x \rightarrow \sqrt{\frac{3-\sqrt{5}}{2}}^-} f(x) \approx 0.024$ .

*Step 2: bracing the optimal points.* Rewrite equation  $\frac{\partial}{\partial N} t_{MD}(N, p) = 0$  as follows:

$$-\frac{q}{p} \left(N^2 + \frac{1}{\ln q} \left(\frac{1}{q}\right)^N\right) + \frac{1}{\ln q} = N,$$

and denote the function on the rhs by  $f(N, p)$ . Then each fixed point of  $N \mapsto f(N, p)$  is the zero of  $\frac{\partial}{\partial N} t_{MD}(N, p)$ . In particular, the statement applies to  $N_*^{MD}(p)$ . In this step, by making use of this observation, we show that  $N_*^{MD}(p)$  (the detailed explanation of the latter fact is given in *Step 3*) does not deviate a lot from  $\sqrt{p^{-1}}$ . To achieve the goal, we consider points

$$N = N(\theta) = \left(\frac{1}{\sqrt{p}} + 1 - \frac{5}{2}p\right)\theta + \left(\frac{1}{\sqrt{p}} - p\right)(1 - \theta) = \\ \frac{1}{\sqrt{p}} - p + \theta \left(1 - \frac{3}{2}p\right), \quad \theta \in [0, 1], \quad (5.7)$$

and demonstrate that

$$\begin{aligned} \forall p \in \left(0, \frac{3-\sqrt{5}}{2}\right) \exists! \theta \in [0, 1] : f(N(\theta), p) = N(\theta) \iff \\ \forall p \in \left(0, \frac{3-\sqrt{5}}{2}\right) \exists! \theta \in [0, 1] : p^{\frac{3}{2}}(f(N(\theta), p) - N(\theta)) = 0. \quad (5.8) \end{aligned}$$

For convenience, put  $h(\theta, p) = p^{\frac{3}{2}}(f(N(\theta), p) - N(\theta))$ . Calculating derivative yields

$$\begin{aligned} \frac{\partial}{\partial \theta} h(\theta, p) &= p^{\frac{3}{2}} \left( \frac{\partial}{\partial N} f(N, p) - 1 \right) \frac{\partial}{\partial \theta} N(\theta) = \\ &= p^{\frac{3}{2}} \left( -\frac{q}{p} \left( 2N(\theta) - \left(\frac{1}{q}\right)^{N(\theta)} \right) - 1 \right) \left( 1 - \frac{3}{2}p \right) \quad (5.9) \end{aligned}$$

and then, since  $\frac{\partial}{\partial \theta} N(\theta)$  does not depend on  $\theta$ ,

$$\begin{aligned} \frac{\partial^2}{\partial \theta^2} h(\theta, p) &= p^{\frac{3}{2}} \left( \frac{\partial^2}{\partial N^2} f(N, p) \right) \left( \frac{\partial}{\partial \theta} N(\theta) \right)^2 = \\ &= -q\sqrt{p} \left( 2 + \left(\frac{1}{q}\right)^{N(\theta)} \ln q \right) \left( 1 - \frac{3}{2}p \right)^2. \quad (5.10) \end{aligned}$$

For a fixed  $p$ , consider the term  $\left(2 + \left(\frac{1}{q}\right)^{N(\theta)} \ln q\right)$ . Since  $N(\theta) \uparrow$ , it is upper bounded by  $\left(2 + \left(\frac{1}{q}\right)^{N(1)} \ln q\right)$ . By making use of the second derivative test, one can show that  $p \mapsto \left(2 + \left(\frac{1}{q}\right)^{N(1)} \ln q\right)$  is strictly concave and has negative derivative on  $\left(0, \frac{3-\sqrt{5}}{2}\right)$ . Therefore, it is lower bounded by its left limit at  $\frac{3-\sqrt{5}}{2}$ . The value of the latter is approximately equal to 0.93. It follows then from (5.10) that, for any fixed  $p \in \left(0, \frac{3-\sqrt{5}}{2}\right)$ ,  $\theta \mapsto \frac{\partial^2}{\partial \theta^2} h(\theta, p)$  is negative on  $[0, 1]$ . Hence, for any fixed  $p \in \left(0, \frac{3-\sqrt{5}}{2}\right)$ ,  $\theta \mapsto h(\theta, p)$  is concave and has therefore a decreasing derivative. Inspection of  $p \mapsto \frac{\partial}{\partial \theta} h(0, p)$  reveals that  $\theta \mapsto \frac{\partial}{\partial \theta} h(\theta, p)$  is negative for any fixed  $p$  meaning that the range of  $\theta \mapsto h(\theta, p)$  is equal to  $[h(1, p), h(0, p)]$ . Finally, omitting the details of a tedious exercise of verification that the range of  $p \mapsto h(1, p)h(0, p)$  lies in  $(-\infty, 0)$ , we finish proof of this step and conclude that (5.8) indeed holds.

*Step 3: the end of the proof.* By Step 2, for each  $p \in \left(0, \frac{3-\sqrt{5}}{2}\right)$ , there exists  $N(p) \in \left(\frac{1}{\sqrt{p}} - p, \frac{1}{\sqrt{p}} + 1 - \frac{5}{2}p\right)$  which solves  $\frac{\partial}{\partial N} t_{MD}(N, p) = 0$ . By Step 1,  $N(p) \in A^{MD}$ . Therefore, by Theorem 5.1.1,  $N(p)$  is the unique global minimizer of  $[1, \infty) \ni N \mapsto t_{MD}(N, p)$ , i.e.,  $N(p) = N_*^{MD}(p)$ . This implies that  $N_{opt}^{MD}(p) \in \{\lfloor \sqrt{p^{-1}} \rfloor - 1, \lfloor \sqrt{p^{-1}} \rfloor, \lfloor \sqrt{p^{-1}} \rfloor + 1\}$ , and it remains to exclude the point  $\lfloor \sqrt{p^{-1}} \rfloor - 1$ . The route is as follows. First, by making exactly the same technique as in Step 2, show that  $N(p) \in \left(\frac{1}{\sqrt{p}}, \frac{1}{\sqrt{p}} + 1 - \frac{5}{2}p\right)$  for  $p \in (0, 0.3)$ . Second, note that integer parts of  $\frac{1}{\sqrt{p}} - p$  and  $\frac{1}{\sqrt{p}}$  coincide for  $p \in [0.3, \frac{3-\sqrt{5}}{2}]$ . Figure 5.4 provides a good visual summary of the whole proof and explains the need of this workaround in particular.  $\square$

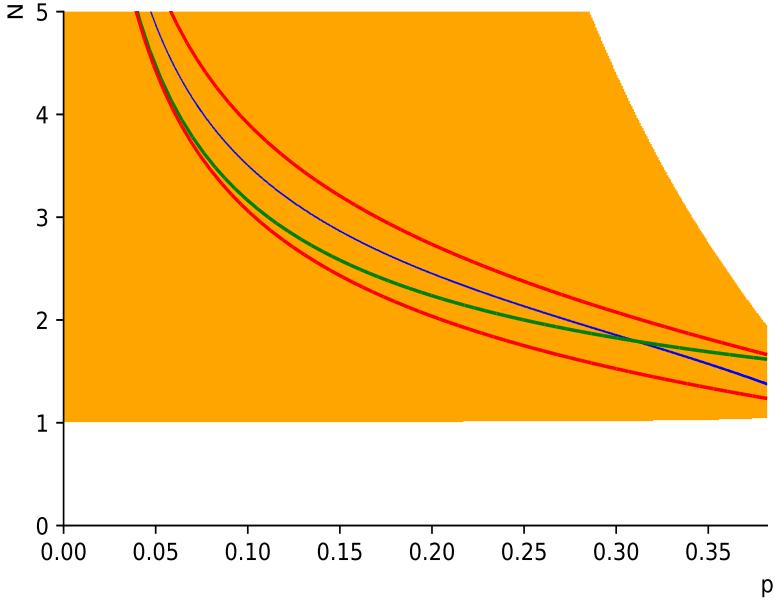


Figure 5.4: orange area corresponds to  $\{(N, p) : t_{MD}(N, p) < 1\}$ ; upper and lower red lines show bracing functions  $p \mapsto \sqrt{p^{-1}} + (1 - (5/2)p)$  and  $p \mapsto \sqrt{p^{-1}} - p$  respectively; green line shows function  $p \mapsto \sqrt{p^{-1}}$ ; blue line shows  $p \mapsto N_*^{MD}(p)$

We note that in the proof of Theorem 3.1.3 below, Step 1 is a repetition of Lemma 1 in [38]. Nonetheless, we have decided to rewrite it here

because it is very short and, along the way, some notions used in the sequel appear.

*Proof of Theorem 3.1.3.* *Step 1.* We first show that, for any fixed  $n \in (2, \infty)$ , there exists a unique  $q_n \in (0, 1)$  such that

$$g(q_n, n) = 0, \quad g(q, n) < 0 \quad \forall q \in (q_n, 1) \quad \text{and} \quad g(q, n) > 0 \quad \forall q \in (0, q_n).$$

To this end, note that

$$\begin{aligned} \frac{\partial}{\partial q} g(q, n) &= -2nq^{n-1} + (2n-1)q^{2n-2} = \\ &\quad -2nq^{n-1} \left(1 - \frac{2n-1}{2n}q^{n-1}\right) < 0 \quad \forall n \in (2, \infty). \end{aligned}$$

Thus, given  $n \in (2, \infty)$ ,  $q \mapsto g(q, n)$  is decreasing on  $(0, 1)$ . Since  $g(0+, n) = \frac{2}{n} > 0$  and  $g(1-, n) = \frac{2}{n} - 1 < 0$ , the claim holds true.

*Step 2.* From *Step 1* it follows that, for any fixed  $q \in (0, 1)$ , we have a well defined function  $n \mapsto q_n$  which is given implicitly by equation  $g(q_n, n) = 0$ . Since this function is continuous, its range  $I \subset (0, 1)$  is an interval. Next, for  $\varepsilon \in (0, 1)$  and  $n = 2 + \varepsilon$ ,

$$\begin{aligned} g(1 - \varepsilon, n) &= \frac{2}{n} - 2(1 - n\varepsilon + O(\varepsilon^2)) + (1 - (2n-1)\varepsilon + O(\varepsilon^2)) = \\ &\quad \frac{-\varepsilon}{2 + \varepsilon} + \varepsilon + O(\varepsilon^2) > 0, \end{aligned}$$

provided  $\varepsilon$  is small enough. Hence, analysis accomplished in *Step 1* implies that  $q_{2+\varepsilon} \in (1 - \varepsilon, 1)$ . Therefore,  $I = (q_*, 1)$  for some  $q_* \in (1/2, 1)$ . To justify the lower bound  $1/2$ , note that

$$g(1/2, n) > 0 \iff 2^n > n \left(1 - \frac{1}{2^n}\right).$$

Since the rhs holds true for all  $n > 2$ , it follows that  $\forall n > 2 \ g(1/2, n) > 0$ . Also note that  $\forall (q, n) \in (0, q_*) \times (2, \infty) \ g(q, n) > 0$  since an opposite contradicts the definition of  $q_*$ .

*Step 3.* Fix  $q \in (q_*, 1)$ . By *Step 1–Step 2*,  $g(q, n) = 0$  has at least one solution  $n = n(q)$  (suffices it to take  $n$  such that  $q_n = q$ ). To show that there are two solutions, put  $c = \frac{1}{2q}$ , make a change of variable  $q^n = x$ , and rewrite  $g(q, n) = 0$  in a form

$$-\ln q = -\ln x(x - cx^2). \tag{5.11}$$

Consider function  $h(x) = -(x - cx^2) \ln x$  for  $x \in (0, 1)$ . Note that

$$\frac{d}{dx} h(x) = -((1 - cx) + (1 - 2cx) \ln x) = 0 \iff -\ln x = \frac{1 - cx}{1 - 2cx}. \quad (5.12)$$

Since  $1 - cx > 1 - c > 0$  and  $-\ln x > 0$  for all  $x \in (0, 1)$ , it follows that  $1 - 2cx > 0$  as well, provided  $x$  solves (5.12). Therefore, the range of possible solutions of (5.12) shrinks to  $(0, q)$ . Moreover, relationships  $\lim_{x \rightarrow 0^+} \frac{d}{dx} h(x) = \infty$ ,  $\lim_{x \rightarrow 1^-} \frac{d}{dx} h(x) = c - 1 < 0$  imply that (5.12) has at least one solution. Since  $\frac{1-cx}{1-2cx} + \ln x = 1 + \frac{cx}{1-2cx} + \ln x$  increases on  $(0, q)$ , the solution is unique. Denote it  $x_0$ . Based on the sign of the derivative, we have that  $h \uparrow$  on  $(0, x_0)$  and  $h \downarrow$  on  $(x_0, 1)$ . Hence, at  $x_0$ ,  $h$  attains its maximum and (5.11) admits exactly two solutions if  $h(x_0) > -\ln q$ , one solution if  $h(x_0) = -\ln q$ , and has no solutions if  $h(x_0) < -\ln q$ . By the choice of  $q$  (recall that  $q > q_*$ ), the last case can not hold. To exclude the second one, note that the function  $\left(\frac{1}{2}, \frac{1}{2q_*}\right) \ni c \mapsto x_0(c)$  is well defined and decreasing since

$$\begin{aligned} -\ln x_0 &= \frac{1 - cx_0}{1 - 2cx_0} \Rightarrow \\ -\frac{d}{dc} \ln x_0(c) &= -\frac{\frac{d}{dc} x_0(c)}{x_0(c)} = \frac{d}{dc} \frac{1 - cx_0}{1 - 2cx_0} = \frac{x_0 + c \frac{d}{dc} x_0(c)}{(1 - 2cx_0)^2} \Rightarrow \\ \frac{d}{dc} x_0(c) &= \frac{-x_0^2}{(1 - 2cx_0)^2 + cx_0} < 0. \end{aligned}$$

Therefore,  $(q_*, 1) \ni q \mapsto x_0(q)$  is increasing. Taking into account that  $q \mapsto -\ln q$  is decreasing, we finally deduce that  $h(x_0) > -\ln q$  for all  $q \in (q_*, 1)$ . The monotonicity of  $q \mapsto x_0(q)$  and  $q \mapsto -\ln q$  also leads to the conclusion that  $q_*$  can be solved from equation

$$-\ln q_* = h(x_0(q_*))$$

along with a unique  $n_* \in (2, \infty)$ . Hence (i).

*Step 4.* Assume the setting of *Step 3*. Let  $0 < x_L = q^{n_U} < x_U = q^{n_L} < q$  denote two solutions of (5.11). By above,  $h(x) > -\ln q \iff x \in (x_L, x_U)$ . Reverting to  $(0, \infty) \ni n \mapsto g(q, n)$ , this reads as  $g(q, n) < 0 \iff n \in (n_L, n_U)$ . Note that  $n \mapsto g(q, n) > 0$  in the neighborhood of  $\infty$ . Also, from *Step 1–Step 2*, we have that  $n \mapsto g(q, n) > 0$  in the right neighborhood of 2 and that  $n_U \in (2, \infty)$ . Therefore, continuity of  $n \mapsto g(q, n) > 0$  yields that  $n_L \in (2, \infty)$  as well. Finally, it is clear that  $n_L < n_* < n_U$  (see figure 5.5 for a graphical illustration). Hence (ii).

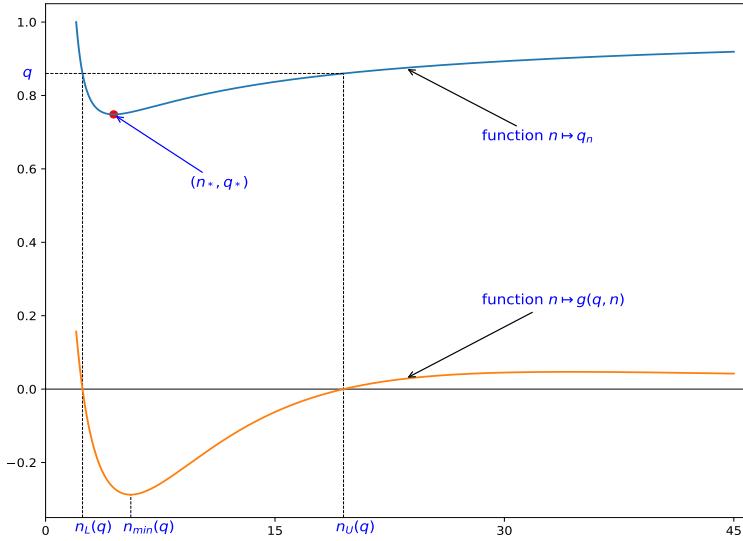


Figure 5.5: Graph illustrating relationships of  $n \mapsto q_n$  and related quantities. The lower curve corresponds to  $q = 0.86$ .

*Step 5.* In this step, we identify the number and location of zeroes of the derivative of  $(2, \infty) \ni n \mapsto g(q, n)$  having fixed  $q \in (q_*, 1)$ . From analysis given in *Step 4*, it follows that  $(2, \infty) \ni n \mapsto g(q, n)$  has at least two extremes: there must be a minimum in  $(n_L, n_U)$  (since function is negative here), and maximum in  $(n_U, \infty)$  (since the function is positive here and  $\lim_{n \rightarrow \infty} g(q, n) = 0+$ ). To see that there are no other extremes, consider an equation

$$\frac{\partial}{\partial n} g(q, n) = -\frac{2}{n^2} - 2q^n \ln q + 2q^{2n-1} \ln q = 0, \quad (5.13)$$

and rewrite it by making use of notions introduced in *Step 3* as follows:

$$-x \ln x = \sqrt{-\ln q} \sqrt{\frac{x}{1-2cx}}, x \in \left(0, \frac{1}{2c}\right) = (0, q). \quad (5.14)$$

Next, note that:

- $h_1(x) = -x \ln x$  is strictly convex-up and positive on  $(0, 1)$  with  $\lim_{x \rightarrow 0+} h_1(x) = \lim_{x \rightarrow 1-} h_1(x) = 0+$ ;

- $h_2(x) = \sqrt{\frac{x}{1-2cx}}$  is strictly positive and increasing on  $(0, \frac{1}{2c}) = (0, q)$ , it has one inflection point and  $\lim_{x \rightarrow 0^+} h_2(x) = 0$ ,  $\lim_{x \rightarrow q^-} h_2(x) = \infty$ .

Taking this information into account, we conclude that (5.14) can have at most two solutions and confirm thereby the assertion stated above.

*Step 6.* It remains to justify expression (3.4). Let

$$n(q, t) = \frac{1}{p^{\frac{2}{3}}} + \frac{1}{2p^{\frac{1}{3}}} + 0.2 + 3p^2 + t, \quad t \in [0, 1]. \quad (5.15)$$

It suffices to prove that, for all  $q \in [0.755, 1)$ , the following statements hold true:

(a)  $\max(g(q, n(q, 0)), g(q, n(q, 1))) < 0$ ; and

(b)  $\exists t \in [0, 1] : \left. \frac{\partial}{\partial n} g(q, n) \right|_{n=n(q,t)} = 0$ .

Analytical calculations behind (a) and (b) are standard yet very lengthy and tedious. Therefore, we omit the details and end up with a graphical proof and a sketch of the analytical one.

Figures 5.6–5.7 show graph of  $q \mapsto \max(g(q, n(q, 0)), g(q, n(q, 1)))$  from which it is evident that (a) holds. Analytical proof consists of the following steps.

(s1) Calculate  $\frac{\partial}{\partial q} g(q, n(q, i))$ ,  $i = 0, 1$ .

(s2) Check that  $\frac{\partial}{\partial q} g(q, n(q, i)) < 0$ ,  $i = 0, 1$  on  $[0.755, 1)$  and deduce that  $g(q, n(q, i))$  decrease on  $[0.755, 1)$ .

(s3) Conclude that (a) indeed holds since

$$g(0.755, n(0.755, 0)) \approx -0.002258, \quad g(0.755, n(0.755, 1)) \approx -0.013690.$$

Turning to (b), first rewrite (5.13) as follows:

$$-n^2 q^n \ln q (1 - q^{n-1}) = 1.$$

Next, consider function  $h(t, q) = -n^2(q, t)q^{n(q,t)} \ln q (1 - q^{n(q,t)-1}) - 1$  with  $n(q, t)$  given by (5.15) and  $q \in [0.755, 1)$ . Since  $t \mapsto h(t, q)$  is continuous, it suffices to show that, for any  $q \in [0.755, 1)$ ,  $h(0, q) < 0$  and  $h(1, q) > 0$ . Figure 5.8 shows graphs of  $q \mapsto h(0, q)$ ,  $q \mapsto h(1, q)$ . These confirm (b). Considering analytical part, the following is the suggested route.

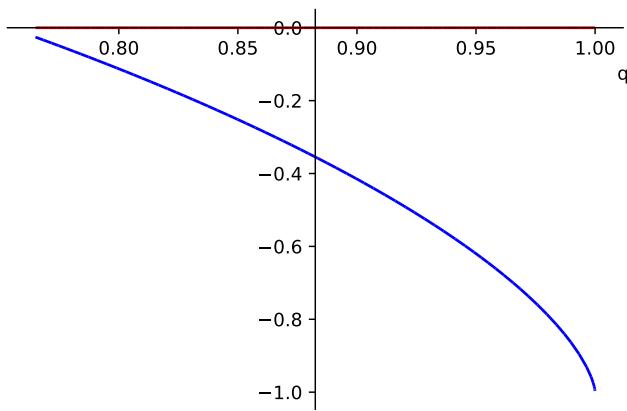


Figure 5.6: Graph of  $q \mapsto \max(g(q, n(q, 0)), g(q, n(q, 1)))$  for  $q \in [0.765, 1)$ . For reference, a function identically equal to 0 is plotted in red.

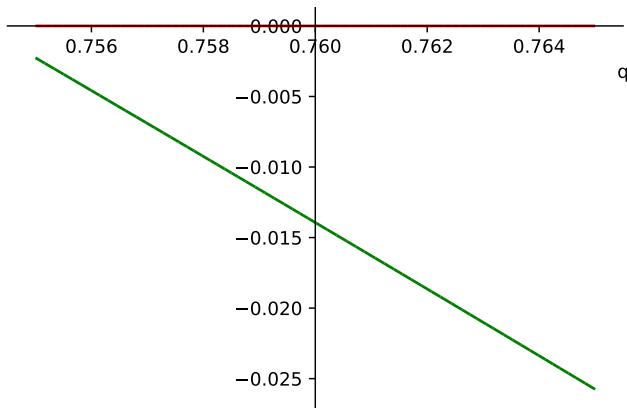


Figure 5.7: Graph of  $q \mapsto \max(g(q, n(q, 0)), g(q, n(q, 1)))$  for  $q \in [0.755, 0.765]$ . For reference, a function identically equal to 0 is plotted in red.

- (s1) Calculate  $\frac{\partial^2}{\partial q^2} h(i, q), i = 0, 1$ .
- (s2) Check that  $\frac{\partial^2}{\partial q^2} h(0, q) > 0$  whereas  $\frac{\partial^2}{\partial q^2} h(1, q) < 0$  on  $[0.755, 1)$  and

deduce that  $h(0, q)$  is convex downwards whereas  $h(1, q)$  is convex upwards on  $[0.755, 1]$ .

- (s3) By making use of Taylor's expansion, check that  $\lim_{q \rightarrow 1^-} h(0, q) = 0-$  and  $\lim_{q \rightarrow 1^-} h(1, q) = 0+$ .

- (s4) Conclude that (b) indeed holds since

$$h(0, 0.755) \approx -0.2645889 \quad \text{and} \quad h(1, 0.755) \approx 0.081749. \quad \square$$

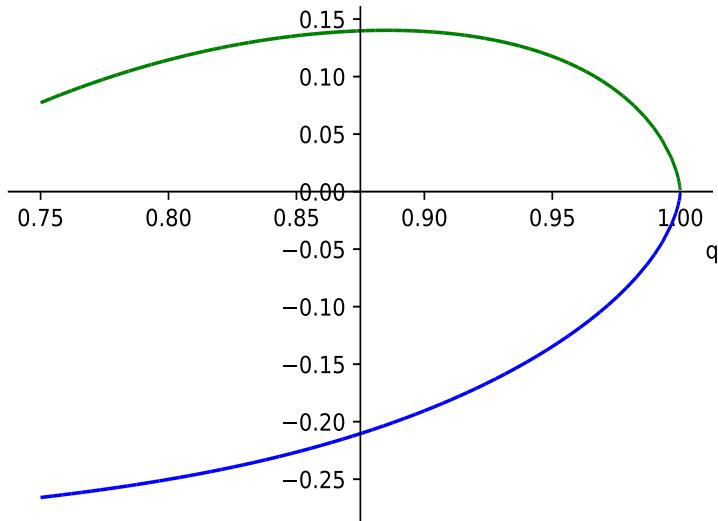


Figure 5.8: Graphs of  $q \mapsto h(0, q), q \in [0.755, 1]$  (plotted in blue) and  $q \mapsto h(1, q), q \in [0.755, 1]$  (plotted in green).

*Proof of Corollary 3.1.4.* Uniqueness of  $q_5$  was established in Step 1 of the proof of Theorem 3.1.3. Hudgens and Kim [38] (Lemmas 2, 7, and 14) have demonstrated that  $n_{opt}(q) \notin \{2, 3, 4\} \forall q \in (0, 1)$ . From their results we also have that  $\forall q \in (q_*, 0.755] n_{opt}(q) = 5$ . It is straightforward to verify that

$$\forall q \in (q_5, 0.755] \left[ \frac{1}{p^{\frac{2}{3}}} + \frac{1}{2p^{\frac{1}{3}}} + 3p^2 + 1.2 \right] = 5.$$

Hence the claim in the region  $(q_5, 0.755]$ . For  $q \in [0.755, 1)$ , it follows from Theorem 3.1.3 by noting that at least one of numbers in the set (3.5) belongs to  $\{n \in (2, \infty) : t_{A2}(n, q) < 1\}$  because of (a) in Step 6 of proof of Theorem 3.1.3.  $\square$

## 5.2 Proofs of results stated in Section 3.2

*Proof of Proposition 3.2.1.* Take  $p \in (0, p_{X,c})$ . By the definition of  $p_{X,c}$ , there exists  $N \in (c, \infty) : t_X(N, p) < 1$ . Hence,  $X$  makes sense for that  $p$ . Since  $p \in (0, p_{X,c})$  was arbitrary, it holds true for all  $p \in (0, p_{X,c})$ .

Next, assume that  $p_{X,c} < \text{UCP}$  (otherwise implication (3.11) is obvious) and take  $p \in (p_{X,c}, \text{UCP}]$ . Case " $\exists N \in (c, \infty) : t_X(N, p) < 1$ " contradicts the definition of  $p_{X,c}$ . Hence,  $\forall N \in (c, \infty) t_X(N, p) \geq 1$ . Assuming that  $t_X(N, p) = 1$  for some  $N \in (c, \infty)$  again leads to contradiction. Indeed, take  $p' \in (p_{X,c}, p)$  and employ (M2) to deduce that

$$\begin{aligned} t_X(N, p') < t_X(N, p) = 1 \Rightarrow p_{X,c} = \\ \sup\{p \in (0, \text{UCP}) \mid \exists N \in (c, \infty) : t_X(N, p) < 1\} \geq p' > p_{X,c}. \quad \square \end{aligned} \tag{5.16}$$

*Proof of Proposition 3.2.2.* First, note that  $t_X(N, p_{X,c}) \geq 1$  for all  $N \in (c, \infty)$  provided  $p_{X,c} < \text{UCP}$ . Assuming an opposite leads to the same contradiction as in (5.16). With this in view, we proceed to the analysis of the distinct types of possible bifurcations.

*Case (b0).* Since  $p_{X,c}$  solves  $t_X(N, p_{X,c}) = 1$  for some  $N \in (c, \infty)$ , (M3) implies that  $p_{X,c} < \text{UCP}$ . Further, note that

$$\forall (N, p) \in (c, \infty) \times (0, p_{X,c}) t_X(N, p) < 1. \tag{5.17}$$

must hold since the existence of  $(N, p) \in (c, \infty) \times (0, p_{X,c})$  s.t.  $t_X(N, p) \geq 1$  contradicts the premise<sup>2</sup> " $p_{X,c}$  is the only value of the control parameter for which dynamical system (3.13) admits fixed points in  $(c, \infty)$ ".

Fix arbitrary  $(N_1, p_1) \in (c, \infty) \times (0, p_{X,c})$  and take  $p_2 \in (p_{X,c}, \text{UCP})$ . By Prop. 3.2.1,  $t_X(N_1, p_2) > 1$ . By (M2),  $[p_1, p_2] \ni p \mapsto t_X(N_1, p)$  is strictly increasing. Therefore, there exists unique  $p_0 \in (p_1, p_2)$  s.t.

---

<sup>2</sup> $t_X(N, p) = 1$  is clearly impossible; assuming strict inequality  $t_X(N, p) > 1$ , to reach a contradiction, employ (M2) and (M4)

$t_X(N_1, p_0) = 1$ . By the premise,  $p_0 = p_{X,c}$ . Since  $N_1$  was arbitrary, it follows that

$$\forall N \in (c, \infty) t_X(N, p_{X,c}) = 1 \quad (5.18)$$

and  $p_{X,c}$  is the unique value of the control parameter  $p \in (0, \text{UCP})$  obeying this property. Differentiating both sides of (5.18) with respect to  $N$  one finds out that (3.13) holds as well.

*Case (b1)–(b2).* Assume that there exists  $(N_l, p_l) \in (c, \infty) \times (0, p_{X,c})$  s.t.  $t_X(N_l, p_l) = 1$ . Take arbitrary  $p \in (p_l, p_{X,c})$ . Since  $p > p_l$ , it follows that  $t_X(N_l, p) > 1$  because of (M2). On the other hand, by the definition of  $p_{X,c}$  there exists  $N_p \in (c, \infty)$  s.t.  $t_X(N_p, p) < 1$ . Therefore, by the Intermediate Value Theorem and continuity of  $N \mapsto t_X(N, p)$ , there exists  $N_1 \in (\min(n_p, N_l), \max(n_p, N_l))$  s.t.  $t_X(N_1, p) = 1$ . Since this holds for any  $p \in [p_l, p_{X,c}]$ , we have (b1) provided  $\forall N \in (c, \infty) t_X(N, p_{X,c}) > 1$ . Otherwise, we have (b2) and  $p_{X,c}$  then can't be equal to UCP because of (M3).

*Inversion of the bifurcation curve.* When dealing with (b0), we have already shown that the bifurcation curve defines a constant map  $(c, \infty) \ni N \mapsto p_N \equiv p_{X,c}$ . As for (b1)–(b2), note that, for any fixed  $N \in (c, \infty)$ , the following applies:

- by the said in the very beginning of the proof and (M3),

$$\forall N t_X(N, p_{X,c}) \geq 1; \quad (5.19)$$

- by (M4),

$$\forall N \exists p \in (0, p_{X,c}) : t_X(N, p) < 1; \quad (5.20)$$

- (5.19)–(5.20) and (M2) imply existence of a unique  $p_N \in (0, p_{X,c}]$  s.t.  $t_X(N, p_N) = 1$ .

Therefore, we have a well defined map  $(c, \infty) \ni N \mapsto p_N \in (0, p_{X,c}]$ . By (M1)–(M2),  $p \mapsto t_X(N, p)$  is differentiable and increasing for any  $N \in (c, \infty)$ . Therefore,  $\forall N \in (c, \infty) \frac{\partial}{\partial p} t_X(N, p) > 0$  and one can apply the Implicit Function Theorem to  $\varphi(N, p) = t_X(N, p) - 1$  to deduce that  $N \mapsto p_N$  is differentiable as well. Moreover, differentiating both sides of  $t_X(N, p_N) = 1$  and applying the chain rule, we have that

$$\frac{\partial}{\partial N} t_X(N, p_N) + \frac{\partial}{\partial p} t_X(N, p_N) \frac{\partial}{\partial N} p_N = 0 \Rightarrow \frac{\partial}{\partial N} p_N = -\frac{\frac{\partial}{\partial N} t_X(N, p_N)}{\frac{\partial}{\partial p} t_X(N, p_N)}.$$

(5.21)

Therefore, looking for extremes of  $N \mapsto p_N$  one has to solve

$$\frac{\partial}{\partial N} t_X(N, p_N) = 0 \iff \frac{\partial}{\partial N} p_N = 0$$

with respect to  $N$ . Since  $p_N$  also solves  $t_X(N, p) = 1$ , extremes and corresponding values can be obtained by solving (3.13). For bifurcations of type (b0) and (b2), maximal value  $p_{X,c}$  is attained at some inner point(s)  $N_c \in (c, \infty)$ ; for the bifurcation of type (b1), the maximal value lies on the boundary of its domain.  $\square$

### 5.3 Proofs of results stated in Section 3.3

*Proof of Proposition 3.3.1.* By the description of the testing procedure,

$$\begin{aligned} \Theta_N = & (1 + \Theta_{N-2}) \mathbb{1}\{Y_N + Y_{N-1} = 0\} + \\ & (2 + \Theta_{N-1}) \mathbb{1}\{Y_N + Y_{N-1} > 0\} Y_N + \\ & (2 + \Theta_{N-2}) \mathbb{1}\{Y_N + Y_{N-1} > 0\} \bar{Y}_N = \\ & \bar{Y}_N(1 + Y_{N-1}) + 2Y_N + \Theta_{N-1}Y_N + \Theta_{N-2}\bar{Y}_N \end{aligned} \quad (5.22)$$

since

$$\begin{aligned} \mathbb{1}\{Y_N + Y_{N-1} = 0\} &= \bar{Y}_N \bar{Y}_{N-1}, \quad \mathbb{1}\{Y_N + Y_{N-1} > 0\} Y_N = Y_N, \\ \text{and} \quad \mathbb{1}\{Y_N + Y_{N-1} > 0\} \bar{Y}_N &= \bar{Y}_N Y_{N-1}. \end{aligned}$$

Let

$$\tau_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \quad A_1 = \tau_1, \quad \tau_k \stackrel{k \geq 2}{=} \begin{pmatrix} \Theta_k \\ \Theta_{k-1} \end{pmatrix}, \quad A_k \stackrel{k \geq 2}{=} \begin{pmatrix} \bar{Y}_k(1 + Y_{k-1}) + 2Y_k \\ 0 \end{pmatrix},$$

and let  $B_k$ ,  $k \geq 1$ , be as in the statement of the Proposition. From (5.22) it follows that

$$\begin{aligned} \tau_N &= A_N + B_N \tau_{N-1} = \dots = \\ A_N + \sum_{k=1}^{N-2} B_N B_{N-1} \dots B_{N-k+1} A_{N-k} + B_N \dots B_2 \tau_1 &= \\ A_N + \sum_{j=3}^N B_N \dots B_j A_{j-1} + B_N \dots B_2 \tau_1 &= A_N + \sum_{j=2}^N B_N \dots B_j A_{j-1}. \end{aligned}$$

Define

$$M_0 = \begin{pmatrix} 1 & 0 \\ 1 & 0 \end{pmatrix}, M_1 = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}, M_2 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, M_3 = \begin{pmatrix} 0 & 1 \\ 0 & 1 \end{pmatrix},$$

$$S = \{M_0, M_1, M_2, M_3\}, \quad (5.23)$$

Then  $S$  forms a semi-group with respect to ordinary matrix multiplication since

$$M_0^2 = M_0, M_0M_1 = M_3, M_0M_3 = M_3, M_1M_0 = M_0, M_1^2 = M_2,$$

$$M_1M_3 = M_3, M_3M_0 = M_0, M_3M_1 = M_0, M_3^2 = M_3. \quad (5.24)$$

Let  $J_i = \{j \in \{2, \dots, n\} \mid Y_j = i\}$ ,  $i = 0, 1$ . Note that  $\forall i B_i = Y_i M_0 + \bar{Y}_i M_1 \in S$  and that  $M_0$  is an absorbing element of  $S$ . Therefore, by (5.24)

$$\sum_{j \in J_1} B_N \dots B_j A_{j-1} = \sum_{j \in J_1} M_0 A_{j-1} = \sum_{j=2}^N Y_j M_0 A_{j-1}$$

and

$$\sum_{j \in J_0} B_N \dots B_j A_{j-1} = \sum_{j=2}^N \bar{Y}_j \left( \mathbb{1}_{\{B_N \dots B_j M_1 = M_0\}} M_0 + \mathbb{1}_{\{B_N \dots B_j M_1 = M_1\}} M_1 + \mathbb{1}_{\{B_N \dots B_j M_1 = M_2\}} M_2 + \mathbb{1}_{\{B_N \dots B_j M_1 = M_3\}} M_3 \right) A_{j-1}.$$

To extract  $\Theta_N$  from  $\tau_N$ , it suffices to multiply  $\tau_N$  by (1 0) from the left. Since

$$(1 \ 0) M_i A_{j-1} = \begin{cases} 0, & \text{for } i = 1, 3 \text{ and all } j; \\ \bar{Y}_{j-1}(1 + Y_{j-2}) + 2Y_{j-1}, & \text{for } i = 0, 2 \text{ and } j \geq 3; \\ 1, & \text{for } i = 0, 2 \text{ and } j = 2, \end{cases}$$

after the collection of terms, we finally end up with an expression (3.22).

□

*Proof of Theorem 3.3.2. Step 1: auxiliary recurrence.* For  $\lambda_1, \lambda_2 \in \mathbb{R}$ , let

$$M_{i,N}(\lambda_1, \lambda_2) = E \left( e^{\lambda_1 \Theta_N + \lambda_2 \Theta_{N-1}} \mid X_N = i \right), \quad i = 0, 1. \quad (5.25)$$

By equation (5.22),

$$\begin{aligned}
M_{0,N}(\lambda_1, \lambda_2) &= \mathbb{E} \left( e^{\lambda_1(1+X_{N-1}+\Theta_{N-2})+\lambda_2\Theta_{N-1}} \right) = \\
&\quad p \mathbb{E} \left( e^{\lambda_1(2+\Theta_{N-2})+\lambda_2\Theta_{N-1}} \mid X_{N-1} = 1 \right) + \\
&\quad q \mathbb{E} \left( e^{\lambda_1(1+\Theta_{N-2})+\lambda_2\Theta_{N-1}} \mid X_{N-1} = 0 \right) = \\
&\quad pe^{2\lambda_1} M_{1,N-1}(\lambda_2, \lambda_1) + qe^{\lambda_1} M_{0,N-1}(\lambda_2, \lambda_1); \\
M_{1,N}(\lambda_1, \lambda_2) &= \mathbb{E} \left( e^{\lambda_1(2+\Theta_{N-1})+\lambda_2\Theta_{N-1}} \right) = \\
&\quad e^{2\lambda_1} (pM_{1,N-1}(\lambda_1 + \lambda_2, 0) + qM_{0,N-1}(\lambda_1 + \lambda_2, 0)). \tag{5.26}
\end{aligned}$$

For  $\lambda \in \mathbb{R}$ , let

$$\begin{aligned}
m_{1,N} &= m_{1,N}(\lambda) = M_{1,N}(\lambda, 0), & m_{2,N} &= m_{2,N}(\lambda) = M_{0,N}(\lambda, 0), \\
m_{3,N} &= m_{3,N}(\lambda) = M_{1,N}(0, \lambda), & m_{4,N} &= m_{4,N}(\lambda) = M_{0,N}(0, \lambda); \\
A &= A(\lambda) = e^{2\lambda} \begin{pmatrix} p & q \\ 0 & 0 \end{pmatrix}, & B &= B(\lambda) = \begin{pmatrix} 0 & 0 \\ e^{2\lambda}p & e^\lambda q \end{pmatrix}, \\
C &= \begin{pmatrix} p & q \\ p & q \end{pmatrix}, & O &= \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}. \tag{5.27}
\end{aligned}$$

From (5.26) it then follows that  $m_N = m_N(\lambda) = (m_{1,N}, m_{2,N}, m_{3,N}, m_{4,N})^\top$  satisfies recurrent equation

$$m_N = \begin{pmatrix} A & B \\ C & O \end{pmatrix} m_{N-1} = \cdots = \begin{pmatrix} A & B \\ C & O \end{pmatrix}^{N-1} m_1. \tag{5.28}$$

Writing

$$\begin{pmatrix} A & B \\ C & O \end{pmatrix}^N = \begin{pmatrix} A_N & B_N \\ C_N & D_N \end{pmatrix}$$

and applying inductive argument, one finds out that the  $2 \times 2$  blocks  $A_N, B_N, C_N, D_N$  satisfy

$$\begin{cases} A_N = AA_{N-1} + BC_{N-1}, \\ C_N = CA_{N-1}; \end{cases} \tag{5.29}$$

$$\begin{cases} B_N = AB_{N-1} + BD_{N-1}, \\ D_N = CB_{N-1}; \end{cases} \tag{5.30}$$

with  $A_0 = D_0 = Id$  and  $C_0 = B_0 = O$ . Consider system (5.29). Since  $A = \begin{pmatrix} e^{2\lambda} & 0 \\ 0 & 0 \end{pmatrix} C$ , we have that

$$A_N = \begin{pmatrix} e^{2\lambda} & 0 \\ 0 & 0 \end{pmatrix} CA_{N-1} + BC_{N-1} = \begin{pmatrix} e^{2\lambda} & 0 \\ 0 & 0 \end{pmatrix} C_N + BC_{N-1}. \quad (5.31)$$

Therefore,

$$C_N = C \left( \begin{pmatrix} e^{2\lambda} & 0 \\ 0 & 0 \end{pmatrix} C_{N-1} + BC_{N-2} \right). \quad (5.32)$$

Let  $\kappa_N$  be defined by (3.24). We claim that  $C_N = \kappa_N C$  solves (5.32). For  $N = 2$  (as well as  $N = 0, 1$ ) the claim holds by the direct check. Assume it holds for  $k \leq N$  with  $N \geq 2$ . Noting that

$$C \begin{pmatrix} e^{2\lambda} & 0 \\ 0 & 0 \end{pmatrix} C = pe^{2\lambda} C, \quad CBC = qe^\lambda(q + pe^\lambda)C$$

and then applying inductive assumption and multiplication yields

$$\begin{aligned} C_{N+1} &= C \left( \kappa_N \begin{pmatrix} e^{2\lambda} & 0 \\ 0 & 0 \end{pmatrix} C + \kappa_{N-1} BC \right) = \\ &\quad (pe^{2\lambda}\kappa_N + qe^\lambda(q + pe^\lambda)\kappa_{N-1})C = \kappa_{N+1}C \end{aligned}$$

since an expression for  $\kappa_N$  given in (3.24) is precisely the solution of the second order linear difference equation

$$\kappa_{N+1} = pe^{2\lambda}\kappa_N + qe^\lambda(q + pe^\lambda)\kappa_{N-1}, \quad \kappa_1 = 1, \quad \kappa_0 = 0.$$

Substituting  $C_N = \kappa_N C$  to (5.31), we obtain an expression for  $A_N$ .

System (5.30) is handled in the same way by noting that it is identical to (5.29) and only the initial conditions differ leading thereby to the following solution:

$$D_N = \kappa_{N-1}D_2, \quad B_N = \begin{pmatrix} e^{2\lambda} & 0 \\ 0 & 0 \end{pmatrix} D_N + BD_{N-1} \text{ for } N \geq 1. \quad (5.33)$$

Step 2: final expression. From the results of Step 1, we obtain an expression for  $m_n$  given by (5.28) since  $m_1$  is readily available and equal to<sup>3</sup>  $(e^\lambda, e^\lambda, 1, 1)^\top$ :

$$m_n = \begin{pmatrix} (e^\lambda A_{N-1} + B_{N-1}) \begin{pmatrix} 1 \\ 1 \end{pmatrix} \\ (e^\lambda C_{N-1} + D_{N-1}) \begin{pmatrix} 1 \\ 1 \end{pmatrix} \end{pmatrix}.$$

---

<sup>3</sup>note that  $\Theta_1 \equiv 1, \Theta_0 \equiv 0$

Noting that

$$\mathbb{E} e^{\lambda \Theta_N} = p \mathbb{E} (e^{\lambda \Theta_N} | X_N = 1) + q \mathbb{E} (e^{\lambda \Theta_N} | X_N = 0) = pm_{1,N} + qm_{2,N},$$

we finally arrive to expression (3.25).  $\square$

*Proof of Corollary 3.3.3.* Recall that the  $k$ -th derivative of the moment generating function evaluated at zero yields the  $k$ -th moment. Therefore, to obtain the announced formulae, one simply needs to differentiate expression (3.25). Though conceptually an exercise is trivial, the calculations require tedious work. Therefore, we provide key steps and some intermediate quantities, yet omit the detailed listing in order not to overwhelm the Thesis with trivial content. For the sake of convenience, we make a change of variables  $x = e^\lambda$  and work with the probability generating function  $G(x) = \mathbb{E} x^{\Theta_N} = M_{\Theta_N}(\ln x)$ . By (3.24)–(3.25) and a slight abuse of notation,

$$\begin{aligned} G(x) &= g_1(x)\kappa_{N-2}(x) + g_2(x)\kappa_{N-3}(x) \text{ with} \\ g_1(x) &= ((1-q)^2x^3 + q(1-q)^2x^2 + q(1-q^2)x + q^2)x^2, \\ g_2(x) &= q((1-q)^2x^3 + q(1-q)(2-q)x^2 + 2q^2(1-q)x + q^3)x^2, \\ \alpha_i &= \alpha_i(x) = \frac{1}{2} \left( px^2 + (-1)^i \sqrt{p^2x^4 + 4qx(q+px)} \right), \text{ for } i = 0, 1, \\ \text{and } \kappa_N &= \kappa_N(x) = \frac{\alpha_0^N(x) - \alpha_1^N(x)}{\alpha_0(x) - \alpha_1(x)} \text{ for } N \geq 0. \end{aligned} \quad (5.34)$$

Then

$$\begin{aligned} \mathbb{E} \Theta_N &= G'(1) = g'_1(1)\kappa_{N-2}(1) + g_1(1)\kappa'_{N-2}(1) + \\ &\quad g'_2(1)\kappa_{N-3}(1) + g_2(1)\kappa'_{N-3}(1) \end{aligned} \quad (5.35)$$

and

$$\begin{aligned} \mathbb{E} \Theta_N(\Theta_N - 1) &= G''(1) = g''_1(1)\kappa_{N-2}(1) + 2g'_1(1)\kappa'_{N-2}(1) + \\ &\quad g_1(1)\kappa''_{N-2}(1) + g''_2(1)\kappa_{N-3}(1) + 2g'_2(1)\kappa'_{N-3}(1) + g_2(1)\kappa''_{N-3}(1). \end{aligned} \quad (5.36)$$

Therefore,  $\text{Var } \Theta_N = G''(1) + G'(1) - (G'(1))^2$  and to verify the an-

nounced formulae, one needs to check the validity of the equalities

$$\begin{aligned}
\alpha_0(1) &= 1, \quad \alpha_1(1) = -q, \quad \alpha'_0(1) = \frac{2-q^2}{1+q}, \quad \alpha'_1(1) = -\frac{q^2}{1+q}, \\
\alpha''_0(1) &= 4\frac{1-q}{q+1} - \frac{2}{(q+1)^3}, \quad \alpha''_1(1) = -\frac{2(1-q)^2}{q+1} + \frac{2}{(q+1)^3}, \\
g_1(1) &= 1, \quad g_2(1) = q, \\
g'_1(1) &= q^3 - q^2 - 3q + 5, \quad g'_2(1) = -q(q^2 + 2q - 5), \\
g''_1(1) &= 6q^3 - 2q^2 - 22q + 20, \quad g''_2(1) = 2q(q^3 - 2q^2 - 8q + 10), \\
\kappa_N(1) &= \frac{1 - (-q)^N}{1+q}, \\
\kappa'_N(1) &= N\frac{2-q^2}{(1+q)^2} + \frac{(-q)^N(2-q(1+q)N) - 2}{(1+q)^3}, \\
\kappa''_N(1) &= \frac{2N(1-q)\left(2 + (1-q)(-q)^{N-1}\right)}{(q+1)^2} + \\
&\frac{N(N-1)\left((2-q^2)^2 - (-q)^{N+2}\right) - 2(1-q)(3-q)\left(1 - (-q)^N\right)}{(q+1)^3} - \\
&\frac{2N\left(-2q^2 + 5 + (-q)^{N-1}(2q^2 + 1)\right)}{(q+1)^4} + 12\frac{1 - (-q)^N}{(q+1)^5},
\end{aligned}$$

substitute them into (5.35)–(5.36), and carefully compute the terms.  $\square$

*Proof of corollary 3.3.4.* Step 1: expansions. Applying Taylor's formula, we obtain the following equalities (for  $\lambda \rightarrow 0$ ):

$$\begin{aligned}
p^2 e^{4\lambda} + 4qe^\lambda (q + pe^\lambda) &= \\
(1+q)^2 \left[ 1 + \frac{4\lambda}{(1+q)^2} + 2\lambda^2 \left( \frac{2-q}{1+q} \right)^2 + O(\lambda^3) \right]; \\
\sqrt{p^2 e^{4\lambda} + 4qe^\lambda (q + pe^\lambda)} &= \\
(1+q) \left( 1 + \frac{2\lambda}{(1+q)^2} + \lambda^2 \left( \left( \frac{2-q}{1+q} \right)^2 - \frac{2}{(1+q)^4} \right) + O(\lambda^3) \right); \\
\alpha_0 &= 1 + \lambda \frac{2-q^2}{1+q} + \frac{\lambda^2}{2} \left( 2(1-q) + \frac{(2-q)^2}{(1+q)} - \frac{2}{(1+q)^3} \right) + O(\lambda^3); \\
\alpha_1 &= -q - \lambda \frac{q^2}{1+q} + \frac{\lambda^2}{2} \left( 2(1-q) - \frac{(2-q)^2}{(1+q)} + \frac{2}{(1+q)^3} \right) + O(\lambda^3).
\end{aligned} \tag{5.37}$$

Let  $c_{ij}$  denote a coefficient near  $\lambda^j$  in the expansion of  $\frac{\alpha_i}{(-q)^i}$  for  $j = 0, 1, 2$  and  $i = 0, 1$ . Then

$$\ln \left( \frac{\alpha_i}{(-q)^i} \right)^N = N \left( c_{i1}\lambda + \left( c_{i2} - \frac{c_{i1}^2}{2} \right) \lambda^2 \right) + O(N\lambda^3). \quad (5.38)$$

Consequently,

$$\begin{aligned} (\alpha_0 - \alpha_1)\kappa_N(\lambda) &= \alpha_0^N - \alpha_1^N = e^{N \ln \alpha_0} - (-1)^N e^{N \ln(q \frac{\alpha_1}{-q})} = \\ &\exp \left\{ N \left( c_{01}\lambda + \left( c_{02} - \frac{c_{01}^2}{2} \right) \lambda^2 \right) + O(N\lambda^3) \right\} - \\ &(-1)^N \exp \left\{ N \left( c_{11}\lambda + \left( c_{12} - \frac{c_{11}^2}{2} \right) \lambda^2 \right) + N \ln q + O(N\lambda^3) \right\}. \end{aligned} \quad (5.39)$$

Finally, let  $g_i(x)$  denote the same polynomials as given in (5.34). Taylor expanding yields

$$g_1(e^\lambda) = 1 + O(\lambda), \quad g_2(e^\lambda) = q + O(\lambda).$$

Combining all above, we then obtain the following asymptotic expansion for the moment-generating function:

$$M_{\Theta_N}(\lambda) = \frac{1 + O(\lambda)}{1 + q + O(\lambda)} ((1 + O(\lambda))\kappa_{N-2}(\lambda) + (q + O(\lambda))\kappa_{N-3}(\lambda)) \quad (5.40)$$

with asymptotic expressions for  $\kappa_{N-2}, \kappa_{N-3}$  stemming from (5.39).

Step 2: LLN. To prove relationship  $\frac{\Theta_N}{N} \xrightarrow{L_2} \frac{2-q^2}{1+q}$ , note that, by Corollary 3.3.3,

$$\begin{aligned} E \left( \frac{\Theta_N}{N} - \frac{2-q^2}{1+q} \right)^2 &= E \left( \frac{\Theta_N}{N} - E \frac{\Theta_N}{N} \right)^2 + \left( E \frac{\Theta_N}{N} - \frac{2-q^2}{1+q} \right)^2 = \\ &\frac{1}{N^2} \left( \text{Var } \Theta_N + \left( \frac{q^2+q-1}{(1+q)^2} (1 - (-q)^N) \right)^2 \right) = O \left( \frac{1}{N} \right). \end{aligned}$$

To prove a.s. convergence, we bound the probability

$$\begin{aligned} P \left( \left| \frac{\Theta_N}{N} - \frac{2-q^2}{1+q} \right| > \gamma \frac{\ln N}{\sqrt{N}} \right) &= \\ P \left( \frac{\Theta_N}{\sqrt{N}} - \sqrt{N} \frac{2-q^2}{1+q} > \gamma \ln N \right) &+ \\ P \left( \frac{\Theta_N}{\sqrt{N}} - \sqrt{N} \frac{2-q^2}{1+q} < -\gamma \ln N \right), \end{aligned}$$

where  $\gamma > 0$  is arbitrary yet fixed constant. By Markov's inequality,

$$\begin{aligned} \mathrm{P}\left(\frac{\Theta_N}{\sqrt{N}} - \sqrt{N} \frac{2-q^2}{1+q} > \gamma \ln N\right) &\leq \\ \mathrm{e}^{-\sqrt{N} \frac{2-q^2}{1+q} - \gamma \ln N} \mathrm{E} \mathrm{e}^{\frac{\Theta_N}{\sqrt{N}}} &= \mathrm{e}^{-\sqrt{N} \frac{2-q^2}{1+q} - \gamma \ln N} M_{\Theta_N}\left(\frac{1}{\sqrt{N}}\right). \end{aligned}$$

From results obtained in *Step 1* and after some rearrangement, it follows that

$$\begin{aligned} M_{\Theta_N}\left(\frac{1}{\sqrt{N}}\right) &= \left(1 + O\left(\frac{1}{\sqrt{N}}\right)\right) \left(\mathrm{e}^{c_{01}\sqrt{N} + c_{02} - \frac{c_{01}^2}{2} + O(\frac{1}{\sqrt{N}})} - \right. \\ &\quad \left. (-1)^N \mathrm{e}^{c_{11}\sqrt{N} + c_{12} - \frac{c_{11}^2}{2} + N \ln q + O(\frac{1}{\sqrt{N}})}\right). \end{aligned}$$

Since  $c_{01} = \frac{2-q^2}{1+q}$  and

$$c_{11}\sqrt{N} - \frac{2-q^2}{1+q}\sqrt{N} + c_{12} - \frac{c_{11}^2}{2} + N \ln q = N \ln q \left(1 + O\left(\frac{1}{\sqrt{N}}\right)\right),$$

we obtain that

$$\mathrm{e}^{-\sqrt{N} \frac{2-q^2}{1+q} - \gamma \ln N} M_{\Theta_N}\left(\frac{1}{\sqrt{N}}\right) = \mathrm{e}^{-\gamma \ln N} O(1) \leq \frac{C_q}{N^\gamma}$$

for some constant  $C_q \in (0, \infty)$  independent of  $\gamma$ . In the same way,

$$\mathrm{P}\left(\frac{\Theta_N}{N} - \frac{2-q^2}{1+q} < -\gamma \frac{\ln N}{\sqrt{N}}\right) \leq \mathrm{e}^{\sqrt{N} \frac{2-q^2}{1+q} - \gamma \ln N} M_{\Theta_N}\left(-\frac{1}{\sqrt{N}}\right) \leq \frac{C_q}{N^\gamma},$$

provided  $C_q$  in the previous inequality was chosen large enough. Taking  $\gamma > 1$ , we then have that

$$\sum_{N=2}^{\infty} \mathrm{P}\left(\left|\frac{\Theta_N}{N} - \frac{2-q^2}{1+q}\right| > \gamma \frac{\ln N}{\sqrt{N}}\right) \leq 2C_q \sum_{N=1}^{\infty} \frac{1}{N^\gamma} < \infty.$$

Hence the claim.

Step 3: CLT. It suffices to show that

$$M_{\sqrt{N}\left(\frac{\Theta_N}{N} - \frac{2-q^2}{1+q}\right)}(t) \xrightarrow[N \rightarrow \infty]{} M_\xi(t), \quad \xi \sim N(0, \sigma^2)$$

for some fixed  $\varepsilon > 0$  and any fixed  $t \in (-\varepsilon, \varepsilon)$ . Applying expansions obtained in the *Step 1* and the reasoning similar to that of *Step 2*, we have

$$\begin{aligned} M_{\sqrt{N}\left(\frac{\Theta_N}{N} - \frac{2-q^2}{1+q}\right)}(t) &= e^{-t\sqrt{N}\frac{2-q^2}{1+q}} M_{\Theta_N}\left(\frac{t}{\sqrt{N}}\right) = \\ &e^{-tc_{01}\sqrt{N}} \left(1 + O\left(\frac{1}{\sqrt{N}}\right)\right) \left[ e^{tc_{01}\sqrt{N} + t^2\left(c_{02} - \frac{c_{01}^2}{2}\right) + O\left(\frac{1}{\sqrt{N}}\right)} - \right. \\ &\left. (-1)^N e^{tc_{11}\sqrt{N} + t^2\left(c_{12} - \frac{c_{11}^2}{2}\right) + N \ln q + O\left(\frac{1}{\sqrt{N}}\right)} \right] = \\ &\left(1 + O\left(\frac{1}{\sqrt{N}}\right)\right) e^{t^2\left(c_{02} - \frac{c_{01}^2}{2}\right) + O\left(\frac{1}{\sqrt{N}}\right)} + O\left(q^N\right) \xrightarrow[N \rightarrow \infty]{} e^{t^2\left(c_{02} - \frac{c_{01}^2}{2}\right)}. \end{aligned}$$

Direct calculations show that  $c_{02} - \frac{c_{01}^2}{2} = \frac{\sigma^2}{2}$ .

Step 4: LDP. To prove the final claim, we apply Gärtner-Ellis (GE) Theorem (see [16], Section 2.3) to  $Z_N = \frac{\Theta_N}{N}$ . First, note that, for any fixed  $\lambda \in \mathbb{R}$ ,

$$\begin{aligned} \alpha_0(\lambda) > |\alpha_1(\lambda)| \Rightarrow \lim_{N \rightarrow \infty} \frac{\kappa_{N-3}(\lambda)}{\kappa_{N-2}(\lambda)} = \frac{1}{\alpha_0(\lambda)} \Rightarrow \\ \Lambda(\lambda) := \lim_{N \rightarrow \infty} \frac{1}{N} \ln M_{Z_N}(N\lambda) &= \lim_{N \rightarrow \infty} \frac{1}{N} \ln M_{\Theta_N}(\lambda) = \\ &\lim_{N \rightarrow \infty} \frac{1}{N} \ln \alpha_0^N(\lambda) = \ln \alpha_0(\lambda) \in \mathbb{R}. \end{aligned}$$

Since  $\mathbb{R} \ni \lambda \mapsto \Lambda(\lambda)$  is differentiable at every  $\lambda \in \mathbb{R}$ , it follows that all GE assumptions hold and  $\Theta_N$  satisfies LDP with a good rate function  $I$  equal to the Legendre transform of  $\Lambda$ .  $\square$

# Bibliography

- [1] J. Abrahams. An improved lower bound on the minimum expected number of binomial group tests. *Probability in the Engineering and Informational Sciences*, 7(1):121–124, 1993.
- [2] J. Abrahams. Huffman-type codes for infinite source distributions. In *Proceedings of IEEE Data Compression Conference (DCC'94)*, pages 83–89. IEEE, 1994.
- [3] J. Abrahams. Code and parse trees for lossless source encoding. *Proceedings. Compression and Complexity of SEQUENCES 1997 (Cat. No. 97TB100171)*, pages 145–171, 1997.
- [4] M. Aldridge. Adaptive group testing as channel coding with feedback. In *2012 IEEE International Symposium on Information Theory Proceedings*, pages 1832–1836, 2012. doi: 10.1109/ISIT.2012.6283596.
- [5] M. Aldridge, O. Johnson, J. Scarlett, et al. Group testing: an information theory perspective. *Foundations and Trends® in Communications and Information Theory*, 15(3-4):196–392, 2019.
- [6] A. Allemann. *Improved upper bounds for several variants of group testing*. PhD thesis, Bibliothek der RWTH Aachen, 2003.
- [7] E. Alonderyté. Lithuanian firms to get 30 million euros to test employees, Feb. 2021. URL <https://www.lrt.lt/en/news-in-english/19/1339899/lithuanian-firms-to-get-eur30m-to-test-employees>.
- [8] H. Aprahamian, D. R. Bish, and E. K. Bish. Optimal Risk-Based Group Testing. *Management Science*, 65(9):4365–4384, Sept. 2019. ISSN 0025-1909, 1526-5501. doi: 10.1287/mnsc.2018.3138. URL <https://pubsonline.informs.org/doi/10.1287/mnsc.2018.3138>.
- [9] T. Berger, J. W. Mandell, and P. Subrahmanyam. Maximally efficient two-stage screening. *Biometrics*, 56, 2000. doi: 10.1111/j.0006-341x.2000.00833.x.
- [10] N. Bobkova, Y. Chen, and H. Eraslan. Optimal group testing with heterogeneous risks. *Economic Theory*, 77(1-2):413–444, Feb. 2024. ISSN 0938-2259, 1432-0479. doi: 10.1007/s00199-023-01502-3. URL <https://link.springer.com/10.1007/s00199-023-01502-3>.

1007/s00199-023-01502-3.

- [11] P. Chen, L. Hsu, and M. Sobel. Entropy-based Optimal Group-testing Procedures. *Probability in the Engineering and Informational Sciences*, 1(4):497–509, Oct. 1987. ISSN 0269-9648, 1469-8951. doi: 10.1017/S0269964800000541. URL [https://www.cambridge.org/core/product/identifier/S0269964800000541/type/journal\\_article](https://www.cambridge.org/core/product/identifier/S0269964800000541/type/journal_article).
- [12] X.-F. Chi, X.-Y. Lou, M. C. Yang, and Q.-Y. Shu. An optimal dna pooling strategy for progressive fine mapping. *Genetica*, 135(3): 267–281, 2009.
- [13] W. K. Chow and C. L. Chow. A discussion on implementing pooling detection tests of novel coronavirus (sars-cov-2) for a large population. *Epidemiology and Infection*, 149:e17, 2021. doi: 10.1017/S0950268820003155.
- [14] W. contributors. List of countries implementing pool testing strategy against COVID-19, Oct. 2023. URL [https://en.wikipedia.org/w/index.php?title=List\\_of\\_countries\\_implementing\\_pool\\_testing\\_strategy\\_against\\_COVID-19&oldid=1179921187](https://en.wikipedia.org/w/index.php?title=List_of_countries_implementing_pool_testing_strategy_against_COVID-19&oldid=1179921187). Page Version ID: 1179921187.
- [15] R. Cunningham, J. L. Northwood, C. D. Kelly, E. H. Boxall, and N. J. Andrews. Routine antenatal screening for hepatitis b using pooled sera: validation and review of 10 years experience. *Journal of Clinical Pathology*, 51(5):392–395, 1998. ISSN 0021-9746. doi: 10.1136/jcp.51.5.392. URL <https://jcp.bmjjournals.org/content/51/5/392>.
- [16] A. Dembo and O. Zeitouni. *Large Deviations Techniques and Applications*. Stochastic Modelling and Applied Probability. Springer Berlin Heidelberg, 2009. ISBN 9783642033117.
- [17] P. Desmos Studio. Desmos graphing calculator. <https://www.desmos.com/calculator>, 2024. Online; Accessed July 2024.
- [18] R. Dodd, E. Notari, and S. Stramer. Current prevalence and incidence of infectious disease markers and estimated window-period risk in the American Red Cross blood donor population. *Transfusion*, 42(8):975–979, Aug. 2002. ISSN 0041-1132, 1537-2995.
- [19] R. Dorfman. The detection of defective members of large populations. *The Annals of Mathematical Statistics*, 14(4):436–440, 1943. doi: 10.1214/aoms/1177731363.
- [20] D. Du and F. Hwang. *Combinatorial group testing and its applications*. Number vol. 12 in Series on applied mathematics. World Scientific, Singapore ; River Edge, NJ, 2nd ed. edition, 2000. ISBN 9789810241070.
- [21] T. A. Ebert, R. Brlansky, and M. Rogers. Reexamining the Pooled

- Sampling Approach for Estimating Prevalence of Infected Insect Vectors. *Annals of the Entomological Society of America*, 103(6):827–837, 11 2010. ISSN 0013-8746. doi: 10.1603/AN09158. URL <https://doi.org/10.1603/AN09158>.
- [22] T. S. Ferguson and C. Tatsuoka. An optimal strategy for sequential classification on partially ordered sets. *Statistics & probability letters*, 68(2):161–168, 2004.
  - [23] E. Fernandez, L. Rodrigo, S. Garcia, S. Riestra, and C. Blanco. Hepatitis b surface antigen detection using pooled sera. a cost-benefit analysis. *Revista Espanola de Enfermedades Digestivas*, 98(2):112–121, 2006.
  - [24] H. M. Finucan. The Blood Testing Problem. *Applied Statistics*, 13 (1):43, 1964. ISSN 00359254. doi: 10.2307/2985222. URL <https://www.jstor.org/stable/2985222?origin=crossref>.
  - [25] H. M. Finucan. Errata: The Blood Testing Problem. *Applied Statistics*, 14(2/3):210, 1965. ISSN 00359254. doi: 10.2307/2985344. URL <https://www.jstor.org/stable/2985344?origin=crossref>.
  - [26] S. M. G. Cormode. What’s hot and what’s not: Tracking most frequent items dynamically. *ACM Transactions on Database Systems*, 30(1):249–278, 2005. doi: 10.1145/1061318.1061325.
  - [27] M. Garey and F. Hwang. Isolating a single defective using group testing. *Journal of the American Statistical Association*, 69(345):151–153, 1974.
  - [28] M. A. Ghoneim, F. M. Attia, H. A. Kamaleldin, and M. H. Mohammad. Comparison of human immunodeficiency virus and syphilis results in pooled sera versus individual samples of blood donors attending suez canal university hospital. *The Egyptian Journal of Immunology*, 30(3):56–63, 2023.
  - [29] M. J. Golin. A combinatorial approach to golomb forests. *Theoretical Computer Science*, 263(1-2):283–304, 2001.
  - [30] M. J. Golin and K. K. Ma. Algorithms for infinite huffman-codes. In *Proceedings of the fifteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 758–767, 2004.
  - [31] M. T. Goodrich, M. J. Atallah, and R. Tamassia. Indexing Information for Data Forensics. In D. Hutchison, T. Kanade, J. Kittler, J. M. Kleinberg, F. Mattern, J. C. Mitchell, M. Naor, O. Nierstrasz, C. Pandu Rangan, B. Steffen, M. Sudan, D. Terzopoulos, D. Tygar, M. Y. Vardi, G. Weikum, J. Ioannidis, A. Keromytis, and M. Yung, editors, *Applied Cryptography and Network Security*, volume 3531, pages 206–221. Springer Berlin Heidelberg, Berlin, Heidelberg, 2005. ISBN 9783540262237 9783540315421. doi: 10.1007/11496137\_15. URL [http://link.springer.com/10.1007/11496137\\_15](http://link.springer.com/10.1007/11496137_15).

- [32] V. K. Gurani, P. Shynu, and C. L. Chowdhary. Monitoring application for dos attacks using group-testing. In *2014 International Conference on Electronics and Communication Systems (ICECS)*, pages 1–4. IEEE, 2014.
- [33] Y. Habtesllassie, L. M. Haines, H. Mwambi, and J. Odhiambo. Array-based schemes for group screening with test errors which incorporate a concentration effect. *Journal of Statistical Planning and Inference*, 167:41–57, Dec. 2015. ISSN 03783758. doi: 10.1016/j.jspi.2015.05.009. URL <https://linkinghub.elsevier.com/retrieve/pii/S0378375815001147>.
- [34] C. R. Harris et al. Array programming with NumPy. *Nature*, 585(7825):357–362, Sept. 2020. doi: 10.1038/s41586-020-2649-2. URL <https://doi.org/10.1038/s41586-020-2649-2>.
- [35] J. F. Hayes. An adaptive technique for local distribution. *IEEE Transactions on Communications*, 26(8):1178–1186, 1978. doi: 10.1109/TCOM.1978.1094204.
- [36] E. Hong and R. Ladner. Group testing for image compression. *IEEE Transactions on Image Processing*, 11(8):901–911, Aug. 2002. ISSN 1057-7149. doi: 10.1109/TIP.2002.801124. URL <http://ieeexplore.ieee.org/document/1025164/>.
- [37] L. Hsu. New procedures for group-testing based on the huffman lower bound and shannon entropy criteria. *Lecture Notes-Monograph Series*, 25:249–262, 1995. ISSN 07492170. URL <http://www.jstor.org/stable/4355848>.
- [38] M. G. Hudgens and H.-Y. Kim. Optimal Configuration of a Square Array Group Testing Algorithm. *Communications in Statistics - Theory and Methods*, 40(3):436–448, Jan. 2011. ISSN 0361-0926, 1532-415X. doi: 10.1080/03610920903391303.
- [39] J. D. Hunter. Matplotlib: A 2d graphics environment. *Computing in Science & Engineering*, 9(3):90–95, 2007. doi: 10.1109/MCSE.2007.55.
- [40] F. K. Hwang. On Finding a Single Defective in Binomial Group Testing. *Journal of the American Statistical Association*, 69(345):146–150, Mar. 1974. ISSN 0162-1459, 1537-274X. doi: 10.1080/01621459.1974.10480141. URL <http://www.tandfonline.com/doi/abs/10.1080/01621459.1974.10480141>.
- [41] F. K.-m. Hwang and D.-z. Du. *Pooling designs and nonadaptive group testing: important tools for DNA sequencing*, volume 18. World Scientific, 2006.
- [42] N. L. Johnson, S. Kotz, and X. Wu. *Inspection Errors for Attributes in Quality Control*. Springer US, Boston, MA, 1991. ISBN 9780412387708 9781489931962. doi: 10.1007/978-1-4899-3196-2.
- [43] S. M. Karimian and A. G. Straatman. Benchmarking of a 3d, unstructured, finite volume code for incompressible navier-stokes

- equation on a cluster of distributed-memory computers. In *19th International Symposium on High Performance Computing Systems and Applications (HPCS'05)*, pages 11–16. IEEE, 2005.
- [44] S. Khattab, S. Gobriel, R. Melhem, and D. Mosse. Live Baiting for Service-Level DoS Attackers. In *IEEE INFOCOM 2008 - The 27th Conference on Computer Communications*, pages 171–175, Phoenix, AZ, USA, Apr. 2008. IEEE. ISBN 9781424420261 9781424420254. doi: 10.1109/INFOCOM.2008.43. URL <http://ieeexplore.ieee.org/document/4509638/>.
  - [45] H.-Y. Kim and M. G. Hudgens. Three-Dimensional Array-Based Group Testing Algorithms. *Biometrics*, 65(3):903–910, Sept. 2009. ISSN 0006341X. doi: 10.1111/j.1541-0420.2008.01158.x. URL <http://doi.wiley.com/10.1111/j.1541-0420.2008.01158.x>.
  - [46] H.-Y. Kim, M. G. Hudgens, J. M. Dreyfuss, D. J. Westreich, and C. D. Pilcher. Comparison of Group Testing Algorithms for Case Identification in the Presence of Test Error. *Biometrics*, 63(4):1152–1163, Dec. 2007. ISSN 0006341X. doi: 10.1111/j.1541-0420.2007.00817.x. URL <http://doi.wiley.com/10.1111/j.1541-0420.2007.00817.x>.
  - [47] N. Lagopati, P. Tsoli, I. Mourkioti, A. Polyzou, A. Papaspyropoulos, A. Zafiropoulos, K. Evangelou, G. Sourvinos, and V. G. Gorgoulis. Sample pooling strategies for sars-cov-2 detection. *Journal of Virological Methods*, 289:114044, 2021. ISSN 0166-0934. doi: <https://doi.org/10.1016/j.jviromet.2020.114044>. URL <https://www.sciencedirect.com/science/article/pii/S0166093420302962>.
  - [48] M. Li and M. Xie. Nonparametric and semiparametric regression analysis of group testing samples. *International Journal of Statistics in Medical Research*, 1(1):60–72, 2012.
  - [49] Y. Liu, Y. Yin, M. P. Ward, K. Li, Y. Chen, M. Duan, P. P. Y. Wong, J. Hong, J. Huang, J. Shi, X. Zhou, X. Chen, J. Xu, R. Yuan, L. Kong, and Z. Zhang. Optimization of screening strategies for covid-19: Scoping review. *JMIR Public Health Surveill*, 10:e44349, Feb 2024. ISSN 2369-2960. doi: 10.2196/44349. URL <http://www.ncbi.nlm.nih.gov/pubmed/38412011>.
  - [50] D. S. H. M. T. Goodrich. Improved adaptive group testing algorithms with applications to multiple access channels and dead sensor diagnosis. *Journal of Combinatorial Optimization*, 15(1):95–121, 2008. doi: 10.1007/s10878-007-9087-z.
  - [51] T. Madej. An application of group testing to the file comparison problem. In [1989] Proceedings. *The 9th International Conference on Distributed Computing Systems*, pages 237–243, Newport Beach, CA,

- USA, 1989. IEEE Comput. Soc. Press. ISBN 9780818619533. doi: 10.1109/ICDCS.1989.37952. URL <http://ieeexplore.ieee.org/document/37952/>.
- [52] S. A. Mahmoud, E. Ibrahim, B. Thakre, J. G. Teddy, P. Raheja, S. Ganesan, and W. A. Zaher. Evaluation of pooling of samples for testing sars-cov-2 for mass screening of covid-19. *BMC Infectious Diseases*, 21(1):1–9, 2021.
  - [53] Y. Malinovsky. Conjectures on optimal nested generalized group testing algorithm. *Applied Stochastic Models in Business and Industry*, 36(6):1029–1036, 2020.
  - [54] Y. Malinovsky and P. S. Albert. A note on the minimax solution for the two-stage group testing problem. *The American Statistician*, 69(1):45–52, 2015. doi: 10.1080/00031305.2014.983545. URL <https://doi.org/10.1080/00031305.2014.983545>.
  - [55] Y. Malinovsky and P. S. Albert. Revisiting Nested Group Testing Procedures: New Results, Comparisons, and Robustness. *The American Statistician*, 73(2):117–125, Apr. 2019. ISSN 0003-1305, 1537-2731. doi: 10.1080/00031305.2017.1366367. URL <https://www.tandfonline.com/doi/full/10.1080/00031305.2017.1366367>.
  - [56] Y. Malinovsky and P. S. Albert. Nested group testing procedures for screening. *arXiv preprint arXiv:2102.03652*, 2021.
  - [57] Y. Malinovsky and P. S. Albert. *Nested Group Testing Procedures for Screening*, pages 1–8. John Wiley & Sons, Ltd, 2021. ISBN 9781118445112. doi: <https://doi.org/10.1002/9781118445112.stat08363>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/9781118445112.stat08363>.
  - [58] S. May, A. Gamst, R. Haubrich, C. Benson, and D. M. Smith. Pooled Nucleic Acid Testing to Identify Antiretroviral Treatment Failure During HIV Infection. *JAIDS Journal of Acquired Immune Deficiency Syndromes*, 53(2):194–201, Feb. 2010. ISSN 1525-4135. doi: 10.1097/QAI.0b013e3181ba37a7. URL <https://journals.lww.com/00126334-201002010-00006>.
  - [59] A. Meurer et al. Sympy: symbolic computing in python. *PeerJ Computer Science*, 3:e103, Jan. 2017. ISSN 2376-5992. doi: 10.7717/peerj-cs.103. URL <https://doi.org/10.7717/peerj-cs.103>.
  - [60] L. Novack, B. Sarov, R. Goldman-Levi, V. Yahalom, J. Safi, H. Soliman, M. Orgel, A. Yaari, N. Galai, J. Pliskin, and E. Shinar. Impact of pooling on accuracy of hepatitis B virus surface antigen screening of blood donations. *Transactions of The Royal Society of Tropical Medicine and Hygiene*, 102(8):787–792, 08 2008. ISSN 0035-9203. doi: 10.1016/j.trstmh.2008.04.005. URL <https://doi.org/10.1016/j.trstmh.2008.04.005>.

1016/j.trstmh.2008.04.005.

- [61] P. J. Peters, E. Westheimer, S. Cohen, L. B. Hightow-Weidman, N. Moss, B. Tsoi, L. Hall, C. Fann, D. C. Daskalakis, S. Beagle, P. Patel, A. Radix, E. Foust, R. P. Kohn, J. Marmorino, M. Pandori, J. Fu, T. Samandari, and C. L. Gay. Screening Yield of HIV Antigen/Antibody Combination and Pooled HIV RNA Testing for Acute HIV Infection in a High-Prevalence Population. *JAMA*, 315(7):682–690, 02 2016. ISSN 0098-7484. doi: 10.1001/jama.2016.0286. URL <https://doi.org/10.1001/jama.2016.0286>.
- [62] C. G. Pfeifer and P. Enis. Dorfman-Type Group Testing for a Modified Binomial Model. *Journal of the American Statistical Association*, 73(363):588–592, Sept. 1978.
- [63] R. M. Phatarfod and A. Sudbury. The use of a square array scheme in blood testing. *Statistics in Medicine*, 13, 1994. doi: 10.1002/sim.4780132205.
- [64] A. Pikovski and K. Bentele. Pooling of coronavirus tests under unknown prevalence. *Epidemiology and Infection*, 148:e183, 2020. doi: 10.1017/S0950268820001752.
- [65] W. Preiser and G. U. van Zyl. Pooled testing: A tool to increase efficiency of infant HIV diagnosis and virological monitoring. *African Journal of Laboratory Medicine*, 9:1 – 7, 00 2020. ISSN 2225-2010. URL [http://www.scielo.org.za/scielo.php?script=sci\\_arttext&pid=S2225-20102020000200005&nrm=iso](http://www.scielo.org.za/scielo.php?script=sci_arttext&pid=S2225-20102020000200005&nrm=iso).
- [66] S. M. Samuels. The Exact Solution to the Two-Stage Group-Testing Problem. *Technometrics*, 20(4):497–500, Nov. 1978. ISSN 0040-1706, 1537-2723. doi: 10.1080/00401706.1978.10489706. URL <http://www.tandfonline.com/doi/abs/10.1080/00401706.1978.10489706>.
- [67] G. Shavit, M. F. Ringenburg, J. West, R. E. Ladner, and E. A. Riskin. Group testing for video compression. In *Data Compression Conference, 2004. Proceedings. DCC 2004*, pages 212–221. IEEE, 2004.
- [68] C. H. Sherlock, S. A. Strathdee, T. Le, D. Sutherland, M. V. O’Shaughnessy, and M. T. Schechter. Use of pooling, and outpatient laboratory specimens in an anonymous seroprevalence survey of hiv infection in british columbia, canada. *Aids*, 9(8):945–950, 1995.
- [69] V. Skorniakov, R. Leipus, G. Juzeliūnas, and K. Staliūnas. Group testing: Revisiting the ideas. *Nonlinear Analysis: Modelling and Control*, 26(3):534–549, May 2021. doi: 10.15388/namc.2021.26.23933. URL <https://www.journals.vu.lt/nonlinear-analysis/article/view/23933>.
- [70] M. Sobel. Group Testing to Classify Efficiently all Defectives in a Binomial Sample. In R. E. Machol, editor, *Information and Decision Processes*, pages 127–161. New York, McGraw Hill, 1960.

- [71] M. Sobel. Optimal group testing. In *Proceedings of the Colloquium on Information Theory*, pages 411–488, Debrecen (Hungary), 1967. Organized by the Bolyai Mathematical Society.
- [72] M. Sobel and P. A. Groll. Group testing to eliminate efficiently all defectives in a binomial sample. *Bell System Technical Journal*, 38: 1179–1252, 1959.
- [73] A. Sterrett. On the Detection of Defective Members of Large Populations. *The Annals of Mathematical Statistics*, 28(4):1033–1036, Dec. 1957. ISSN 0003-4851. doi: 10.1214/aoms/1177706807. URL <http://projecteuclid.org/euclid.aoms/1177706807>.
- [74] S. H. Strogatz. *Nonlinear dynamics and chaos: with applications to physics, biology, chemistry, and engineering*. Westview Press, a member of the Perseus Books Group, Boulder, CO, second edition edition, 2015. ISBN 9780813349107. OCLC: ocn842877119.
- [75] K. Tang and J. Tang. Design of screening procedures: a review. *Journal of quality technology*, 26(3):209–226, 1994.
- [76] C. Tatsuoka and T. Ferguson. Sequential classification on partially ordered sets. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 65(1):143–157, 2003.
- [77] The pandas development team. pandas-dev/pandas: Pandas, Feb. 2020. URL <https://doi.org/10.5281/zenodo.3509134>.
- [78] E. Triesch. A group testing problem for hypergraphs of bounded rank. *Discrete Applied Mathematics*, 66(2):185–188, Apr. 1996. ISSN 0166218X. doi: 10.1016/0166-218X(95)00120-G. URL <https://linkinghub.elsevier.com/retrieve/pii/0166218X9500120G>.
- [79] P. Ungar. The cutoff point for group testing. *Communications on Pure and Applied Mathematics*, 13:49–54, 1960. doi: 10.1002/cpa.3160130105.
- [80] P. Virtanen et al. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, 17:261–272, 2020. doi: 10.1038/s41592-019-0686-2.
- [81] Wes McKinney. Data Structures for Statistical Computing in Python. In Stéfan van der Walt and Jarrod Millman, editors, *Proceedings of the 9th Python in Science Conference*, pages 56 – 61, 2010. doi: 10.25080/Majora-92bf1922-00a.
- [82] J. Wolf. Born again group testing: Multiaccess communications. *IEEE Transactions on Information Theory*, 31(2):185–191, Mar. 1985. ISSN 0018-9448. doi: 10.1109/TIT.1985.1057026. URL <http://ieeexplore.ieee.org/document/1057026/>.
- [83] M. Xie. Regression analysis of group testing samples. *Statistics in medicine*, 20(13):1957–1969, 2001.
- [84] M. Xie, K. Tatsuoka, J. Sacks, and S. S. Young. Group testing with

- blockers and synergism. *Journal of the American Statistical Association*, 96(453):92–102, 2001.
- [85] Y. Yao and F. Hwang. On optimal nested group testing algorithms. *Journal of Statistical Planning and Inference*, 24(2):167–175, Feb. 1990. ISSN 03783758. doi: 10.1016/0378-3758(90)90039-W. URL <https://linkinghub.elsevier.com/retrieve/pii/037837589090039W>.
- [86] Y. C. Yao and F. K. Hwang. A fundamental monotonicity in group testing. *SIAM Journal on Discrete Mathematics*, 1, 05 1988. doi: 10.1137/0401026.
- [87] Ying Xuan, Incheol Shin, M. T. Thai, and T. Znati. Detecting Application Denial-of-Service Attacks: A Group-Testing-Based Approach. *IEEE Transactions on Parallel and Distributed Systems*, 21(8):1203–1216, Aug. 2010. ISSN 1045-9219. doi: 10.1109/TPDS.2009.147. URL <http://ieeexplore.ieee.org/document/5232807/>.
- [88] U. Čižikovienė and V. Skorniakov. On a couple of unresolved group testing conjectures. *Communications in Statistics - Theory and Methods*, 52(8):2448–2460, 2023. doi: 10.1080/03610926.2021.1953531. URL <https://doi.org/10.1080/03610926.2021.1953531>.
- [89] U. Čižikovienė and V. Skorniakov. On the optimal configuration of a square array group testing algorithm. *Statistics and Its Interface*, 16(4):579–591, 2023. ISSN 1938-7997. doi: 10.4310/22-SII746. URL <https://www.intlpress.com/site/pub/pages/journals/items/sii/content/vols/0016/0004/a008/abstract.php>.
- [90] U. Čižikovienė and V. Skorniakov. On the generic cut-point detection procedure in the binomial group testing, 2023. URL <https://arxiv.org/abs/2304.07263>.
- [91] U. Čižikovienė and V. Skorniakov. On the optimal pairwise group testing algorithm. *Brazilian Journal of Probability and Statistics*, 38(2):253 – 265, 2024. doi: 10.1214/24-BJPS603. URL <https://doi.org/10.1214/24-BJPS603>.

# Santrauka (Summary in Lithuanian)

## Tyrimų sritis

Grupinio testavimo<sup>4</sup> (GT) modeliai taikomi tais atvejais, kai reikia klasifikuoti objektus pagal tyrėją dominantį požymį į dvi grupes — turinčius požymį ir neturinčius požymio. Nuo išprasto ištisinio testavimo, kai testuojami visi objektai, GT modeliai skiriasi tuo, kad individualių objektų testavimas atskirose viso testavimo proceso fazėse keičiamas jų grupių (iš čia pavadinimas) testavimu taip bandant suraupoti bendrą testų skaičių ir sumažinti su jais susijusias išlaidas. Tradiciškai literatūroje metodologijos pradininku laikomas R. Dorfman, straipsnyje [19] aprašęs pirmą grupinio testavimo procedūrą (toliau ją vadinsime D procedūra), skirtą tirti JAV karių, dalyvavusių antrame pasaulyiniame kare, krauko mèginius sifilio užkratui rasti. Jo idéja buvo tokia. Turint  $N$  kraujų mèginių, reikia juos sumaišyti ir ištirti gautą grupės mèginių. Jei testas užkrato nerodo, visi  $N$  tiriamujų sveiki; jei rodo — reikia ištenuoti kiek-vieną tiriamajį individualiai. Reguliariomis sąlygomis sifilis nèra dažna liga, todèl tokia metodika leidžia suraupoti nemažai testų, nes testuojant masiškai  $N$  dydžio grupėmis dominuos sveikų tiriamujų grupės, kurioms vietoje  $N$  testų pakaks vieno testo, ir tik retais atvejais pasitaikys grupės su infekuotais asmenimis (tada grupei bus sunaudotas  $N + 1$  testas).

Per 80 metų, praëjusių nuo [19] straipsnio publikavimo, GT idéja išplito į įvairias žmonijos veiklos sritis, kuriose reikia testuoti ir identifikuoti požymį turinčius objektus. Šiandien šie modeliai sutinkami medicinoje (identifikuojant pacientus, sergančius infekcinémis (ir kitomis) ligomis kaip ŽIV [58, 61, 65], hepatitas B [15, 23, 28, 60] ar COVID-

---

<sup>4</sup>kitaip — kaupinių testavimo

19 [13, 47, 49, 52, 64]), informacinių technologijų sektoriuje (aktyvių vartotojų [35, 82] ar neveikiančių jutiklių [50] komunikaciniuose tinkluose paieška; taikymai kibernetinio saugumo [31, 32, 44, 51, 87], duomenų bazių valdymo [26], kodavimo teorijos ir duomenų spūdos [36, 67] srityse), kokybės kontrolėje [42] ir kitur. Tai dinamiška, aktyviai tebevystoma ir tiriama sritis.

## Tyrimų objektas ir uždaviniai

### Savokos, žymenys ir prielaidos

Tyrėja dominantis požymis priklauso nuo taikymo srities, todėl ji turintys objektai skirtingose srityse vadinami skirtingai. Gamyboje — tai defektiniai gaminiai, medicinoje — infekuoti pacientai, komunikaciniai tinkle — pasyvūs mazgai ir t.t. Apibrėžtumo dėlei visoje tolimesnėje santraukos dalyje požymį turinčius objektus vadinsime defektiniais, o jo neturinčius — gerais, nedefektiniais arba be defekto.

$N$  raide žymėsime testuoojamas grupės dydį. Kai kuriuose GT modeiliuose jis gali būti parametrizuojamas natūraliu parametru  $n$  (konkretus pavyzdys pateikiama tolimesniame poskyryje); tada  $N$  bus  $n$  funkcija:  $\mathbb{N} \ni n \mapsto N(n) \in \mathbb{N}$ .

Nagrinėsime binarinius testavimo modelius. Juose tariama, kad grupės testo rezultatas gali būti tik dvireikšmis: jei grupėje yra bent vienas defektinis objektas, testo rezultatas lygus 1; jei visi objektai be defekto — rezultatas lygus 0. Jei grupei, kurioje yra daugiau nei vienas objektas (vieno objekto grupės irgi galimos), pritaikytas testas teigiamas, tai jis nenusako kurie objektai turi defektą.

Pagrindinės prielaidos, kuriomis remiasi dissertacijoje nagrinėti modeliai, vadinamos *binominio testavimo prielaidomis* (BTP) ir formuluojamos taip:

(BTP1) testavimo pradžioje kiekvienas objektas gali būti defektinis su ta pačia eigoje nekintarčia tikimybe  $p \in (0, 1)$ ;

(BTP2) objektai yra nepriklausomi;

(BTP3) naudojamas (fizinis) testas yra tobulas (jautrumas ir specifišumas yra lygūs 100%) ir nepriklauso nuo testuoojamos grupės dydžio.

Remdamiesi BTP galime griežtai apibrėžti disertacijoje nagrinėtą bendrą modelį.

- Tegu  $Y_1, \dots, Y_N$  yra nepriklausomi vienodai pasiskirstę Bernulio atsitiktiniai dydžiai (toliau a.d.);
- $Y_i = 1 \iff i\text{-asis objektas defektinis}$  (toki atveju  $Y_i = 0 \iff i\text{-asis objektas yra be defekto}$ );
- Bet kuriai testavimo procedūrai  $X$  ir bet kokiai netuščiai  $A \subset \{Y_1, \dots, Y_N\}$  teisinga lygybė  $X(A) = \mathbb{1}\{\sum_{Y_i \in A} Y_i > 0\}$ .

Turint grupę  $C = \{Y_1, \dots, Y_N\}$  ir testavimo procedūrą  $X$ ,  $T_X$  žymės (atsitiktinių) testų skaičių, reikalingą tam, kad identifikuoti visus defektinius aibės  $C$  elementus;  $\theta_X(N, p) = E T_X$  žymės a.d.  $T_X$  vidurkį. Atkreipsime dėmesį, kad  $\theta_X$  yra dviejų argumentų ( $N \in \mathbb{N}$  ir  $p \in (0, 1)$ ) funkcija. Visoje santraukoje  $q$  žymės dydį  $1 - p$ .

## Nagrinėti modeliai

Šiame poskyryje pateikiami konkrečių disertacijoje tirtų GT procedūrus, tenkinančių BTP, algoritmai.

**Dorfman procedūra D** Šią procedūrą jau aprašėme įvade. Ją sudaro du žingsniai.

Žingsnis nr. 1: testuojame visos grupės jungtinį mèginį  $JM$ ;

Žingsnis nr. 2: jei  $JM$  testas neigiamas, baigiamo; priešingu atveju pakartotinai testuojame kiekvieną objektą individualiai.

Iš procedūros aprašymo išplaukia lygybė

$$T_D = 1 + N \mathbb{1}\{Y_1 + \dots + Y_N > 0\};$$

todėl,

$$\theta_D(N, p) = 1 + N E \mathbb{1}\{Y_1 + \dots + Y_N > 0\} = 1 + N(1 - q^N). \quad (\text{S.1})$$

**Modifikuota Dorfman procedūra MD** Sobel ir Groll darbe [72] pasrebėjo, kad D procedūroje galimas perteklinis testas ir pasiūlė tokią pataisą: jeigu pradinis jungtinio mèginio testas buvo teigiamas ir pakartotinai ištestavus  $N - 1$  objektų defektinis vis dar nebuvo aptiktas, tai paskutinio objekto testuoti neberekia (ir taip aišku, kad jis defektinis). Remiantis procedūros aprašymu

$$T_{MD} = 1 + (N - 1)\mathbb{1}\{Y_1 + \dots + Y_N > 0\} + \mathbb{1}\{Y_1 + \dots + Y_{N-1} > 0\}.$$

Taigi,

$$\theta_{MD}(N, p) = 1 + (N - 1)(1 - q^N) + 1 - q^{N-1} = 1 - pq^{N-1} + N(1 - q^N). \quad (\text{S.2})$$

**Sterrett procedūra ST** Sterrett [73] pasiūlė kitą D procedūros modifikaciją, nusakomą tokiu algoritmu.

Žingsnis nr. 1: testuojame visos grupės jungtinį mèginį  $JM$ ;

Žingsnis nr. 2: jei  $JM$  testas neigiamas, baigiamo; priešingu atveju vykdome žingsnį nr. 3;

Žingsnis nr. 3: pakartotinai testuojame po vieną objektą tol, kol aptinkame pirmą defektinį; jei ištestuota visa grupė, baigiamo; priešingu atveju likusią netestuotą grupės dalį laikome nauja pradine grupe ir jai taikome algoritmą pradėdami nuo žingsnio nr. 1.

Darbe [72] buvo parodyta, kad

$$\theta_{ST}(N, p) = 2q - p^{-1}(1 - q^{N+1}) + (2 - q)N. \quad (\text{S.3})$$

**Porinė testavimo procedūra PT** Ši procedūra nagrinėta Yao ir Hwang darbe [85]. Ji nusakoma žemiau pateikiamu algoritmu, kurio kiekvienam žingsnyje aibė  $C$  žymi testuojamų objektų aibę.

Žingsnis nr. 1: jei aibėje  $C$  yra vienintelis elementas, atliekame jo testavimą ir baigiamo; jei elementų išvis nėra, baigiamo; priešingu atveju vykdome žingsnį nr. 2;

Žingsnis nr. 2: iš aibės  $C$ , turinčios  $N \geq 2$  objektų, parenkame du; suformuojame jungtinį mèginį ir ištestave ji vykdome žingsnį nr. 3;

### Žingsnis nr. 3:

- jei jungtinio mèginio testas neigiamas, priskiriame abu objektus grupei be defekto,  $N = N - 2$ ,  $C = C \setminus \{\text{testuota pora}\}$ ;
- jei jungtinio mèginio testas teigiamas, parenkame bet kuri iš dviejų objektų ir pakartotinai ištetsuojame; jei jo testas neigiamas, priskiriame netestuotą objektą grupei su defektu,  $N = N - 2$ ,  $C = C \setminus \{\text{testuota pora}\}$ ; jei jo testas teigiamas,  $N = N - 1$ ,  $C = C \setminus \{\text{pakartotinai testuotas objektas}\}$ .

### Žingsnis nr. 4: pradëti viską iš naujo nuo žingsnio nr. 1.

Darbe [85] išvesta ši a.d.  $T_{PT}$  vidurkio formulė:

$$\theta_{PT}(N, p) = N \frac{2 - q^2}{1 + q} + \frac{q^2 + q - 1}{(1 + q)^2} (1 - (-q)^N). \quad (\text{S.4})$$

**Kvadratinës matricos procedûra A2** Procedûra<sup>5</sup> pasiûlyta Phatarfod and Sudbury [63]; vëliau apibendrinta Berger, Mandell ir Subrahmany [9]. Norint taikyti procedûrą, reikia, kad bendras mèginių skaičius  $N$  bûtų pavidalo  $N = n^2, n \in \mathbb{N}$ . Galiojant šiai sąlygai, procedûra nusakoma žemiau pateikiamu algoritmu.

### Žingsnis nr. 1: išdëstome turimus mèginius ant $n \times n$ kvadratinës matricos;

Žingsnis nr. 2: suformuojame  $n$  junginių mèginių, atitinkančių eilutes, ir  $n$  junginių mèginių, atitinkančių stulpelius; ištetsuojame šiuos  $2n$  mèginius;

Žingsnis nr. 3: jei visi testai neigiami, baigiamo; priešingu atveju pakartotinai testuojame objektus  $I_{ij}$ , tenkinančius sąlygą „eilutes  $i$  stulpelio  $j$  testai teigiami“.

Phatarfod ir Sudbury [63] apskaičiavo vidutinį A2 testų skaičių:

$$\begin{aligned} \theta_{A2}(N, p) &= 2n + n^2 (1 - 2q^n + q^{2n-1}) = \\ &= 2\sqrt{N} + N \left(1 - 2q^{\sqrt{N}} + q^{2\sqrt{N}-1}\right). \quad (\text{S.5}) \end{aligned}$$

---

<sup>5</sup>žymuo A2 nuo angl. square array

## Disertacijoje spręsti Grupinio testavimo uždaviniai

Prieš įvardindami disertacijoje spręstus uždavinius, aptarsime tipinius GT uždavinius.

### Tipiniai Grupinio testavimo uždaviniai

Tegu  $X$  — fiksuota GT procedūra. Kadangi pagrindinis GT tikslas yra minimizuoti  $\theta_X(N, p)$  ir paprastai tikimybė  $p$  galima (bent jau trumpuoju laikotarpiu) laikyti nekintančia, vienas tipinių GT teorijos uždavinių — rasti optimalų testuojamos grupės dydį  $N$ , kai tikimybė  $p$  laikoma nekintančia. Formaliai problema nusakoma pasitelkiant funkciją

$$\mathbb{N} \times (0, 1) \ni (N, p) \mapsto t_X(N, p) := \frac{\theta_X(N, p)}{N}, \quad (\text{S.6})$$

žyminčią vidutinį testų skaičių, tenkantį vienam objektui tuo atveju, kai testuojamos grupės dydis yra  $N$ . Bet kuris globalus šios funkcijos minimumo taškas  $N \in \arg \min_{N \in \mathbb{N}} t_X(N, p)$  vadinamas optimalia konfigūracija ir žymimas  $N_{opt}^X$ . Kadangi  $p$  yra fiksuota,  $N_{opt}^X = N_{opt}^X(p)$  priklauso nuo  $p$ . Iš to išplaukia keli svarbūs pastebėjimai.

Funkcija  $(0, 1) \ni p \mapsto N_{opt}^X(p)$  rodo kaip kinta optimalios grupės dydis, kai testuojame maksimizuodami vidutinį išlošį

$$G_X(p) := 1 - t_X(N_{opt}^X(p), p) \quad (\text{S.7})$$

ilgoje testavimo (grupėmis) serijoje. Intuityviai aišku, kad bet kuriai tipinei testavimo procedūrai  $X$  funkcija  $(0, 1) \ni p \mapsto N_{opt}^X(p)$  bus nedidėjanti,  $G_X(p) \xrightarrow[p \rightarrow 0+]{} 1-$ , o  $N_{opt}^X(p) \xrightarrow[p \rightarrow 0+]{} \infty$ . Kai kuriuose praktiniuose taikymuose egzistuoja natūralūs apribojimai testuojamų grupių dydžiams, kuriuos peržengus testo charakteristikos tampa nebepriimtinės — drastiškai sumažėja jautrumas ir/arba specifišumas. Žinant maksimalią slenkstinę vertę  $N_{max}$ , kurią peržengus stebimas minėtas efektas, galima rasti sritį  $R_X = \{p \in (0, 1) : N_{opt}^X(p) \leq N_{max}\}$ ; tam reikia išreikštinio funkcijos  $(0, 1) \ni p \mapsto N_{opt}^X(p)$  pavidalo. Išreikštinis šios funkcijos pavidas leidžia palyginti kelias procedūras ir pasirinkti labiausiai tinkančią turimai problemai: jei  $X_i, i = 1, \dots, k$ , yra problemai tinkančios procedūros, galima apskaičiuoti išlošius  $G_{X_i}(p)$  ir grupių dydžius  $N_{opt}^{X_i}(p)$  skirtingoms  $p$  reikšmėms bei kiekvienai  $p$  reikšmei atsirinkti labiausiai tinkančią procedūrą.

Su funkcija  $N_{opt}^X$  tam priai susijęs dar vienas tipinis GT uždavinys — optimalaus tikimybinio slenksčio (OTS)  $p_c^X \in (0, 1)$  radimas. Šio uždavinio apibréžimą iš pradžių paaškinsime neformaliai. Kaip minėta įvade D procedūros atveju ir kaip sufleruoja intuicija, GT metodika turėtų pasiteisinti tik tada, kai defekto tikimybė  $p$  pakankamai maža. Pvz., prisiminus Dorfman procedūrą aišku, kad didelėms  $p$  reikšmėms grupėje iš  $N$  objektų dažnai pasitaikys bent vienas defektinis ir tokiais atvejais, užuot naudojė  $N$  testų testuodami individualiai, taikydami D procedūrą naudosime  $N + 1$  testą; todėl OTS radimo uždavinys formuluojamas taip: fiksavus procedūrą  $X$  reikia rasti tokią  $p_c^X \in (0, 1)$  reikšmę, kad  $\forall p \in (p_c^X, 1) N_{opt}^X(p) = 1$ . Bendrame BTP kontekste šią problemą sprendė Ungar. Darbe [79] jis gavo ši fundamentalų GT rezultata.

**Teorema S.0.1.** Binominių GT procedūrų taikymas turi prasmę tada ir tik tada, kai  $p \in \left(0, \frac{3-\sqrt{5}}{2}\right)$ : jei  $p \notin \left(\frac{3-\sqrt{5}}{2}, 1\right)$ , tai neegzistuoja GT procedūros, kuri vidutiniškai naudotų mažiau nei vieną testą objektui (t.y.,  $t_X(N_{opt}^X(p), p) \geq 1$  su bet kokia procedūra  $X$ ); jei  $p \in \left(0, \frac{3-\sqrt{5}}{2}\right)$ , tai atsiras bent viena tokia procedūra  $X$ , kad  $t_X(N_{opt}^X(p), p) < 1$ .

Binominio GT literatūroje skaičius  $\frac{3-\sqrt{5}}{2}$  dažnai vadinamas universaliu tikimybiniu slenksčiu (toliau UCP nuo angl Universal Cut-Point). Kai kurioms GT procedūroms jis sutampa su  $p_c^X$ , t.y.

$$G_X(p) = 0 \iff p \geq \frac{3 - \sqrt{5}}{2}.$$

Kita vertus, egzistuoja tokios procedūros, kurioms

$$G_X(p) = 0 \iff p \geq p_c^X \text{ su } p_c^X < \frac{3 - \sqrt{5}}{2}.$$

Taigi, tiriant konkrečią procedūrą  $X$ , visų pirma svarbu surasti  $p_c^X$ . Atsižvelgiant į teoremą S.0.1, visoje disertacijoje (ir santraukoje), funkcijos, priklausančios nuo  $p$ , nagrinėjamos tik intervale  $\left(0, \frac{3-\sqrt{5}}{2}\right)$ .

## Disertacijoje spręsti uždaviniai

- Trims GT procedūroms — MD, ST ir A2 — spręstas išreikštinio  $N_{opt}^X$  pavidalo radimo uždavinys.

- Specifiniam procedūrų poklasiui, tenkinančiam BTP, konstruotas  $p_c^X$  radimo algoritmas.
- Ieškotas būdas, leidžiantis charakterizuoti a.d.  $T_{PT}$  skirstinį taip išvertinant PT procedūros tikimybines savybes.

## Disertacijos struktūra

Disertaciją sudaro penki skyriai. Pirmasis skirtas įvadui. Jame aptariama GT idėja, apžvelgiama raidos istorija, pateikiami taikymų pavyzdžiai iš įvairių veiklos sričių. Antrame skyriuje apibrėžiamos disertacijoje vartojamos sąvokos ir žymėjimai, detaliai aprašomas disertacijos tyrimo objektas. Trečiajame skyriuje suformuluoti teoriniai tyrimų rezultatai, pateikiami taikymų pavyzdžiai. Ketvirtasis skyrius skirtas rezultatų aptarimui ir baigiamosioms pastaboms. Paskutiniame skyriuje surinkti teorinių rezultatų įrodymai.

## Rezultatų apžvalga

Kiekvienam disertacijoje spręstam uždaviniui skiriamas atskiras poskyris, kuriamo suformuluoti gauti teoriniai rezultatai ir aptariama susijusi literatūra.

### Optimalios Modifikuotos Dorfman, Sterrett ir Kvadratinės matricos procedūrų konfigūracijos

Mūsų žiniomis, darbų, kuriuose ieškota analizinė  $N_{opt}^X$  išraiška, nėra daug. D procedūra nagrinėta Samuels straipsnyje [66]. Jame parodyta, kad<sup>6</sup>  $N_{opt}^D(p) \in \{\lfloor \sqrt{p^{-1}} \rfloor + 1, \lfloor \sqrt{p^{-1}} \rfloor + 2\}$ , kai  $p \in (0, 1 - (1/3)^{1/3}) \approx (0, 0.31)$  ir  $N_{opt}^D(p) = 1$ , kai<sup>7</sup>  $p \geq 1 - (1/3)^{1/3}$ . Modifikuotos Dorfman procedūros MD, Sterrett procedūros ST ir kvadratinės matricos procedūros A2 atvejais tikslios analizinės išraiškos  $N_{opt}^X$  nebuvvo žinomos ilgą laiką. Tiksliau tariant, MD ir ST atveju, Malinovsky ir Albert [55]

<sup>6</sup>čia ir toliau  $\lfloor x \rfloor$  žymi sveikają skaičiaus  $x \in \mathbb{R}$  dalį;  $\lceil x \rceil = x$ , kai  $x \in \mathbb{Z}$  ir  $\lceil x \rceil = \lfloor x \rfloor + 1$ , kai  $x \in \mathbb{R} \setminus \mathbb{Z}$

<sup>7</sup>pastarasis rezultatas taip pat rodo, kad procedūros D atveju  $p_c^D = 1 - (1/3)^{1/3}$  yra griežtai mažesnis už universalų slenkštį  $UCP = \frac{3-\sqrt{5}}{2}$

skaitiškai nuspėjo šias išraiškas gana plačiame  $p$  reikšmių intervale ir iškélė hipotezę, kad jos turėtų būti tokios pat visoms  $p \in (0, UCP)$  reikšmėms. Mūsų darbe [88] mes patvirtinome jų hipotezes įrodydami žemiau pateiktiamas teoremas.

**Teorema S.0.2.** *Tegu*

$$g_0(p) := \frac{1}{q} \left( \frac{1 - 2pq}{q \left( 1 - \ln q \sqrt{\frac{2}{p}} \right)} \right)^{\sqrt{\frac{p}{2}}} \quad \text{su } p \in \left( 0, \frac{3 - \sqrt{5}}{2} \right). \quad (\text{S.8})$$

Aibė  $g_0^{-1}(\{1\})$  sudaryta iš vieno taško  $p_* \approx 0.1711$ .

$$\begin{aligned} N_{opt}^{ST}(p) &\in \left\{ \lfloor \sqrt{2p^{-1}} \rfloor, \lfloor \sqrt{2p^{-1}} \rfloor + 1 \right\}, \text{ kai } p \in \left( p_*, \frac{3 - \sqrt{5}}{2} \right); \\ N_{opt}^{ST}(p) &\in \left\{ \lfloor \sqrt{2p^{-1}} \rfloor, \lfloor \sqrt{2p^{-1}} \rfloor + 1, \lfloor \sqrt{2p^{-1}} \rfloor + 2 \right\}, \text{ kai } p \in (0, p_*]. \end{aligned}$$

**Teorema S.0.3.** *Su visomis  $p \in \left( 0, \frac{3 - \sqrt{5}}{2} \right)$  reikšmėmis*

$$N_{opt}^{MD}(p) \in \left\{ \lfloor \sqrt{p^{-1}} \rfloor, \lfloor \sqrt{p^{-1}} \rfloor + 1 \right\}. \quad (\text{S.9})$$

A2 procedūra tirta Hudgens ir Kim darbe [38]. Autoriams pavyko gauti gana tikslius apatinius ir viršutinius  $N_{opt}^{A2}(p)$  režius, tačiau analizinė šios funkcijos išraiška liko nežinoma. Jie taip pat nagrinėjo OTS problemą ir surado  $p_c^{A2}$  reikšmę. Darbe [89] mums pavyko išvesti analizinę  $N_{opt}^{A2}(p)$  išraišką ir pateikti papildomų ižvalgų apie  $p_c^{A2}$  (žr. pastabą S.0.6 apačioje). Prieš formuluodami gautus rezultatus, priminsime, kad A2 atveju testuojamos grupės dydis  $N$  parametrizuojamas natūraliu parametru  $n$ :  $N(n) = n^2$ . Be to, A2 atveju mums buvo patogiau traktuoti  $t_{A2}$  kaip funkciją  $[2, \infty) \times (UCP, 1) \ni (n, q) \mapsto t_{A2}(n, q)$  su tolydžiu argumentu  $n \in [2, \infty)$  užuot laikius ją argumentų  $(N, p)$ , kintančių aibėje  $\{n^2 : n \in \mathbb{N}\} \times (0, UCP)$ , funkcija. Žemiau pateikiamuose teiginiuose šis susitarimas galioja.

**Teorema S.0.4.** *Tegu  $g(q, n) = \frac{2}{n} - 2q^n + q^{2n-1} = t_{A2}(n, q) - 1$ .*

(i) *Srityje  $(q, n) \in (1/2, 1) \times (2, \infty)$  egzistuoja vienintelis sistemos*

$$\begin{cases} 1 = nq^n \left( 1 - \frac{q^{n-1}}{2} \right) \\ n \ln q = -\frac{\left( 1 - \frac{q^{n-1}}{2} \right)}{(1 - q^{n-1})} \end{cases}$$

sprendinys  $(q_*, n_*) \approx (0.748416, 4.453524)$ .

- (ii) *n atžvilgiu bet kokiai fiksuotai  $q \in (q_*, 1)$  reikšmei lygtis  $g(q, n) = 0$  turi du sprendinius  $n_L, n_U : 2 < n_L < n_* < n_U < \infty$ . Intervale  $(n_L, n_U)$  funkcija  $n \mapsto g(q, n)$  įgyja neigiamas reikšmes, o aibėje  $(2, \infty) \setminus [n_L, n_U]$  — teigiamas.*
- (iii) *Fiksuotai  $q \in (q_*, 1)$  reikšmei A2 procedūra efektyvi intervale  $(n_L, n_U)$ , t.y.  $t_{A2}(q, n) < 1$  su visais  $n \in (n_L, n_U)$ . Šiame intervale egzistuoja vienintelis (ir todėl globalus) funkcijos  $(2, \infty)$   $\ni n \mapsto t_{A2}(q, n)$  minimumo taškas  $n_{min}$ . Jei  $q \in [0.755, 1)$ , tai*
- $$n_{min} = \frac{1}{p^{\frac{2}{3}}} + \frac{1}{2p^{\frac{1}{3}}} + 0.2 + 3p^2 + t_*$$
- su  $t_* \in [0, 1]$ .*
- (iv) *Fiksuotai  $q \in (q_*, 1)$  reikšmei intervale  $(n_U, \infty)$  funkcija  $(2, \infty) \ni n \mapsto t_{A2}(n, q)$  taip pat turi vienintelį (ir todėl globalų) maksimumo tašką. Bet kokiai  $q \in (0, q_*)$  reikšmei A2 nėra optimali: funkcija  $(2, \infty) \ni n \mapsto t_{A2}(n, q)$  įgyja reikšmes intervale  $(1, \infty)$ .*

**Išvada S.0.5.** Tegu  $g(q, n)$  tokia pati kaip teoremoje S.0.4. Lygtis  $g(q, 5) = 0$  turi vienintelį sprendinį  $q_5 \approx 0.750209961$ . Visoms  $q \in (q_5, 1)$  reikšmėms  $n_{opt}(q)$  priklauso aibei

$$\left\{ \left\lfloor \frac{1}{p^{\frac{2}{3}}} + \frac{1}{2p^{\frac{1}{3}}} + 3p^2 + 0.2 \right\rfloor + i : i = 0, 1, 2 \right\}.$$

**Pastaba S.0.6.** Gali susidaryti įspūdis, kad pateikti rezultatai nėra išsamūs (neaišku kokia yra  $n_{opt}(q)$  reikšmė, kai  $q \in (q_*, q_5)$ ). Taip nėra: remiantis teorema S.0.4, bet kokiai  $q \in (q_*, q_5)$  reikšmei  $t_{A2}(n_{min}, q) < 1$  su  $n_{min} = n_{min}(q)$  iš (iii) dalies; iš disertacijoje pateikiamo įrodymo išplaukia, kad  $n_{min}(q) \in (4, 5)$  ir  $\min(t_{A2}(4, q), t_{A2}(5, q)) > 1$ , kai  $q \in (q_*, q_5)$ . Hudgens ir Kim [38] analizavo A2 diskrečioje skalėje, t.y. tardami, kad  $n \in \mathbb{N}$ . Jie įrodė, kad:

- imant  $n = 2, 3, 4$  ir bet kokią  $q$  reikšmę procedūra A2 nėra optimali;
- $\forall q \in (q_5, 1)$  funkcija  $\{5, 6, \dots\} \ni n \mapsto t_{A2}(n, q)$  įgyja reikšmes mažesnes už 1.

Faktiškai jų rezultatai reiškia, kad  $(q_5, 1)$  — tai tas intervalas, kuriame A2 efektyvi arba kitaip — kad  $p_c^{A2} = 1 - q_5$ . Kaip matome, mūsų analizė pateikia papildomą ižvalgą apie  $p_c^{A2}$  ir procedūros A2 elgesį tolydžioje skalėje. Ižvalgos įdomios iš teorinės pusės, tačiau pridėtinės praktinės vertės nesuteikia.  $\square$

## Optimalaus tikimybinio slenksčio paieškos algoritmas

OTS paieška  $p_c^X$  — svarbi problema: konkrečioje situacijoje žinodami, kad  $p > p_c^X$ , iškart žinome, kad procedūros  $X$  taikymas prasmės neturi. Vis dėlto, mūsų žiniomis, teorema S.0.1 yra vienintelis bendro pobūdžio rezultatas, skirtas šiai problemai. Kitų autorių darbuose  $p_c^X$  radimas siejamas su konkretios procedūros analize. Darbe [90] mes pasiūlėme algoritma, leidžiantį surasti apytikrę OTS reikšmę ir dažnais atvejais rekonstruoti tikslią reikšmę  $p_c^X$  gana plačiai BTP tenkinančių procedūrų klasei. Analizės metu mes taip pat atskleidėme įdomų ryšį tarp GT ir dinaminiių sistemų bifurkacijų teorijos. Pasiūlytas algoritmas tinkamai testavimo procedūroms, tenkinančiomis šiuos apribojimus.

(M0)  $\exists c \geq 2$  tokis, kad a-priori žinoma, jog procedūra  $X$  neefektyvi, kai  $N \in [1, c]$ .

(M1) Funkcija  $\mathbb{N} \times (0, UCP] \ni (N, p) \mapsto \theta_X(N, p)$  gali būti traktuojama kaip tolydžiai kintančių argumentų funkcija aibėje  $[c, \infty) \times (0, UCP]$  ir yra diferencijuojama šios aibės viduje.

(M2)  $\forall N \in (c, \infty)$  funkcija  $(0, UCP] \ni p \mapsto \theta_X(N, p)$  yra griežtai didėjanti.

(M3)  $\forall N \in (c, \infty) t_X(N, UCP) > 1$ .

(M4)  $\forall N \in (c, \infty) \exists p \in (0, UCP) : t_X(N, p) < 1$ .

Mūsų rezultatai pateikiami dviejuose žemiau suformuluotuose teiginiuose. Pirmasis charakterizuojas OTS savybes.

**Teiginys S.0.1.** Tarkime, kad tenkinamos prielaidos (M0)–(M4). Tegu  $p_{X,c} = \sup\{p \in (0, UCP) \mid \exists N \in (c, \infty) : t_X(N, p) < 1\}$ . Tada  $\forall p \in (0, p_{X,c})$  procedūra  $X$  efektyvi tolydžioje skalėje;  $\forall p \in (p_{X,c}, UCP]$  ji neefektyvi jokia prasme, t.y.

$$(N, p) \in (c, \infty) \times (p_{X,c}, UCP] \Rightarrow t_X(N, p) > 1.$$

Antrasis teiginys rodo, kad galiojant (M0)–(M4) egzistuoja  $p_{X,c}$  radimo algoritmas, kuris gali būti aprašomas bifurkacijų teorijos terminais. Iš tikrujų, traktuodami  $p \in (0, UCP]$  kaip kontroliuojamą parametra, o  $N \in (c, \infty)$  — kaip latentinio tolydaus kintamojo funkcija, įveskime dinaminę sistemą, apibrėžiamą lygtimi

$$\dot{N} = t_X(N, p) - 1. \quad (\text{S.10})$$

**Teiginys S.0.2.** Tarkime, kad tenkinamos prielaidos (M0)–(M4); tada  $p_{X,c}$  yra sistemos (S.10) bifurkacijos taškas ir galima išskirti tris žemiau aprašytus bifurkacijų tipus.

(b0)  $p_{X,c}$  yra vienintelė kontroliuojamo parametro reikšmė su kuria (S.10) turi fiksotų taškų intervale  $(c, \infty)$ . Šiuo atveju  $p_{X,c} < UCP$  ir bet kuris  $N \in (c, \infty)$  yra lygties  $t_X(N, p_{X,c}) = 1$  sprendinys.

Jei egzistuoja tokia reikšmė  $p_l \in (0, p_{X,c})$ , kad (S.10) turi fiksotą tašką  $N \in (c, \infty)$ , tai galimi du atvejai:

(b1) (S.10) turi fiksotų taškų intervale  $(c, \infty)$  su visomis  $p \in [p_l, p_{X,c})$  reikšmėmis, tačiau nėra fiksotų taškų, atitinkančių reikšmę  $p_{X,c}$ ;

(b2) (S.10) turi fiksotų taškų intervale  $(c, \infty)$  su visomis  $p \in [p_l, p_{X,c}]$  reikšmėmis išskaitant  $p_{X,c}$  reikšmę, kuri šiuo atveju yra griežtai mažesnė už  $UCP$ .

Visais atvejais bifurkacijos kreivė indukuoja diferencijuojamą atvaizdį  $(c, \infty) \ni N \mapsto p_N \in (0, p_{X,c}]$  ir  $p_{X,c}$  lygus jo maksimumui. Turint (b0) ir (b2) tipų bifurkacijas, šis taškas gali būti surastas sprendžiant dviejų kintamųjų sistemą

$$\begin{cases} t_X(N, p) = 1, \\ \frac{\partial}{\partial N} t_X(N, p) = 0 \end{cases} \quad (\text{S.11})$$

(t.y.  $N$  ir  $p$  atžvilgiu) ir parenkant maksimalią  $p$  reikšmę iš aibės  $S = \{(N, p) \in (c, \infty) \times (0, UCP] \mid (N, p) \text{ yra sistemos (S.10) sprendinys}\}$ . Turint (b1) tipo bifurkaciją  $p_{X,c} = \max(\lim_{N \rightarrow c+} p_N, \lim_{N \rightarrow \infty} p_N)$ ; pastaroji lygybė galioja ir tuo atveju, kai sistema (S.11) neturi sprendinių aibėje  $(c, \infty) \times (0, UCP]$ .

Pateiksime kelias aiškinamąjas pastabas.

**Pastaba S.0.7.** Remiantis teiginiu S.0.2 galima rasti OTS tolydžioje skalėje (TOTS). Praktikoje testuojamos grupės dydis  $N$  yra sveikaskaitinis; todėl dirbama diskrečioje skalėje. Paprastai diskretus OTS (toliau DOTS), iki šiol žymėtas ir toliau žymimas tuo pačiu simboliu  $p_c^X$ , yra mažesnis už  $TOTS = p_{X,c}$  iš teiginio S.0.1, tačiau skirtumai dažniausiai nežymūs (pavyzdžiai pateikiami disertacijoje). Kadangi  $p_c^X$  radimas dažnai sudėtingas,  $p_{X,c}$  yra nebloga aproksimacija. Be to, kaip iliustruoja disertacijoje pateikiami pavyzdžiai, (b2) atveju DOTS neretai galima surasti naudojantis tokiu algoritmu:

- parenkame tokį  $N_c$ , kad pora  $(N_c, p_{X,c})$  yra (S.11) sistemos sprendinys;
- imame  $DOTS = \max(p_{\lfloor N_c \rfloor}, p_{\lceil N_c \rceil})$ .

(b1) atveju DOTS dažnai sutaps su TOTS.  $\square$

**Pastaba S.0.8.** Procedūros A2 pavyzdys rodo, kad kai kuriais atvejais grupės dydis  $N = N(n)$  yra argumento  $n \in \mathbb{N}$  funkcija ir patogiau traktuoti  $\theta_X, t_X$  bei kitas susijusias funkcijas kaip argumento  $(n, p)$  (o ne  $(N, p)$ ) funkcijas. Keičiant  $N$  į  $n$  salygose (M0)–(M4) ir visose susijusiose funkcijose, nekeičia tvirtinimų S.0.2–S.0.1 išvadą jei  $(c, \infty) \ni n \mapsto N(n)$  yra diferencijuojama ir griežtai didėjanti.  $\square$

**Pastaba S.0.9.** Mūsų manymu, bifurkacijos (b1)–(b2) yra dominuojančios — mes nesugebėjome rasti pavyzdžio, atitinkančio tipą (b0). Kita vertus, mes nesugebėjome eliminuoti šio atvejo teoriškai.  $\square$

Baigdami trumpai aptarsime salygų (M0)–(M1) prasmę ir stiprumą.

(M0) galima traktuoti kaip salyga, reikalingą apriboti dinaminės sistemos (S.10) kitimo sričią. Jos neapribojus, bifurkacijos kreivės forma galėtų smarkiai pasikeisti (atitinkamas pavyzdys pateikiamas disertacijoje). Kartu su salyga (M3), ji apibrėžia kraštinę reikšmę  $p_{X,c}$  radimui tais atvejais, kai  $p_{X,c} = UCP$ . Tipinėje situacijoje galima imti  $c = 2$ ; tada (M0) bus tenkinama, o lygypė  $c = 2$  jokio esminio apribojimo nededa, nes nagrinėjamame kontekste ji tereiškia, kad testuojant aibę iš vieno elemento jokia GT nereikalinga — vienam elementui visada reikia vieno testo.

Salyga (M1) yra stipriausia. Ją tenkina toli gražu ne visos procedūros. Pvz., procedūrai PT ji negalioja — iš lygties (S.4) matome, kad  $\theta_{PT}(N, p)$

neišeina pratęsti iki funkcijos aibėje  $[c, \infty) \times (0, UCP)$ , išyjančios reikšmes intervale  $[0, \infty)$ . Iš esmės ši sąlyga ir apibrėžia klasę, kuriai mūsų pasiūlytas metodas tinkta, nes kita sąlygas galima traktuoti kaip „natūralias“ ir tinkančias daugeliui GT procedūrų, tenkinančių BTP.

(M2) reiškia, kad didėjant defekto tikimybei vidutinis testų skaičius, tenkantis  $N$  dydžio grupei, irgi turėtų augti. Salygos pagrįstumui galiama paminėti fundamentalų Yao ir Hwang darbą [86], kuriame parodyta, kad  $\forall N \in \mathbb{N}$  funkcija  $(0, UCP] \ni p \mapsto \inf_X \theta_X(N, p)$ , kurioje infimumas imamas pagal visas galimas BTP tenkinančias procedūras, yra griežtai didėjanti.

Gali pasiroyti, kad (M3) eliminuoja visas procedūras, kurioms  $p_{X,c} = UCP$ . Disertacijoje pateikiami pavyzdžiai rodo, kad taip nėra. Mūsų manymu, (M3) eliminuoja optimalias procedūras<sup>8</sup>. Tai nėra trūkumas, nes optimalioms procedūroms lygybė  $p_{X,c} = UCP$  labai tikėtina.

Galiausiai (M4) techniškai išreiškia faktą, kad mes nagrinėjame tik tas procedūras, kurios savo efektyvioje srityje turi prasmę bet kokiam testuojamų objektų skaičiui bent jau atskiroms defektyvumo tikimybių reikšmėms. Ši sąlyga galioja daugeliui procedūrų, nes, kaip minėta anksčiau, dažniausiai  $N_{opt}^X(p) \xrightarrow[p \rightarrow 0+]{} \infty$ .

## Porinės testavimo procedūros analizė

Tarp disertacijoje tirtų procedūrų PT užima išskirtinę vietą: darbe [85] buvo parodyta, kad intervale  $p \in \left[1 - \frac{1}{\sqrt{2}}, UCP\right]$  PT procedūra yra globaliai optimali įdėtuju procedūru (angl. nested procedures) klasėje<sup>9</sup>. Kitaip tariant, imant bet kokią įdėtają procedūrą  $X$  ir bet kokį  $N \in \mathbb{N}$ ,  $\theta_{PT}(N, p) \leq \theta_X(N, p)$  tolygiai  $p \in \left[1 - \frac{1}{\sqrt{2}}, UCP\right]$  atžvilgiu. Nepaisant to, PT procedūra nebuvo plačiai nagrinėjama GT literatūroje. Nusprendę ją tirti ir peržiūrėjė [85] citavusius straipsnius<sup>10</sup> e-platformoje Google Scholar, aptikome, kad iš penkiolikos surastų straipsnių [1–3, 12, 22, 29,

<sup>8</sup>t.y. tokias, kurioms vidurkis  $\theta_X(N, p)$  minimalus kažkuriame  $(a, b) \subset (0, UCP)$

<sup>9</sup>Procedūra vadinama įdėta, jei testai atliekami nuosekliai vienas po kito ir ji tenkina šiuos reikalavimus: 1) kiekviename testavimo žingsnyje galima naudotis visa ankstesne testavimo metu sukaupta informacija; 2) žinant, kad aibė, turinti daugiau nei viena elementą, yra defektinė, kitame žingsnyje testuojamas tikrinis jos poaibis. Ši klasė labai plati. MD, ST ir PT procedūros jai priklauso.

<sup>10</sup>sarašas generuotas 2022 metų birželio 28 d.; straipsniai parašyti ne anglų kalba nenagrinėti

[30, 43, 48, 53, 55, 56, 76, 83, 84] vienintelis Malinovsky [53] nagrinėjo problemą, turinčią tiesioginį ryšį su PT procedūra. Visi likę tyrejai citavo Yao ir Hwang darbą [85] tik kaip turintį ryšį su jų sprendžiamais uždaviniais. Peržiūrėta literatūra lémė pasirinkimą pateikti detalesnę tikimybinię PT procedūros charakterizaciją. Darbe [91] mums pavyko išvesti a.d.  $T_{PT}$  momentų generuojančios funkcijos (MGF) išraišką ir jos pagalba irodyti kelias ribines teoremas — centrinę ribinę (CRT), didžiųjų skaičių dėsnį (DSD) ir didelių nuokrypių principą (DNP). Prieš pateikdami rezultatų formuliuotes, ivesime keliis pažymėjimus.

Tegu  $\Theta_N \equiv T_{PT}$  žymi bendrą testų skaičių, kurio reikia norint identifikuoti visus defektinius objektus dydžio  $N$  aibėje. Tegu  $Y_i \sim Be(p)$ ,  $i = 1, \dots, N$  žymi  $i$ -ojo objekto būsenos indikatorių (1 atitinka defektinį). Galiausiai, tegu  $\bar{Y}_i := 1 - Y_i$  ir

$$M_0 = \begin{pmatrix} 1 & 0 \\ 1 & 0 \end{pmatrix}, \quad M_1 = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}. \quad (\text{S.12})$$

Pirmasis rezultatas nusako  $\Theta_N$  struktūrą įvestų dydžių terminais.

**Teiginys S.0.3.** Tegu  $A = \left\{ \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}, \begin{pmatrix} 0 & 1 \\ 0 & 1 \end{pmatrix} \right\}$  ir  $B_k = \begin{pmatrix} Y_k & \bar{Y}_k \\ 1 & 0 \end{pmatrix}$ ,  $k = 1, \dots, N$ . Tada  $\Theta_2 = 3Y_2 + \bar{Y}_2(1 + Y_1)$ ,  $\Theta_3 = 2 + \bar{Y}_3Y_2 + Y_3\Theta_2$ ,

$$\begin{aligned} \Theta_N &= 1 + Y_N(\bar{Y}_{N-1}Y_{N-2} + 2) + Y_{N-1} + \\ &\sum_{j=3}^{N-1} (\bar{Y}_{j-1}Y_{j-2} + Y_{j-1} + 1) (Y_j + \bar{Y}_{j-1}\mathbb{1}\{B_N B_{N-1} \cdots B_{j+1} \in A\}) + \\ &Y_2 + \bar{Y}_2\mathbb{1}\{B_N B_{N-1} \cdots B_3 \in A\}, \text{ kai } N \geq 4. \end{aligned}$$

Kitas rezultatas aprašo minėtą išreikštinį MGF pavidał.

**Teorema S.0.10.** Tegu  $M_{\Theta_N}(\lambda)$  žymi a.d.  $\Theta_N$  MGF taške  $\lambda \in \mathbb{R}$ . Apibrėžkime

$$\alpha_i = \alpha_i(\lambda) = \frac{1}{2} \left( pe^{2\lambda} + (-1)^i \sqrt{p^2 e^{4\lambda} + 4qe^\lambda(q + pe^\lambda)} \right), \quad i = 0, 1; \quad (\text{S.13})$$

$$\kappa_N = \kappa_N(\lambda) = \frac{\alpha_0^N - \alpha_1^N}{\alpha_0 - \alpha_1} \text{ su } N \geq 0. \quad (\text{S.14})$$

$\forall N \geq 3$  teisinga lygybė  $M_{\Theta_N}(\lambda) =$

$$e^{2\lambda} \left[ \left( (1-q)^2 e^{3\lambda} + q(1-q)^2 e^{2\lambda} + q(1-q^2)e^\lambda + q^2 \right) \kappa_{N-2} + q \left( (1-q)^2 e^{3\lambda} + q(1-q)(2-q)e^{2\lambda} + 2q^2(1-q)e^\lambda + q^3 \right) \kappa_{N-3} \right].$$

Like rezultatai išplaukia iš teoremos [S.0.10](#).

**Išvada S.0.11.**  $E \Theta_N = N \frac{2-q^2}{1+q} + \frac{q^2+q-1}{(1+q)^2} (1 - (-q)^N),$

$$\begin{aligned} \text{Var } \Theta_N &= N \frac{(1-q)}{(q+1)^3} \left( q (q^3 + 3q^2 + 5q + 4) + \right. \\ &\quad \left. (-q)^N (2q+4) (q^2 + q - 1) \right) + \\ &\quad \frac{\left( 1 - (-q)^N \right)}{(q+1)^4} \left( q (5q^2 + 3q - 7) + (-q)^N (q^2 + q - 1)^2 \right), N \geq 3. \end{aligned}$$

**Išvada S.0.12.** Jei  $N \rightarrow \infty$ , tai a.d.  $\Theta_N$  tenkina žemiau nurodytus sąryšius.

DSD:  $\frac{\Theta_N}{N} \xrightarrow{L_2} \frac{2-q^2}{1+q}$  ir  $\frac{\Theta_N}{N} \xrightarrow{b.v.} \frac{2-q^2}{1+q}$ .

CRT:  $\sqrt{N} \left( \frac{\Theta_N}{N} - \frac{2-q^2}{1+q} \right) \xrightarrow{d} N(0, \sigma^2)$ ,  $\sigma^2 = \frac{q(1-q)(q^3+3q^2+5q+4)}{(q+1)^3}$ .

DNP:  $\frac{\Theta_N}{N}$  tenkina didelių nuokrypių principą (DNP) su nuokrypių funkcija  $I$ , lygia  $\mathbb{R} \ni \lambda \mapsto \ln \alpha_0(\lambda)$  Ležandro transformacijai ( $\alpha_0(\lambda)$  apibrėžiama [\(S.13\)](#)). Kitaip tariant, bet kokiai uždarai  $C \subset \mathbb{R}$  ir bet kokiai atvirai  $O \subset \mathbb{R}$ ,

$$\limsup_{N \rightarrow \infty} \frac{1}{N} \ln P \left( \frac{\Theta_N}{N} \in C \right) \leq - \inf_{x \in C} I(x)$$

ir

$$- \inf_{x \in O} I(x) \leq \liminf_{N \rightarrow \infty} \frac{1}{N} \ln P \left( \frac{\Theta_N}{N} \in O \right),$$

su  $I(x) = \sup_{\lambda \in \mathbb{R}} (x\lambda - \ln \alpha_0(\lambda))$ .

## Išvados

**Apie disertacijos turinį** Šiame disertaciame darbe ištirtos keturios grupinio testavimo procedūros — modifikuota Dorfman (MD), Sterrett (ST), kvadratinės matricos (A2) ir porinio testavimo (PT). MD ir ST procedūroms matematiškai pagrįstos iki tol literatūroje hipotetinėmis laikytos analizinės optimalių konfigūracijų išraiškos. A2 procedūros optimalios konfigūracijos analizinės išraiškos pavidalas buvo nežinomas. Disertacijoje jis surastas ir matematiškai pagrįstas.

PT procedūrai surasta testų skaičiaus skirstinio momentų generuojanti funkcija. Pasinaudojant ja (tinkamai transformuotam testų skaičiui) įrodytos trys ribinės teoremos — centrinė ribinė teorema, didžiųjų skaičių dėsnis ir didelių nuokrypių principas.

Be jau aprašytų uždavinių išspręstas dar vienas — pasiūlytas procedūros apytikrio optimalaus tikimybinio slenksčio radimo algoritmas, tinkantis gana plačiai binominio testavimo procedūrų klasei.

Visais atvejais gauti rezultatai iliustruoti taikymu pavyzdžiais.

**Apie rezultatų reikšmę** Disertacijoje nagrinėtos tikimybinio testavimo procedūros, tenkinančios binominio testavimo prielaidas. Vienuose tai- kymuose šios prielaidos yra pateisinamos (pvz., gamyba, kompiuterių mokslas), kituose per daug ribojančios (pvz., kai kurie medicininiai ar socialiniai taikymai), todėl disertacijos rezultatai tinka būtent pirmųjų atvejų analizei. Nepaisant to, kad grupinio testavimo procedūros, tinkerančios specifiniams atvejams, aktyviai tebevystomos, disertacijoje nagrinėtosios nėra pasenusios ir dažnai naudojamos masiniame testavime. Geras pavyzdys — neseniai praėjusi COVID-19 pandemija, kurios metu Dorfman procedūra buvo naudojama ir Lietuvoje testuojant moksleivius mokyklose [7] bei kitose šalyse [14]. Atsižvelgiant į išsakytas pastabas, binominio testavimo procedūrų tyrimai išlieka aktualūs ir disertacijoje gautus rezultatus galima laikyti naudingu indėliu į binominio testavimo teoriją.

## Rezultatų naujumas

Autorės žiniomis, visi disertacijoje gauti rezultatai yra nauji ir iki šiol literatūroje jokia forma nepublikuoti.

## Aprobacija

Disertacijos rezultatai pristatyti keturiose tarptautinėse mokslinėse konferencijose.

1. Čižikovienė, Ugnė; Skorniakov, Viktor. *On the optimal configuration of a square array group testing algorithm* // NBBC21: 8th Nordic-Baltic biometrics virtual conference, 7-10 June 2021, Helsinki, Finland.
2. Čižikovienė, Ugnė; Skorniakov, Viktor. *On the optimal configuration of the modified Dorfman and Sterrett group testing schemes* // IBC2022: 31st International Biometric Conference, 10 - 15 July 2022, Riga, Latvia.
3. Čižikovienė, Ugnė; Skorniakov, Viktor. *On the optimal Pairwise Group Testing Algorithm* // ECMI 2023: 22nd ECMI conference on industrial and applied mathematics, June 26 – 30, 2023, Wrocław, Poland.
4. Čižikovienė, Ugnė; Skorniakov, Viktor. *On the Generic Cut-Point Detection Procedure in the Binomial Group Testing* // The international scientific conference dedicated to the 160th anniversary of Prof. Dr. Hermann Minkowski, June 20 – 22, 2024, Kaunas, Lithuania.

## Publikacijos

Disertacija parašyta keturių žemiau išvardintų straipsnių pagrindu. Trys pirmieji publikuoti **Web of Science** indeksuojamuose žurnaluose. Ketvirtasis priimtas spaudai į Web of Science indeksuojamą žurnalą TWMS *Journal of Applied and Engineering Mathematics*.

1. Čižikovienė, Ugnė, and Viktor Skorniakov. "On a Couple of Unresolved Group Testing Conjectures". *Communications in Statistics - Theory and Methods*, vol. 52, no. 8, Apr. 2023, pp. 2448–60. DOI.org (Crossref), <https://doi.org/10.1080/03610926.2021.1953531>.
2. Čižikovienė, Ugnė, and Viktor Skorniakov. "On the Optimal Configuration of a Square Array Group Testing Algorithm". *Statistics and Its Interface*, vol. 16, no. 4, 2023, pp. 579–91. DOI.org (Crossref), <https://doi.org/10.4310/SII.2023.V16.N4.A1>.

3. Čižikovienė, Ugnė, and Viktor Skorniakov. "On the Optimal Pairwise Group Testing Algorithm". *Brazilian Journal of Probability and Statistics*, vol. 38, no. 2, June 2024. DOI.org (Crossref), <https://doi.org/10.1214/24-BJPS603>.
4. Čižikovienė, Ugnė, and Viktor Skorniakov. "On the Generic Cut-Point Detection Procedure in the Binomial Group Testing". *arXiv*: 2304.07263, *arXiv*, 14 Apr. 2023. arXiv.org, <https://doi.org/10.48550/arXiv.2304.07263>.

## Trumpos žinios apie autorię

### Išsilavinimas ir kvalifikacija

- 2000 - 2012** Vilniaus Gabijos gimnazija;
- 2012 - 2016** Vilniaus Universitetas, Matematikos ir Informatikos fakultetas, ekonometrijos studijų programa, statistikos bakalaureas;
- 2016 - 2018** Vilniaus Universitetas, Ekonomikos fakultetas, ekonominės analizės ir planavimo programa, ekonomikos magistras;
- 2020 - 2024** Vilniaus Universitetas, Matematikos ir Informatikos fakultetas, matematikos doktorantūra.

### Mokslinio ir pedagoginio darbo patirtis

- 2016** UAB "ATEA" praktikantė finansų analitikos pozicijoje;
- 2022 - 2023** Vilniaus Universiteto jaunesnioji asistentė.

### Darbo patirtis

- 2014 - 2021** Vilniaus miesto savivaldybės visuomenės sveikatos biuras, sveikatos stebėsenos specialistė;
- 2017 - 2021** UAB „Coface Baltics Services“, kredito rizikos vertinimo jaunesnioji analitikė;
- 2021 -** UAB „Agrochema“, verslo duomenų analitikė.

# Curriculum Vitae

## ASMENINĖ INFORMACIJA

Vardas, Pavardė | **Ugnė Čižikovienė**  
Gimimo vieta | Vilnius

## IŠSILAVINIMAS

Išsilavinimas Data Įstaiga	Vidurinės mokyklos 2000 - 2012 Vilniaus Gabijos gimnazija
Išsilavinimas Data Profesija Laipsnis Įstaiga Fakultetas	Aukštasis universitetinis 2012 - 2016 Ekonometrija Statistikos bakalauros Vilniaus Universitetas Matematikos ir Informatikos fakultetas
Išsilavinimas Data Profesija Laipsnis Įstaiga Fakultetas	Aukštasis universitetinis 2016 - 2018 Ekonominių analizės ir planavimų (anglų k.) Ekonomikos magistras Vilniaus Universitetas Ekonomikos fakultetas
Išsilavinimas Data Profesija Laipsnis Įstaiga Fakultetas	Doktorantūra (PhD) 2020 - 2024 Matematika Matematikos mokslų daktaras Vilniaus Universitetas Matematikos ir Informatikos fakultetas

## **DARBO PATIRTIS**

Darbovieta Data Pareigos	UAB „Agrochema“ 2021 - Verslo duomenų analitikė
Darbovieta Data Pareigos	Vilniaus universitetas 2022 - 2023 Jaunesnioji asistentė
Darbovieta Data Pareigos	UAB „Coface Baltics Services“ 2017 – 2021 Kredito rizikos vertinimo departamento jaunesnioji analitikė
Darbovieta Data Pareigos	UAB „ATEA“ 2016 Praktikantė finansų analitiko pozicijoje
Darbovieta Data Pareigos	Vilniaus miesto savivaldybės visuomenės sveikatos biuras 2014 – 2017 Sveikatos stebėsenos specialistė

## **DARBO SU KOMPIUTERIU GEBĖJIMAI**

<b>Microsoft Office paketas</b>	Word, Excel, Outlook, PowerPoint
<b>Programavimo kalbos</b>	SQL, Qlik sense, R/RStudio, Python

## **KALBU GEBĖJIMAI**

Gimtoji kalba	Lietuvių
Užsienio kalbos	Anglų, rusų, italų

# Declarations

**Conflict of interest/Competing interests** The author declares that she has no conflicts of interest.

**Usage of Generative Artificial Intelligence tools** The author declares that she did not use any generative artificial intelligence tools for the writing of the text and proofs of this manuscript, nor for the creation of graphics or their corresponding captions.

**Acknowledgements** The author thanks reviewers for pointing out several typos and giving constructive suggestions, which improved the exposition of the presented material.

Ugnė Čižikovienė

Vilnius

29th January 2025

# Publications by the Author

The dissertation is based on the four articles listed below. The first three were published in Web of Science indexed journals. The manuscript of the fourth is accepted for publication in Web of Science indexed journal *TWMS Journal of Applied and Engineering Mathematics*.

1. Čižikovienė, Ugnė, and Viktor Skorniakov. "On a Couple of Unresolved Group Testing Conjectures". *Communications in Statistics - Theory and Methods*, vol. 52, no. 8, Apr. 2023, pp. 2448–60. DOI.org (Crossref), <https://doi.org/10.1080/03610926.2021.1953531>.
2. Čižikovienė, Ugnė, and Viktor Skorniakov. "On the Optimal Configuration of a Square Array Group Testing Algorithm". *Statistics and Its Interface*, vol. 16, no. 4, 2023, pp. 579–91. DOI.org (Crossref), <https://doi.org/10.4310/SII.2023.V16.N4.SII746>.
3. Čižikovienė, Ugnė, and Viktor Skorniakov. "On the Optimal Pairwise Group Testing Algorithm". *Brazilian Journal of Probability and Statistics*, vol. 38, no. 2, June 2024. DOI.org (Crossref), <https://doi.org/10.1214/24-BJPS603>.
4. Čižikovienė, Ugnė, and Viktor Skorniakov. "On the Generic Cut-Point Detection Procedure in the Binomial Group Testing". *arXiv: 2304.07263*, arXiv, 14 Apr. 2023. arXiv.org, <https://doi.org/10.48550/arXiv.2304.07263>.



Ugnė Čižikovienė

Investigations of Binomial Group Testing Models

Doctoral Dissertation

Natural Sciences

Mathematics (N 001)

Thesis Editor: -

Binominių Grupinio Testavimo Modelių Tyrimai

Daktaro disertacija

Gamtos mokslai

Matematika (N 001)

Santraukos redaktorė: -

Vilnius University Press  
9 Saulėtekio Ave., Building III, LT-10222 Vilnius  
Email: [info@leidykla.vu.lt](mailto:info@leidykla.vu.lt), [www.leidykla.vu.lt](http://www.leidykla.vu.lt)  
Print run of 30 copies