

MASTER

**Privacy-preserving biometric databases
from authentication to efficient identification systems**

de Vreede, N.

Award date:
2013

[Link to publication](#)

Disclaimer

This document contains a student thesis (bachelor's or master's), as authored by a student at Eindhoven University of Technology. Student theses are made available in the TU/e repository upon obtaining the required degree. The grade received is not published on the document as presented in the repository. The required complexity or quality of research of student theses may vary by program, and the required minimum study period may vary in duration.

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain

Privacy-preserving Biometric Databases: from Authentication to Efficient Identification Systems

Niels de Vreede

DEPARTMENT OF MATHEMATICS AND COMPUTER SCIENCE
EINDHOVEN UNIVERSITY OF TECHNOLOGY

Graduation Supervisor:
Boris Škorić

Committee:
Boris Škorić
Frans Willems
Berry Schoenmakers

7 May 2013

Abstract

Nowadays, biometrics form an increasingly interesting means of performing authentication as well as identification. Because biometric data contain privacy sensitive information it is important to apply security measures when storing these data. Securely storing these data is complicated by the fact that biometric measurements are noisy.

In this thesis two issues are addressed. First, the matter of storing biometric data in such a manner that privacy can be preserved while still allowing data obtained from noisy measurements to be used to perform authentication. Second, the issue of modifying such a privacy preserving authentication system to be able to perform identification.

The issue of privacy-preserving storage of biometric data is investigated using general helper data systems. The error characteristics of such a system are optimized. A large part of this thesis focuses on systems with the *zero leakage* property. This property ensures that data stored for the purpose of compensating measurement noise do not reveal any information about other any other data that are stored, making it an interesting property for the purpose of privacy preservation. The construction and properties of zero leakage helper data systems are described.

An analysis of the modifications needed to make an authentication system suitable for performing identification is also given in this thesis. The advantage of creating an identification system based on an authentication system is that the privacy properties of the underlying authentication system can largely be preserved. The properties of the resulting identification system are studied analytically where possible and numerically where necessary. Some information theoretic results on this issue are also presented.

Contents

1	Introduction	4
1.1	Motivation	4
1.2	Helper Data Systems	5
1.3	Identification	5
1.4	Independent Dimensional Components	6
1.5	Original Research Questions and Project Goals	6
1.6	Modus Operandi and Short Summary of Results	6
2	Preliminaries	8
2.1	Notation	8
2.2	Definitions	9
2.2.1	Generalized Template System	9
2.2.2	Error Rates	9
2.2.3	Template Database	10
2.2.4	Identification	11
2.2.5	Running Time	13
2.2.6	Defining the Preselection Function	15
2.2.7	Quantization Helper Data System (QHDS)	16
2.2.8	Noise Model	19
2.2.9	Zero Leakage	20
2.2.10	Partitioning Template System (PTS)	20
3	Maximum Likelihood Reproducer	23
3.1	A More Usable Form	23
3.1.1	Discrete \mathcal{W}	24
3.1.2	Continuous \mathcal{W}	25
3.2	Zero Leakage	26
3.3	General Zero Leakage Scheme	28

3.4	Partitioning Template System	28
4	Properties of Partitioning Template Systems	30
4.1	Error Characteristics	30
4.2	Running Time Improvement	34
4.2.1	Gaussian Distribution	35
4.2.2	Numerical Results	35
4.3	Multidimensional Data	40
5	Information Theoretic Approach	42
5.1	Fano's Inequality	42
5.2	Multi-stage Selection	43
5.2.1	General Case	43
5.2.2	Fixed Cardinality	44
5.3	Numerical Results	45
6	Conclusion	48
6.1	Summary	48
6.2	Future Work	48

Chapter 1

Introduction

1.1 Motivation

Biometrics have become a popular solution for authentication or identification, mainly because of its convenience. Unlike other means of authentication, such as passwords or tokens, biometric features cannot be forgotten or lost. Nowadays passports and other identity documents commonly include biometric information extracted from fingerprints, irises or the face of the document holder. Governments maintain biometric databases for fighting crime and biometrics-based user authentication exists for laptops, smart phones and cars.

Even though biometrics are not strictly secret, storing biometric information in an unprotected database carries both security and privacy risks for the users. Security risks include the production of fake biometrics from the features, e.g. rubber fingers [16, 11]. These fake biometrics can be used to obtain unauthorized access to information or services or to leave fake evidence at crime scenes.

Two privacy risks would be that:

- some biometrics are known to reveal diseases and disorders of the user; and
- unprotected storage allows for cross-matching between databases.

Simply encrypting biometric databases does not fully solve the privacy and security issues, as it does not prevent *inside attacks*, i.e., attacks by people who are authorized to access the database, because they possess the decryption keys and therefore have full access to the biometric information.

The issue of privacy preserving storage of biometric data is very similar to that of password storage. The standard solution to securely storing passwords is by using *cryptographic hash functions*, one-way functions that are computationally infeasible to invert. An attacker can not obtain the password even when the hash is known.

The method of hashing passwords cannot be straightforwardly applied to biometric data, however. In order to preserve secrecy, cryptographic hash functions are very sensitive to small deviations in input and should produce very different results for two different input values, regardless of how similar they may be. This property of cryptographic hash functions requires that the data presented during enrollment must be exactly equal to the data presented during later verification. Biometric measurements are inherently noisy, which causes enrollment and verification measurements to differ. Straightforward quantization of the measurement data will not generally solve this issue, as enrollment measurements may lie arbitrarily close to quantization boundaries, which could allow even a slight amount of noise to lead to a different (quantized) verification measurement.

1.2 Helper Data Systems

In order to solve this problem, special algorithms have been developed [10, 3, 9, 8]: Fuzzy Extractors (FE) and Secure Sketches (SS), collectively also referred to as Helper Data systems (HDS).

A HDS works by computing public *helper data* and a high entropy *secret* from the enrollment measurement. The secret can be treated like a password and stored using a cryptographic hash function, along with the helper data. Any later measurements will be close, but not equal to the enrollment measurements. Using the helper data and such a close measurement, the secret can be reconstructed and used for authentication similar to a password by comparing the hash of the reconstructed secret to the stored hash value. A good HDS is characterized by giving, with high probability, the same hash as during enrollment, if a biometric sample from the same person is presented and a different hash is a sample from a different person is presented.

The FE and SS are special cases of the general HDS concept, with different purposes. For the FE the probability distribution of the secret given the helper data is required to be near uniform. The FE is typically used for the extraction of cryptographic keys from noisy sources such as Physical Unclonable Functions [15, 5, 13, 14, 12, 6]. For the SS the secret must have high entropy when conditioned on the helper data. It is typical for the SS's secret to be equal to a discretized version of the enrollment measurement. The SS is very well suited to the biometrics scenario described above. The high entropy requirement ensures that even when the helper data are known, it is difficult to obtain the secret from the hash. Thus, HDS systems achieve *data reconciliation* (error correction) while satisfying security requirements at the same time.

The main privacy question for HDSs is how much, and which information the public data reveal about the biometric. Ideally, the helper data should only contain enough information to enable error correction. This means that the helper data can be considered the ‘noisy part’ the biometric measurement and, as such should reveal only little information about the biometric. In order to quantify such statement, HDSs should be studied information theoretically. Of particular interest is the mutual information contained between the secret and the helper data. It is possible to construct systems where knowledge of the helper data reveals absolutely nothing about the secret. This property will be called *Zero Secrecy Leakage*, or shorter *Zero Leakage (ZL)*. HDSs with the ZL property are interesting for privacy preserving biometric storage. All privacy sensitive information contained in the secret is guaranteed not to be revealed by the helper data.

1.3 Identification

A HDS as described in the previous section can be used to perform biometric *authentication*, i.e., a user's claim of a given identity can be verified. In *identification* it is a priori unknown which (if any) of the database entries corresponds to the person under scrutiny. Simply attempting to authenticate a given biometric against every entry in the database may turn out to be too time consuming, so for an identification system the *running time* is an important issue.

This problem forms the main motivation for this thesis: how to search efficiently in a biometric database when the entries contain only helper data and hashes? In a database that is made purely for authentication, the helper data must, in order to preserve privacy, reveal as little as possible about the identity of an enrolled person. In this case the helper data provides no handle for efficiently searching the database. The cryptographic hashes are, by design, not suitable for searching either. If running time is considered important, clearly more revealing information has to be incorporated in the public data. The question is of course how much more revealing.

Another question is how best to search in a database whose entries are both noisy and of high dimensionality. The noise issue makes a lexicographic ordering much less useful, whereas the dimensionality issue means a sensible lexicographic ordering may not even exist. An interesting

search method was proposed by Voloshynovskiy et al. [18] for searching in the biometric measurement space. A list of (public) special ‘Beacon Points’ is fixed. A lexicographical ordering is then determined according to the distance of the biometric measurement to a beacon point for each of the beacon points. In the biometric measurement space, noise barely affects these distances, allowing for efficient searching. Furthermore, the authors show that relatively short search times can be achieved using far fewer beacon points than the dimensionality of the biometric data, thereby keeping the dimensionality issue in check.

The Beacon Points idea served as inspiration for this thesis. If applied in helper data space instead of biometric measurement space, the method should still work, but with the additional benefit of preserving the privacy characteristics of the underlying authentication system.

1.4 Independent Dimensional Components

For most part of this thesis, all dimensional components of multidimensional biometric data are considered independent of every other component. This can be justified by the use of well-known techniques, such as PCA, during preprocessing of the data. The independence of dimensional components allows for one-dimensional analysis of each individual component. In cases where independence does not suffice to justify one dimensional analysis, often only the one-dimensional case is studied. Some notes are provided on the issue of multidimensionality of the biometric data, but this is largely a subject for future work.

1.5 Original Research Questions and Project Goals

At the start of the project, the above considerations led to the following research questions and project goals.

1. How much entropy must helper data reveal in order to allow for efficient searching of the biometric database?
2. How will the Beacon Points method perform, in terms of running time, when applied in helper data space?
3. Investigate the trade-off between privacy and running time using information-theoretic techniques, and if possible see if the results of [7] follow as a limiting case.

1.6 Modus Operandi and Short Summary of Results

I had weekly discussions with my supervisor Boris Škorić and I joined the informal meetings on Information Theory with Frans Willems, Jean-Paul Linnartz and Joep de Groot at the Signal Processing Systems group of the Electrical Engineering Department. Due to these meetings the focus of my research shifted to Zero Leakage helper data. It has resulted in a publication [2].

The results of this project can be summarized as follows.

- An analysis of how to modify an existing authentication system to obtain an efficient identification system, while retaining the privacy characteristics (Section 2.2.4) and the running times that are expected for such systems (Section 2.2.5).
- An analysis of the issue of optimally reconstructing the enrolled value for an authentication system, both in general (Section 3.1) and under the ZL constraint (Section 3.2).

- How a ZL authentication system can be constructed for generally all continuous data (Sections 3.3 and 3.4)
- An analysis of the properties of a concrete HDS for one-dimensional data, including the expected improvements in running time of such a system relative to a naïve implementation and some research towards the trade-off between such a system’s resilience against errors, the running time and measure of privacy preservation (Chapter 4). The chapter also contains notes on the issues arising as a result of multidimensional data.
- Some information theoretical work towards the original research question of how much entropy must be revealed by the helper data to allow for efficient searching (Chapter 5)
- Contributions towards the formalization and proofs in Sections ‘Zero Leakage Key Extraction’ and ‘Optimal Reconstruction’ of [2], as well as the formalism presented in Section ‘Mismatch Between the Real and Assumed Distribution’.

Chapter 2

Preliminaries

2.1 Notation

Random variables will be denoted using capital letters. The corresponding calligraphic capital will indicate the set from which the random variables are drawn, and corresponding lower case letters are used for the elements of this set.

Due to the large number of different probability density functions used in this thesis, all of these will be denoted in a systematic way. The letter ρ , subscripted with a random variable will indicate that random variable's probability density function, e.g. $X \sim \rho_X$. Similarly, for joint and conditional probability density functions notation like $\rho_{X,Y}$ and $\rho_{X|Y}$ will be used respectively. Cumulative density functions are indicated by the symbol P , with the same subscript as their corresponding probability density function.

For some applications the exact probability density function is unknown and an estimated probability distribution has to be used, denoted by ξ . The symbol Ξ is used for the cumulative density function corresponding to this estimate.

The standard normal probability density function is often used and indicated by the symbol ϕ and its cumulative distribution function by Φ

The notation for probability mass functions is identical to that of probability density functions, to reduce notational clutter in equations. This makes it possible to elegantly describe joint distributions involving both discrete and continuous random variables and equations which hold regardless of particular random variables being discrete or continuous. Furthermore, it allows for the use of variables which have discrete as well as continuous properties. The rules of manipulating probability distributions remain largely unaffected. The symbol \mathbb{E} will be used to indicate the expectation value of a random variable.

Spaces and variables are in general of arbitrary dimensionality, therefore vectors will not be distinguished from scalars. Dimensional components will be indicated with a subscript. E.g., x can be a multidimensional variable, with x_d referring to its d -th dimensional component, which is scalar. The dimensionality of the corresponding space \mathcal{X} is denoted by $\dim \mathcal{X}$.

The symbol $\mathbb{B} = \{\text{True}, \text{False}\}$ will be used for the set of Boolean truth values and $\mathbb{N}_n = \{0, 1, \dots, n-1\}$, where n is a natural number, for finite ranges of natural numbers.

Both the Kronecker and Dirac deltas will be denoted using a lower case delta. The Kronecker delta will have its two discrete arguments in subscripts, e.g., $\delta_{a,b}$, whereas the Dirac delta will be indicated by the difference of two continuous variables in parentheses, e.g., $\delta(a-b)$. In either case, the variables can be multidimensional.

The power set of a set \mathcal{A} , i.e., the set of all subsets of \mathcal{A} will be denoted by $\wp(\mathcal{A})$.

The binary Boolean operator $\stackrel{?}{=}$ will evaluate to True if the operands are equal to each other and to False otherwise. Similarly, the binary Boolean operator $\stackrel{?}{\in}$ will evaluate to True if the left hand operand is a member of the right hand operand and to False otherwise.

2.2 Definitions

2.2.1 Generalized Template System

The research questions are addressed in the framework of an abstract *Generalized Template System* (GTS). The spaces \mathcal{U} and \mathcal{V} represent biometric sample spaces. \mathcal{U} is used for the sample that is to be enrolled in a database, whereas \mathcal{V} indicates a sample that is to be used to match against enrolled values in a database. The spaces \mathcal{U} and \mathcal{V} are distinct so as not to impose superfluous requirements, but could very well be chosen identical. If identical sensors were to be used under the same conditions for both enrollment and matching measurements then $\mathcal{U} = \mathcal{V}$ would hold.

For privacy reasons it is not desirable to store the full enrollment measurement in the database. Instead, only some derived data, or template, from space \mathcal{T} are stored.

The space \mathcal{R} is used as an external source of randomness.

Definition 1 (GTS). *A generalized template system on $\mathcal{U} \times \mathcal{R} \times \mathcal{V} \times \mathcal{T}$ is a pair of functions (Derive, Matches), where*

$$\begin{aligned} \text{Derive: } \mathcal{U} \times \mathcal{R} &\rightarrow \mathcal{T}; \text{ and} \\ \text{Matches: } \mathcal{V} \times \mathcal{T} \cup \{\text{Null}\} &\rightarrow \mathbb{B}. \end{aligned} \tag{2.1}$$

In Definition 1, the function Derive is used to derive so-called *public data*. Ideally, the public data are able to resist inside attacks. The predicate Matches decides whether a later measurement of a biometric matches the enrolled public data. The public data contain, or might entirely consist of, data used for error correction on the biometric measurements. The error correcting data are known as *helper data*.

The special symbol $\text{Null} \notin \mathcal{T}$ is such that $\forall v \in \mathcal{V}: \text{Matches}(v, \text{Null}) = \text{False}$. It serves no purpose other than to simplify the definition of the biometric database below (Definition 4).

2.2.2 Error Rates

Due to the noisy nature of sampling biometrics, it is generally not possible to perform authentication or identification completely without error. Therefore, it is necessary to analyse the rate at which errors are expected to occur.

There are two basic types of error events that need to be distinguished.

- *False reject events* are those events where a sample is not matched to its corresponding enrolled sample, i.e., $\text{Matches}(v, \text{Derive}(u, r)) = \text{False}$ for some $u \in \mathcal{U}$ and $v \in \mathcal{V}$ taken from the same source.
- *False accept events* are those events where a sample is matched to an unrelated enrolled sample, i.e., $\text{Matches}(v, \text{Derive}(u, r)) = \text{True}$ for some $u \in \mathcal{U}$ and $v \in \mathcal{V}$ taken from different sources.

In Section 2.2.5 identification systems will be considered where the false reject events are further subdivided.

The expected probability of an error occurring will be called the *error rate*. In this thesis, the false accept rate is assumed to be negligible compared to the false reject rate, and as such, only false rejects errors shall be considered in the error analysis.

To quantify the error probability, an indicator random variable is introduced.

Definition 2. *The error indicator random variable is defined as:*

$$E = \begin{cases} 0 & \text{if } \text{Matches}(V, \text{Derive}(U, R)) = \text{True} \\ 1 & \text{otherwise} \end{cases}. \quad (2.2)$$

Using E , the expected error rate can be determined.

Definition 3. *The expected error rate, ε , is defined as the expectation value of the indicator random variable:*

$$\varepsilon = \mathbb{E}[E]. \quad (2.3)$$

The error rate for a given verification sample is also defined:

$$\varepsilon(v) = \mathbb{E}[E|V = v]. \quad (2.4)$$

2.2.3 Template Database

The public data derived from a biometric measurement are to be stored in a template database, which allows later retrieval of the stored data to match these against a different biometric sample. A *Template Database* (TDB) is considered merely a mapping of some set of identifiers to a set of public data. The matter of associating these with actual identities is beyond the scope of this project.

Definition 4 (TDB). *A template database is a function $d: \mathbb{N} \rightarrow \mathcal{T}$, with the following additional properties:*

- $|d|$ is the number of entries in the database;
- $\forall n \notin \mathbb{N}_{|d|}: d(n) = \text{Null}$; and
- a GTS is associated with the database, denoted by $(d[\text{Derive}], d[\text{Matches}])$

In this definition the set of identifiers is chosen to be the natural numbers. There is little, if any, need to generalize this. By associating the database with a given GTS, it becomes tied to the given derivation and matching functions and implicitly defines the spaces on which a database exists. The symbol \mathcal{D} shall be used to denote the space of all databases. Associating a GTS with a database also fixes the derivation and matching functions before any enrollment is performed.

There are three operations defined on databases: enrollment, authentication and identification. Enrollment is the process of adding a new entry, containing the derived data of some biometric measurement, to the database. Issues regarding enrollment shall not be considered in this thesis and no formal definition is given. The enrollment operation has the intuitive meaning of inserting data derived from a biometric measurement into the database and associating it with a unique identifier.

Authentication is the process of verifying a biometric sample against the enrolled derived data for a given identifier, defined as:

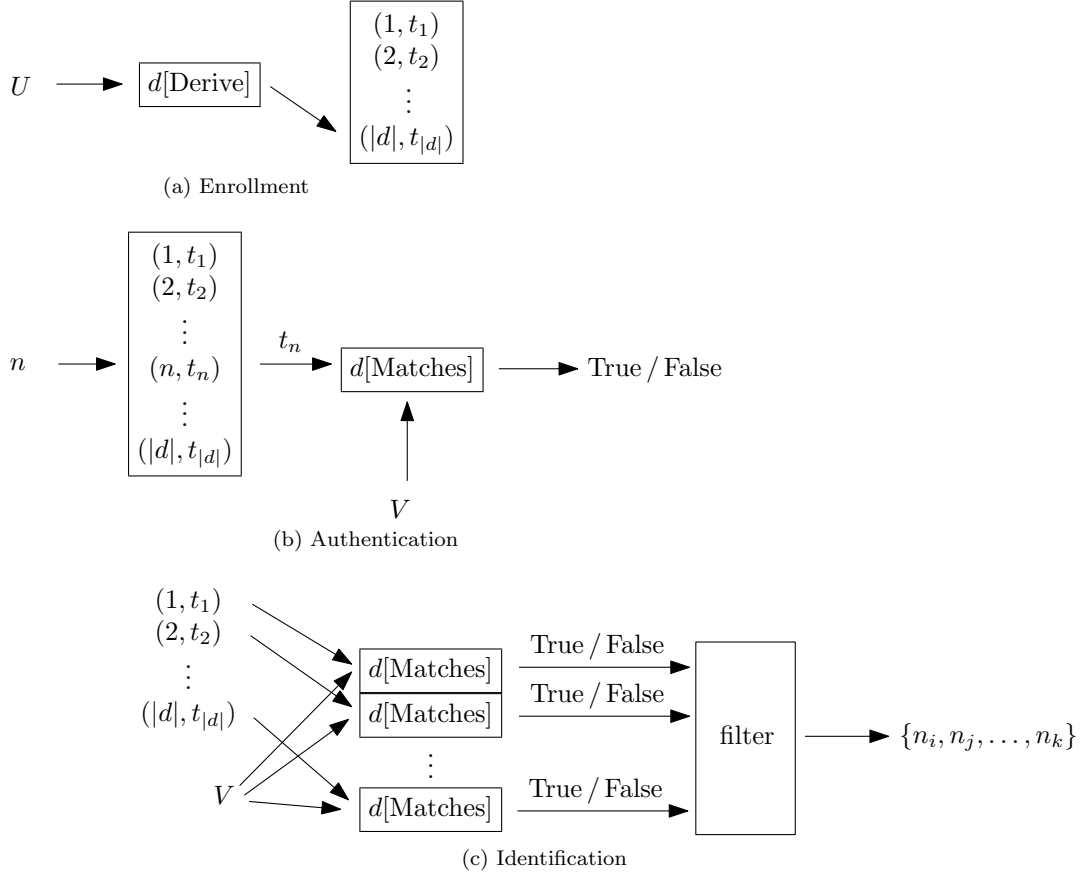


Figure 2.1: The three operations defined on an abstract template database: enrollment, authentication and identification.

Definition 5 (Authentication). *The authentication operation is a function $\text{Authenticate}: \mathcal{V} \times \mathbb{N} \times \mathcal{D} \rightarrow \mathbb{B}$, such that*

$$\text{Authenticate}(v, n, d) = d[\text{Matches}](v, d(n)). \quad (2.5)$$

Identification is similar to authentication, except that no identifier is given. Rather, all identifiers associated with derived data against which a given biometric sample matches are selected from the database:

Definition 6 (Identification). *The identification operation is a function $\text{Identify}: \mathcal{V} \times \mathcal{D} \rightarrow \wp(\mathbb{N})$ such that*

$$\text{Identify}(v, d) = \{n \in \mathbb{N} \mid d[\text{Matches}](v, d(n))\}. \quad (2.6)$$

A schematic representation of these three operations is displayed in Figure 2.1.

2.2.4 Identification

Definitions 5 and 6 imply the following relation between authentication and identification:

$$\forall v \in \mathcal{V}, d \in \mathcal{D}, n \in \mathbb{N}: \text{Authenticate}(v, n, d) \Leftrightarrow n \in \text{Identify}(v, d), \quad (2.7)$$

which raises practical concerns. In an authentication system, one can afford to use a time-consuming matcher, because it will only be applied to a single database entry. If an authentication

system were to be naively used as an identification system by simply performing the authentication procedure on each entry in the database, as suggested by Equation (2.7), the matcher would be executed for each entry in the database, scaling the time consumption by as much. This may be undesirable and demonstrates that matchers suitable for authentication are not necessarily suitable for identification.

In order to implement an identification system based on a GTS suitable for an authentication system that is not subject to this issue, the matching process has to be modified. There are two means of reducing the time consumption: multi-stage selection and first match with ordered evaluation.

Multi-stage Selection

Instead of applying the full matcher to each database entry, a series of matchers is applied. Each matcher in this series has a short execution time compared to the matchers following it, but also makes relatively coarse decisions. Entries that are rejected by a certain matcher are not passed along to any further matchers, so that only a few will remain for the final, time consuming matcher, which is the original matcher the system is based on. The application of all but the last matcher is called preselection. In practice, only two stages might suffice, meaning the preselection consists of a single, coarse matcher.

Note that this multi-stage selection is different from merely optimizing the implementation of an existing matcher. It is possible for the preselection to filter out entries that would have passed the final matcher. This introduces a new type of error and makes the multi-stage matcher functionally distinct from the matcher it is based on.

First Match with Ordered Evaluation

In some applications it may be desirable for the result of identification to contain no more than a single identifier. This suggests extending the enrollment and identification processes (Definition 6) by introducing some kind of matching score on the database entries and selecting the entry with the best score, if any, but in general that would tend to increase the complexity of identifications. Instead, only systems are considered for which identification already satisfies the requirement of, with high likelihood, selecting at most one identifier.

In this case, processing can be stopped as soon as a matching entry is found. However, if the database entries are evaluated in arbitrary order, this only reduces the expected search time by a factor two if a matching entry actually exists and does not reduce the search time at all otherwise. Therefore, for this first match (FM) approach to be useful, some meaningful order for evaluation must exist. Furthermore, the complexity of determining this order must at identification time be significantly smaller than that of the matcher, otherwise its advantage would be lost. Finally, the likelihood of no matching entry existing in the database must be small.

It is not necessary to completely order the database before evaluation. An order can straightforwardly be imposed by using multiple matchers, as in the case of multi-stage selection. In contrast to multi-stage selection, entries that fail on one of the initial matchers are not discarded, but given a lower priority for evaluation by the final matcher. This approach is called *priority selection*.

Precomputation

In general, precomputation can reduce the time consumption of expensive computations such as the identification process. It is important to note that the multi-stage selection and the FM techniques may each be partially precomputed during enrollment for an efficient implementation. This precomputation is not formally included in the framework presented here, as it should only

reduce the time consumption of computations during identification, but does not functionally affect a system.

Composition

Finally, it is noted that the multi-stage selection can be interleaved with the ordering for finding the first match. For example, first some preselection can be made on the entire database, for only the resulting selection to be evaluated by the final matcher in a particular order. This would be suitable in systems where determining an order for the entire database would be infeasible, but relatively fast on only the preselection.

2.2.5 Running Time

The purpose of introducing the multi-stage selection and FM approaches to identification is to reduce the time consumption of the identification process. To quantify the efficiency of a system, the expected running time for a given identification sample is used. If so desired, the average or worst case expected running times can be computed from this.

In general, the running time of an identification process is highly dependent on the particular algorithms used. Although a fairly general analysis is possible, it is considered beyond the scope of this thesis. Therefore, in this section the expected running time of several relatively simple, but very useful systems is discussed.

To distinguish the properties of different matchers, all variables will be subscripted with a name referring to its particular matcher.

Authentication Matcher

The reason for introducing alternative identification matchers is that the naïve application of an authentication matcher is considered to be too time consuming. This will now be quantified.

The authentication matcher is a *simple* matcher, i.e., it is non-composite and does not halt when a match is found. Since a matcher whose running time when applied to a single entry depends on the value of the public data can be considered a composite matcher, the running time of the authentication matcher for a single entry is constant.

In general, the running time of the identification process for a simple matcher is a linear function of N , the number of entries to which it is applied: $T(v, N) = \tau + Nt$, where T is the expected running time of the identification process for a given sample and t is the running time for a single entry. Since part of the computation for a given sample may be the same for each entry, that part of the computation will only have to be carried out once and its running time is represented by τ . The running time of a simple matcher in an authentication scenario can be found by letting $N = 1$. To distinguish matchers and their running time, the variables may be labeled using subscripts.

The need for alternative matchers arises only when the running time of an authentication matcher when used for identification is deemed unsuitable. Since the running time in an authentication scenario is considered acceptable, the constant precomputation time τ_{Auth} is negligible. This gives an expected running time of

$$T_{\text{Auth}}(v, N) = Nt_{\text{Auth}}. \quad (2.8)$$

Authentication Matcher with First Match

When using an authentication matcher for identification, but letting it stop processing as soon as a match is found, rather than continuing to search the database, the running time may be reduced. For this FM approach to be advantageous it is necessary that a match be found at all. Since false accept events are ignored in this thesis, a match can only be found if an entry matching the verification sample exists in the database. Furthermore, the matcher must not fail in matching the verification sample to the corresponding database entry.

To express the running time, two cases will have to be distinguished: whether an entry corresponding to the given verification sample exists in the database or not. If such an entry exists, the verification sample is called *genuine*, otherwise it is called an *attacker* sample. Despite the name, the latter does not necessarily imply malicious intent.

In the attacker scenario, the entire database will be to be searched before concluding no entries match the given sample. The running time in this case will therefore be given by Nt_{Auth} . This holds for any purely FM based approach when false accept events are disregarded.

If the verification sample is genuine, and the authentication matcher does not make an error, processing can be terminated early. Since the database is processed in arbitrary order, an expected half of the database will have to be searched before the match is found. If the matcher does make an error, the entire database will have to be searched before erroneously reporting no match. The expected running time for a genuine sample is then given by

$$T_{\text{Auth, FM}}(v, N) = \frac{1}{2}(1 - \varepsilon(v))Nt_{\text{Auth}} + \varepsilon(v)Nt_{\text{Auth}} = \frac{1}{2}(1 + \varepsilon(v))Nt_{\text{Auth}}. \quad (2.9)$$

This indicates the necessity of imposing an evaluation order for a FM approach to result in significantly reduced running times.

Two Stage Preselection

For two stage preselection, the matcher is composed of a preselection matcher and the authentication matcher:

$$\text{Matches}_{2\text{H}}(v, t) = \text{Matches}_{\text{Pre}}(v, t) \wedge \text{Matches}_{\text{Auth}}(v, t), \quad (2.10)$$

where the authentication matcher is only applied to entries that pass the preselection matcher. The subscript 2H is used to indicate a two stage process where search is halted for entries which do not pass the first stage selection. The expected running time is given by

$$T_{2\text{H}}(v, N) = T_{\text{Pre}}(v, N) + T_{\text{Auth}}(v, \omega_{\text{Pre}}(v)N), \quad (2.11)$$

where $\omega_{\text{Pre}}(v)$ is the fraction of entries expected to pass the preselection, which is defined as follows:

Definition 7. For any matcher Matches , let $\mathcal{T}(v) = \{t \in \mathcal{T} \mid \text{Matches}(v, t)\}$. Then the expected selection fraction is defined as

$$\omega(v) = \int_{\mathcal{T}(v)} dt \rho_T(t). \quad (2.12)$$

Just as the authentication matcher, the preselection matcher is considered to be a simple matcher. Whereas ignoring the constant processing time for the authentication matcher is considered reasonable, for the preselection matcher this may not be the case. However, because the procedure is designed to reduce the total running time, it may be reasonable to disregard the processing time per entry giving running time of

$$T_{\text{Pre}}(v, N) = \tau_{\text{Pre}}. \quad (2.13)$$

Combining Equations (2.8)–(2.13) gives

$$T_{2H}(v, N) = \tau_{\text{Pre}} + \omega_{\text{Pre}}(v)Nt_{\text{Auth}}. \quad (2.14)$$

The preselection process may also be combined with a FM approach. In this case, the running time is reduced to

$$T_{2\text{HFM}}(v, N) = \tau_{\text{Pre}} + \frac{1}{2}(1 + \varepsilon_{\text{Auth}}(v))\omega_{\text{Pre}}(v)Nt_{\text{Auth}}. \quad (2.15)$$

Two Stage Priority Selection

In the preselection process, it is desirable for ω_{Pre} to be small, as this can greatly reduce the running time compared to the authentication matcher. However, this is not generally possible without increasing the error rate. It is possible to decrease the expected running time, while retaining the error rate of the authentication process by continuing to process the remainder of the database when no match is found in the initial selection. This approach is only advantageous when combined with a FM approach. For a given genuine sample, the expected running time is given by

$$\begin{aligned} T_{2\text{CFM}}(v, N) = \tau_{\text{Pre}} + \frac{1}{2} \varepsilon_{0,0} \omega_{\text{Pre}} N t_{\text{Auth}} \\ + \frac{1}{2} \varepsilon_{1,0} (1 + \omega_{\text{Pre}}) N t_{\text{Auth}} \\ + \varepsilon_{\text{Auth}} N t_{\text{Auth}}, \end{aligned} \quad (2.16)$$

where $\varepsilon_{0,0}$ indicates the expected probability of no error being made during either priority selection of the authentication matcher, i.e., $E_{\text{Prio}} = 0$ and $E_{\text{Auth}} = 0$. The expected probability of making an error during priority selection, $E_{\text{Prio}} = 1$, but performing authentication correctly, $E_{\text{Auth}} = 0$ is indicated by $\varepsilon_{1,0}$. The subscript 2CFM indicated a two stage FM based process in which the search continues in entries that do not pass the first selection, should no match be found within the first selection.

2.2.6 Defining the Preselection Function

The purpose of introducing composite matchers is to be able to use an authentication matcher for identification, but to reduce the running time compared to naïvely applying the authentication matcher to all entries in a database. Therefore, suitable preselection and priority selection matchers are those that minimize the running time. Furthermore, it is important for a preselection matcher to retain the authentication matchers error characteristics as much as possible. Although it is desirable for preselection to achieve both a low error rate, and a low selection fraction, it is not generally possible to minimize both quantities.

Because it is not possible to minimize both the expected error probability and the selection fraction, the preselection matcher that minimizes the selection fraction for an imposed expected error probability is used. This defines a relation between the expected error probability and the selection fraction.

In case of priority selection, the expected running time depends on both properties, and the optimal matcher is determined by the running time. In this case, the expected error probability of the preselection matcher has no influence on the expected error of the total matcher. For a two stage selection procedure, the relation between the two characteristics describes the characteristics of which systems are feasible.

2.2.7 Quantization Helper Data System (QHDS)

In this section a somewhat less abstract template system that will form the basis of all further analysis will be introduced. The symbols introduced for this system will be reserved throughout this thesis.

The biometric sample space for enrollment measurements is denoted by \mathcal{X} , the sample space for any later authentication or identification measurements by \mathcal{Y} . The spaces \mathcal{X} and \mathcal{Y} can be, and in reality typically are, multidimensional. It is assumed that each dimensional component is independent of the others. This is generally justified by assuming that the data have been preprocessed, for example, using principal component analysis.

The system to be introduced is inspired by Fuzzy Extractors and Secure Sketches [10, 3, 9, 8]. In short, a Fuzzy Extractor is a pair of functions (Gen, Rep), with $\text{Gen}(X) = (S, W)$ and $\text{Rep}(Y, W) = \hat{S}$, such that, if X and Y are taken from the same biometric, $S = \hat{S}$ with high likelihood. The variables S and W are called the *secret* and *helper data* respectively. Unlike a Fuzzy Extractor, no uniformity requirements are imposed on S in this thesis.

Formally, let some functions $\tilde{s}: \mathcal{X} \rightarrow \mathcal{S}$ and $\tilde{w}: \mathcal{X} \rightarrow \mathcal{W}$ exist for some discrete, finite space \mathcal{S} and some space \mathcal{W} , then \tilde{s} is called a quantization function and \tilde{w} a helper data function. Furthermore, let $\tilde{h}: \mathcal{S} \times \mathcal{Z} \rightarrow \mathcal{H}$ be some cryptographic hash function, where \mathcal{Z} and \mathcal{H} are some discrete, finite spaces. Then, a derivation function and matcher can be defined in terms of these functions. These will be defined separately below and together make up the so-called *Quantization Helper Data System* (QHDS).

QHDS Derivation Function

Definition 8. *The derivation function, Derive, of a QHDS on $(\tilde{s}, \tilde{w}, \tilde{h})$ is defined as*

$$\text{Derive}(x, z) = (\tilde{w}(x); \tilde{h}(\tilde{s}(x), z); z). \quad (2.17)$$

It would have been possible to include some external source of randomness in the domains of \tilde{s} and \tilde{w} , such as is the case for the GTS and is used for, e.g., the Code Offset Method [9], but this would make later analysis less clear, without providing any benefit for the purpose of this thesis.

Similarly, extending the domain of \tilde{s} with \mathcal{Z} and absorbing the hash function into \tilde{s} would hide this detail, but would also complicate later analysis. The need for a hash function is justified by the fact that the quantization function may reveal too much information about a biometric sample. This same issue limits the possible choices of \tilde{w} , but as the helper data are only intended to compensate for noise on the measurements, a good \tilde{w} function should not have $\tilde{w}(x)$ reveal too much information about x .

Altogether, these arguments form the motivation for this particular choice of the derivation function.

QHDS Matcher

There are two straightforward possible definitions for the matcher: the threshold matcher and the reproducer-based matcher.

Definition 9. *The threshold matcher is defined as*

$$\text{Matches}(y, (w, h, z)) = h \stackrel{?}{\in} \{\tilde{h}(s, z) \mid s \in \mathcal{S} \wedge \xi_{S|Y,W}(s|y, w) \geq \theta(y, s, w)\}, \quad (2.18)$$

where θ is some threshold which may depend on the variables involved.

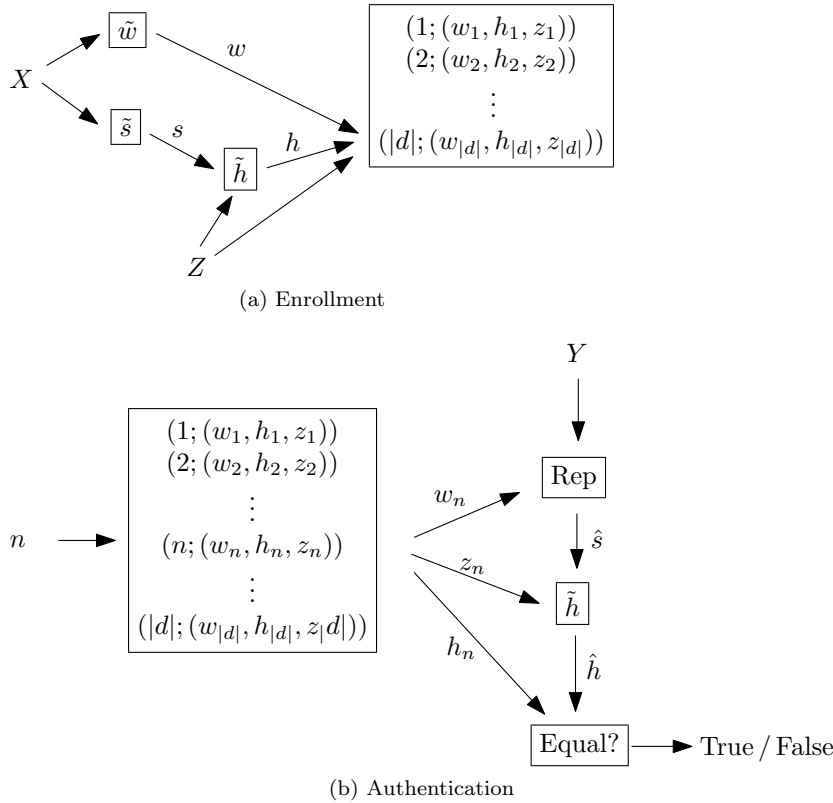


Figure 2.2: The enrollment and authentication procedures as defined on a QHDS.

Definition 10. The reproducer-based matcher is defined as

$$\text{Matches}(y, (w, h, z)) = h \stackrel{?}{=} \tilde{h}(\text{Rep}(y, w), z), \quad (2.19)$$

where Rep is a Fuzzy Extractor-like reproducer function.

The system characteristics of a QHDS using either of these two matchers will depend on the choice of \tilde{s} , \tilde{w} and \tilde{h} , in addition to θ or Rep for the threshold or reproducer-based matchers respectively. The reproducer-based approach is preferred, because it is possible to define the optimal reproducer for given \tilde{s} and \tilde{w} . Using the optimal reproducer to define the optimal reproducer-based matcher makes further analysis of this system possible.

Definition 11. The maximum likelihood reproducer, Rep_{ML} , for (\tilde{s}, \tilde{w}) is defined as

$$\text{Rep}_{\text{ML}}(y, w) = \arg \max_{s \in \mathcal{S}} \xi_{S|Y,W}(s|y, w). \quad (2.20)$$

This reproducer is defined in terms of the approximate distribution $\xi_{S|Y,W}$, as typically only an empirically determined approximation of the true distribution is known. When the approximation approaches the true distribution, this reproducer approaches optimality by definition. Definitions 10 and 11 can be combined to give the *maximum likelihood matcher*.

Combining Definitions 8, 10 and 11 defines the QHDS on $(\tilde{s}, \tilde{w}, \tilde{h})$. A schematic representation of the QHDS is displayed in Figures 2.2 and 2.3.

Although the definition of the maximum likelihood matcher may suggest a brute force implementation, it will be shown that this is not necessary for the concrete template system introduced below.

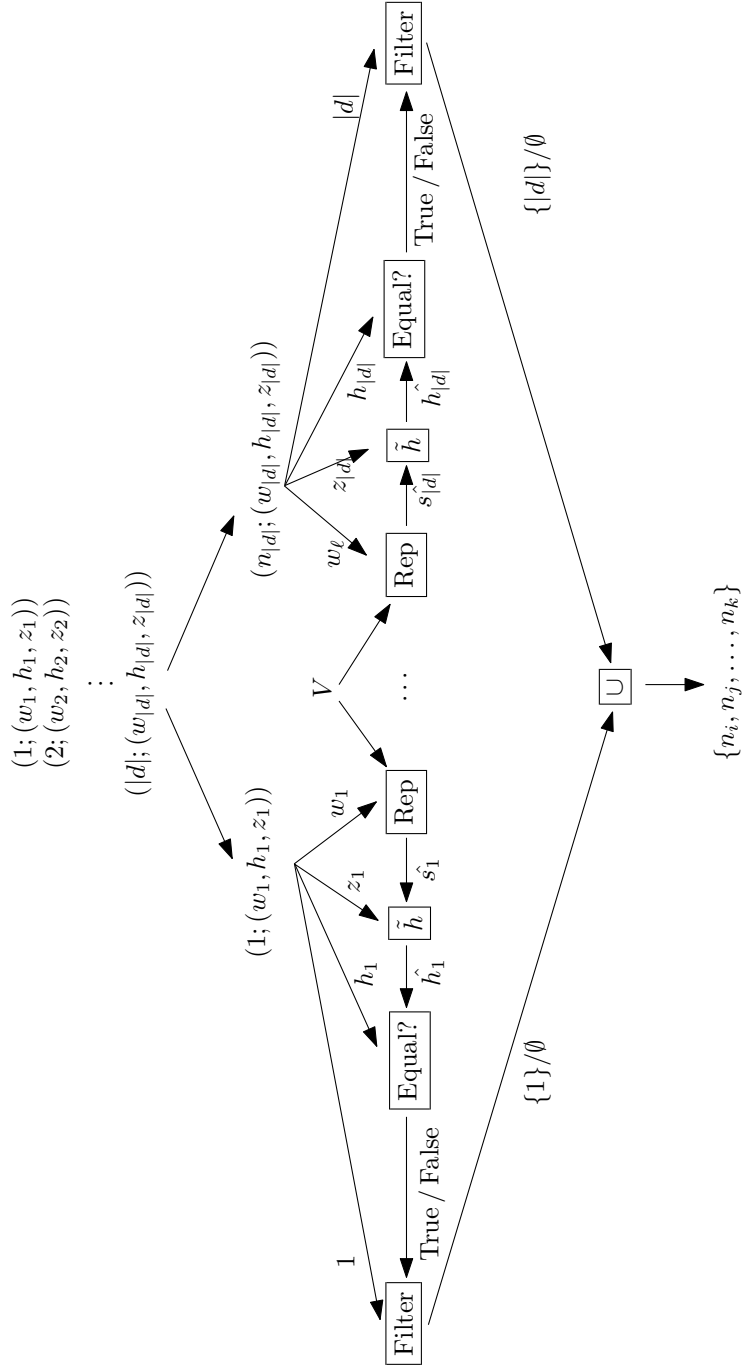


Figure 2.3: The identification procedure for a QHDS.

QHDS Preselection

The purpose of hashing the secret in the QHDS is to make the hashed value usable only for equality testing, while making it prohibitively expensive to perform useful calculations on the hashed value. The hashing operation itself is considered expensive as well in the context of preselection, therefore QHDS preselection processes have to be defined in terms of only the helper data w , rather than the complete public data.

2.2.8 Noise Model

Though much of this thesis applies to noisy systems in general, some analysis requires further specification of the noise model. In these cases, all variable spaces are taken to be \mathbb{R}^D , for some $D \in \mathbb{N}$. It is generally assumed that the data have been preprocessed such that each dimensional component is independent of the others. In this section, the data are considered one-dimensional, but due to the independence between dimensional components, all equations presented here can be readily extended to multidimensional data.

The random variables representing biometric measurements, X and Y , depend on hidden variables representing the true, noiseless biometric, Z , and the measurement noises N_X and N_Y . These variables are never used directly, because they are impossible to determine. However, they impose the correlation between X and Y :

$$\begin{aligned} X &= Z + N_X \\ Y &= Z + N_Y. \end{aligned} \tag{2.21}$$

Introducing $N = N_Y - N_X$, Equation (2.21) gives

$$Y = X + N. \tag{2.22}$$

In this thesis, the distribution of N is usually considered symmetrically fading around 0 and independent of X and Y . This consequently imposes similar conditions on N_X and N_Y .

The use of preprocessing justifies taking the expectation values of all random variables equal to zero. The standard deviations are denoted by σ subscripted with the random variable to which it pertains. σ_X is calculated from σ_Z and σ_{N_X} :

$$\begin{aligned} \sigma_X^2 &= \mathbb{E} [(Z + N_X - \mathbb{E}[Z + N_X])^2] \\ &= \mathbb{E} [(Z + N_X)^2] \\ &= \mathbb{E} [Z^2] + 2\mathbb{E} [ZN_X] + \mathbb{E} [N_X^2] \\ &= \sigma_Z^2 + \sigma_{N_X}^2. \end{aligned} \tag{2.23}$$

Similarly

$$\sigma_Y^2 = \sigma_Z^2 + \sigma_{N_Y}^2.$$

The correlation coefficient ρ is defined as

$$\begin{aligned} \rho &= \frac{\mathbb{E} [(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]}{\sigma_X \sigma_Y} \\ &= \frac{\mathbb{E} [XY]}{\sigma_X \sigma_Y} \\ &= \frac{\mathbb{E} [(Z + N_X)(Z + N_Y)]}{\sigma_X \sigma_Y} \\ &= \frac{\mathbb{E} [Z^2] + \mathbb{E} [ZN_X] + \mathbb{E} [ZN_Y] + \mathbb{E} [N_X N_Y]}{\sigma_X \sigma_Y} \\ &= \frac{\sigma_Z^2}{\sigma_X \sigma_Y}. \end{aligned} \tag{2.24}$$

When concrete distributions are needed for the analysis in Chapter 4, Gaussian distributions will be used for Z , N_X and N_Y . This results in the following joint Gaussian distribution for X and Y :

$$\rho_{X,Y}(x,y) = \frac{1}{2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2}} \exp\left[-\frac{1}{2}\left(\frac{x^2}{\sigma_X^2} + \frac{y^2}{\sigma_Y^2} + \rho\frac{x}{\sigma_X}\frac{y}{\sigma_Y}\right)\right]. \quad (2.25)$$

The accompanying conditional distributions are Gaussian as well:

$$\rho_{Y|X}(y|x) = \frac{\phi\left(\frac{y - \frac{\sigma_Y}{\sigma_X}\rho x}{\sigma_Y\sqrt{1-\rho^2}}\right)}{\sigma_Y\sqrt{1-\rho^2}}; \quad \rho_{X|Y}(x|y) = \frac{\phi\left(\frac{x - \frac{\sigma_X}{\sigma_Y}\rho y}{\sigma_X\sqrt{1-\rho^2}}\right)}{\sigma_X\sqrt{1-\rho^2}}. \quad (2.26)$$

This allows for the conditional expectation value to be defined as

$$\mathbb{E}[Y|X=x] = \frac{\sigma_Y}{\sigma_X}\rho x = \lambda x; \quad \mathbb{E}[X|Y=y] = \frac{\sigma_X}{\sigma_Y}\rho y = \bar{\lambda}y. \quad (2.27)$$

Two special cases can be distinguished:

1. When the noise for enrollment measurements is equal to the noise for verification measurements, i.e., $\sigma_{N_X} = \sigma_{N_Y}$, then $\lambda = \bar{\lambda} = \rho$.
2. When the noise for enrollment measurements is much smaller than the noise for verification measurements, i.e., $\sigma_{N_X} \ll \sigma_{N_Y}$, then $\frac{\sigma_Y}{\sigma_X} \approx \frac{\sigma_Y}{\sigma_Z}$, which means that $\lambda \approx 1$ and $\bar{\lambda} \approx \frac{\sigma_X}{\sigma_Y}$.

2.2.9 Zero Leakage

It may be desirable for the helper data function not to reveal any information about the secrets. This property is called *zero leakage*, and is satisfied when

$$\forall (s,w) \in \mathcal{S} \times \mathcal{W}: \xi_{S|W}(s|w) = \xi_S(s). \quad (2.28)$$

Since $\xi_{S,W}(s,w) = \xi_{S|W}(s|w)\xi_W(w) = \xi_S(s)\xi_W(w)$, the zero leakage property is equivalent to independence of S and W .

2.2.10 Partitioning Template System (PTS)

In this section the quantization and helper data functions of a concrete QHDS will be defined. The properties of the corresponding matcher will be studied in Chapter 3. The dimensional components of \mathcal{X} are treated independently and both the secret and the helper data are derived for each component separately. In the following, all variables are treated as one-dimensional, to avoid the need for subscripts.

Definition 12. *The one-dimensional PTS is a QHDS in which the biometric sample space is $\mathcal{X} \subseteq \mathbb{R}$, the secret space is $\mathcal{S} = \mathbb{N}_n$, where $n \geq 2$, and the helper data space is $\mathcal{W} = [0,1)$. The quantization function is defined as*

$$\tilde{s}(x) = \lfloor n\Xi_X(x) \rfloor \quad (2.29)$$

and the helper data function as

$$\tilde{w}(x) = n\Xi_X(x) - \tilde{s}(x). \quad (2.30)$$

This definition is the continuum limit (i.e., letting the number of subintervals approach infinity) of the system presented in [17].

Definition 12 ensures that S and W are uniformly distributed and satisfy the zero leakage requirement, as evidenced by Lemma 1.

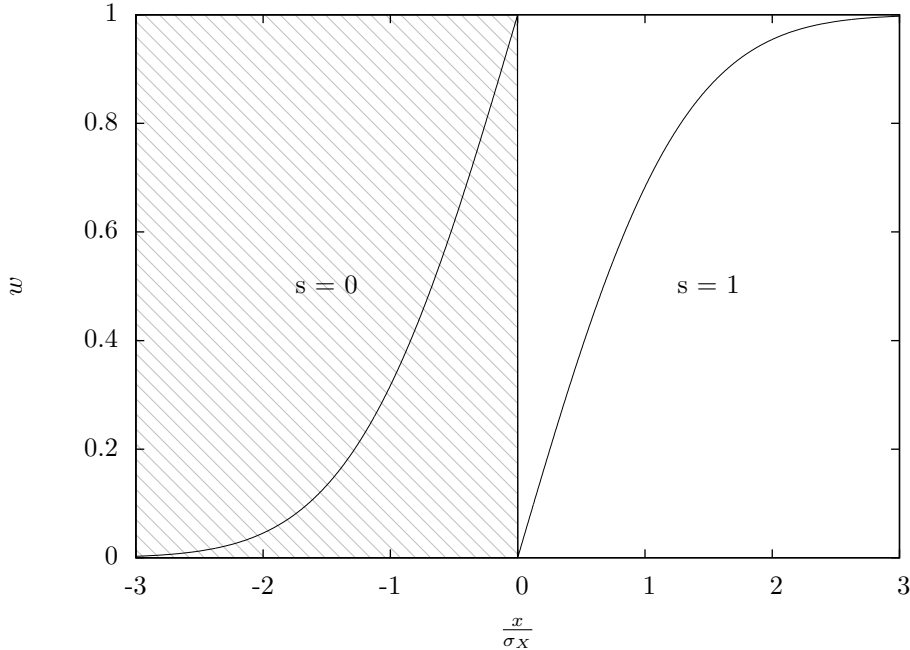


Figure 2.4: Helper data w as a function of biometric sample x for $n = 2$.

Lemma 1. Let $U = P_X(X)$, where X is any continuous random variable such that P_X does not contain any discontinuities. Then the continuous random variable U is uniformly distributed on $[0, 1]$.

Proof.

$$0 \leq P_X(x) \leq 1 \quad (2.31)$$

and

$$\begin{aligned} \rho_U(u(x))d(u(x)) &= \rho_X(x)dx \\ \Leftrightarrow \\ \rho_U(u(x)) &= \frac{\rho_X(x)}{\frac{d(u(x))}{dx}} = \frac{\rho_X(x)}{\rho_X(x)} = 1. \end{aligned} \quad (2.32)$$

□

The number of secrets, n , is determined from the underlying noise model and desired system characteristics. This will be discussed in more detail in Chapter 4.

It is conjectured that this system yields the lowest error probabilities possible for n uniformly distributed secrets under the noise model of Section 2.2.8.

Definition 12 allows for reconstruction of x from given s and w , as $\Xi_X(x) = \frac{s+w}{n}$.

An example in which $n = 2$ using the Gaussian noise model is displayed in Figure 2.4.

Preselection

Given a biometric verification sample y , is it possible to determine a region on \mathcal{X} , such that with high likelihood, the corresponding enrollment sample, x , lies in this region. A preselection can be

made by mapping this region to \mathcal{W} and selecting only database entries whose helper data are a member of this mapped region.

Chapter 3

Maximum Likelihood Reproducer

In this chapter, the maximum likelihood reproducer, Definition 11, is subjected to further analysis. Several alternative forms and the conditions under which they hold are introduced in Section 3.1. The consequences of satisfying the zero leakage requirement (Section 2.2.9) will be detailed in Section 3.2. Section 3.3 demonstrates the construction of a general zero leakage scheme. Finally the results of these sections are combined in Section 3.4 to yield the optimal reproducer for the Partitioning Template System (Definition 12).

3.1 A More Usable Form

To rewrite the approximate maximum likelihood reproducer, first, several ways of partitioning the \mathcal{X} space are introduced.

Definition 13. *The \mathcal{X} space can be partitioned according to \tilde{s} , \tilde{w} or both:*

$$\mathcal{X}_s = \{x \in \mathcal{X} \mid \tilde{s}(x) = s\} \quad s \in \mathcal{S} \quad (3.1)$$

$$\mathcal{X}_w = \{x \in \mathcal{X} \mid \tilde{w}(x) = w\} \quad w \in \mathcal{W} \quad (3.2)$$

$$\mathcal{X}_{s,w} = \{x \in \mathcal{X} \mid \tilde{s}(x) = s \wedge \tilde{w}(x) = w\} \quad (s, w) \in \mathcal{S} \times \mathcal{W}. \quad (3.3)$$

Furthermore, if \mathcal{X}_w does not contain a continuum, its members are called the set of sibling points.

Next, the following two lemmas are needed.

Lemma 2.

$$\xi_{S|X,Y}(s|x, y) = \rho_{S|X}(s|x) = \delta_{\tilde{s}(x), s}. \quad (3.4)$$

Proof. The first equality holds because S is a function of X and only of X . Therefore, the dependence on Y is completely contained in the dependence on X , and the conditional distribution is known exactly. The second holds due to the discrete nature of \mathcal{S} \square

Lemma 3.

$$\xi_{W|X,Y,S}(w|x, y, s) = \rho_{W|X}(w|x). \quad (3.5)$$

Proof. The equality holds because W is a function of X and only of X . Therefore, the dependence on Y and on S is completely contained in the dependence on X . \square

Using Definition 13 and Lemmas 2 and 3, the approximate maximum likelihood reproducer can be rewritten to a more usable form.

Lemma 4.

$$\text{Rep}_{\text{ML}}(y, w) = \arg \max_{s \in \mathcal{S}} \int_{\mathcal{X}_s} dx \xi_{X,Y}(x, y) \rho_{W|X}(w|x). \quad (3.6)$$

Proof. Starting from Definition 11:

$$\begin{aligned} \text{Rep}_{\text{ML}}(y, w) &= \arg \max_{s \in \mathcal{S}} \xi_{S|Y,W}(s|y, w) \\ &= \arg \max_{s \in \mathcal{S}} \frac{\xi_{Y,S,W}(y, s, w)}{\xi_{Y,W}(y, w)} \quad \boxed{\xi_{Y,W}(y, w) \text{ is independent of } s.} \\ &= \arg \max_{s \in \mathcal{S}} \xi_{Y,S,W}(y, s, w) \quad \boxed{\text{'Demarginalized' distribution.}} \\ &= \arg \max_{s \in \mathcal{S}} \int_{\mathcal{X}} dx \xi_{X,Y,S,W}(x, y, s, w) \quad (3.7) \\ &= \arg \max_{s \in \mathcal{S}} \int_{\mathcal{X}} dx \xi_{X,Y}(x, y) \xi_{S|X,Y}(s|x, y) \xi_{W|X,Y,S}(w|x, y, s) \\ &= \arg \max_{s \in \mathcal{S}} \int_{\mathcal{X}} dx \xi_{X,Y}(x, y) \delta_{\bar{s}(x),s} \rho_{W|X}(w|x) \\ &= \arg \max_{s \in \mathcal{S}} \int_{\mathcal{X}_s} dx \xi_{X,Y}(x, y) \rho_{W|X}(w|x). \end{aligned}$$

□

Using Lemma 4, the maximum likelihood reproducer can be expressed in terms of $\xi_{X,Y}(x, y)$ and $\rho_{W|X}(w|x)$. The first of these, $\xi_{X,Y}(x, y)$ is the empirical distribution on which the architecture of the system is based, and which can not be specified more concretely at this point in the analysis. The other, $\rho_{W|X}(w|x)$, allows for further analysis, which is discussed in Sections 3.1.1 and 3.1.2 for discrete and continuous \mathcal{W} respectively.

3.1.1 Discrete \mathcal{W}

If \mathcal{W} is a discrete space, $\rho_{W|X}$ can also be expressed as a Kronecker delta.

Lemma 5. *Let \mathcal{W} be discrete, then*

$$\rho_{W|X}(w|x) = \delta_{\bar{w}(x),w}. \quad (3.8)$$

Proof. The equality holds due to W being a function of X and the discrete nature of \mathcal{W} . □

Lemmas 4 and 5 can be combined into the following theorem.

Theorem 1. *Let \mathcal{W} be discrete, then*

$$\text{Rep}_{\text{ML}}(y, w) = \arg \max_{s \in \mathcal{S}} \int_{\mathcal{X}_{s,w}} dx \xi_{X,Y}(x, y). \quad (3.9)$$

Proof. Starting from Lemma 4:

$$\begin{aligned} \text{Rep}_{\text{ML}}(y, w) &= \arg \max_{s \in \mathcal{S}} \int_{\mathcal{X}_s} dx \xi_{X,Y}(x, y) \rho_{W|X}(w|x) \\ &= \arg \max_{s \in \mathcal{S}} \int_{\mathcal{X}_s} dx \xi_{X,Y}(x, y) \delta_{\bar{w}(x),w} \\ &= \arg \max_{s \in \mathcal{S}} \int_{\mathcal{X}_{s,w}} dx \xi_{X,Y}(x, y). \end{aligned}$$

□

Due to the fact that both \mathcal{S} and \mathcal{W} are discrete, the integration cannot be eliminated in this case.

3.1.2 Continuous \mathcal{W}

In case \mathcal{W} is continuous $\rho_{\mathcal{W}|X}$ can be expressed as a Dirac delta, under the condition that \mathcal{X} and \mathcal{W} are of the same dimensionality, and the Jacobian determinant of \tilde{w} is non-zero for every $x \in \mathcal{X}$. Let J be the Jacobian matrix of \tilde{w} , i.e., $J_{\alpha\beta}(x) = \frac{\partial \tilde{w}_\alpha(x)}{\partial x_\beta}$ and $\det J$ denote its determinant, then:

$$\rho_{\mathcal{W}|X}(w|x) = \delta(w - \tilde{w}(x)). \quad (3.10)$$

These requirements ensure that $\mathcal{X}_{s,w}$ is a discrete set of points.

Theorem 2. *Let \mathcal{W} be continuous, $\dim \mathcal{X} = \dim \mathcal{W}$ and $\det J(x) \neq 0$ for all $x \in \mathcal{X}$,*

$$\text{Rep}_{\text{ML}}(y, w) = \arg \max_{s \in \mathcal{S}} \sum_{x \in \mathcal{X}_{s,w}} \frac{\xi_{X,Y}(x, y)}{|\det J(x)|}. \quad (3.11)$$

Proof. Starting from Lemma 4:

$$\begin{aligned} \text{Rep}_{\text{ML}}(y, w) &= \arg \max_{s \in \mathcal{S}} \int_{\mathcal{X}_s} dx \xi_{X,Y}(x, y) \rho_{\mathcal{W}|X}(w|x) \\ &= \arg \max_{s \in \mathcal{S}} \int_{\mathcal{X}_s} dx \xi_{X,Y}(x, y) \delta(\tilde{w}(x) - w) \quad \boxed{\text{Substitution of variables: } dx = \frac{dw}{|\det J|}.} \\ &= \arg \max_{s \in \mathcal{S}} \int_{\mathcal{W}} dw' \frac{\xi_{X,Y}(\tilde{x}_s(w'), y)}{|\det J(\tilde{x}_s(w'))|} \delta(w' - w) \\ &= \arg \max_{s \in \mathcal{S}} \sum_{x \in \mathcal{X}_{s,w}} \frac{\xi_{X,Y}(x, y)}{|\det J(x)|}. \end{aligned} \quad (3.12)$$

□

If $\mathcal{X}_{s,w}$ contains exactly one point for every s and w , i.e., $\mathcal{X}_{s,w} = \{\tilde{x}_s(w)\}$, the necessary and sufficient condition for which is that \tilde{w} be fully invertible on each \mathcal{X}_s , then the sum in Theorem 2 can be eliminated.

Theorem 3. *Let the requirements of Theorem 2 be satisfied and $\forall (s, w) \in \mathcal{S} \times \mathcal{W}$: $\mathcal{X}_{s,w} = \{\tilde{x}_s(w)\}$, then*

$$\text{Rep}_{\text{ML}}(y, w) = \tilde{s} \left[\arg \max_{x \in \mathcal{X}_w} \frac{\xi_{X,Y}(x, y)}{|\det J(x)|} \right]. \quad (3.13)$$

Proof. Starting from Theorem 2:

$$\begin{aligned} \text{Rep}_{\text{ML}}(y, w) &= \arg \max_{s \in \mathcal{S}} \sum_{x \in \mathcal{X}_{s,w}} \frac{\xi_{X,Y}(x, y)}{|\det J(x)|} \\ &= \arg \max_{s \in \mathcal{S}} \frac{\xi_{X,Y}(\tilde{x}_s(w), y)}{|\det J(\tilde{x}_s(w))|} \\ &= \tilde{s} \left[\arg \max_{x \in \mathcal{X}_w} \frac{\xi_{X,Y}(x, y)}{|\det J(x)|} \right]. \end{aligned} \quad (3.14)$$

□

Theorem 3 implies that, as long as \tilde{s} is such that \tilde{w} is piecewise invertible, the choice of \tilde{s} is of no influence to the system's properties. This is because in this case the maximum likelihood reproducer operates like a Secure Sketch, reconstructing the enrolled X value, from which then the secret is extracted.

It may be desirable for each of the dimensional components of X to have its helper data determined separately and independent from the other dimensions. When the dimensionalities of \mathcal{X} and \mathcal{W} are equal, this results in a diagonal Jacobian matrix and

$$|\det J(x)| = \left| \prod_{d \in \mathbb{N}_{\dim \mathcal{X}}} \frac{d\tilde{w}_d(x_d)}{dx_d} \right|, \quad (3.15)$$

where \tilde{w}_d is a function of x_d only.

3.2 Zero Leakage

Under the zero leakage constraint, $\xi_{S,W}(s, w) = \xi_S(s)\xi_W(w)$, the maximum likelihood reproducer can be rewritten without using integration over \mathcal{X} .

Lemma 6. *When the zero leakage requirement holds,*

$$\text{Rep}_{\text{ML}}(y, w) = \arg \max_{s \in \mathcal{S}} \xi_{Y|S,W}(y|s, w) \xi_S(s). \quad (3.16)$$

Proof. Starting from Definition 11:

$$\begin{aligned} \text{Rep}_{\text{ML}}(y, w) &= \arg \max_{s \in \mathcal{S}} \xi_{S|Y,W}(s|y, w) \\ &= \arg \max_{s \in \mathcal{S}} \frac{\xi_{Y,S,W}(y, s, w)}{\xi_{Y,W}(y, w)} \quad \boxed{\xi_{Y,W}(y, w) \text{ is independent of } S.} \\ &= \arg \max_{s \in \mathcal{S}} \xi_{Y,S,W}(y, s, w) \\ &= \arg \max_{s \in \mathcal{S}} \xi_{S,W}(s, w) \xi_{Y|S,W}(y|s, w) \\ &= \arg \max_{s \in \mathcal{S}} \xi_S(s) \xi_W(w) \xi_{Y|S,W}(y|s, w) \quad \boxed{\xi_W(w) \text{ is independent of } S.} \\ &= \arg \max_{s \in \mathcal{S}} \xi_{Y|S,W}(y|s, w) \xi_S(s). \end{aligned} \quad (3.17)$$

□

Lemma 6 is not quite analogous to Lemma 4, because $\xi_{Y|S,W}$ has not been eliminated. In general, integration is needed to eliminate this factor, but if \tilde{w} is fully invertible on each \mathcal{X}_s , then $\xi_{Y|S,W}(y|s, w) = \xi_{Y|X}(y|\tilde{x}_s(w))$.

Theorem 4. *If the zero leakage requirement holds and \tilde{w} is fully invertible on each \mathcal{X}_s ,*

$$\text{Rep}_{\text{ML}}(y, w) = \tilde{s} \left[\arg \max_{x \in \mathcal{X}_w} \xi_{X,Y}(x, y) \frac{\xi_S(\tilde{s}(x))}{\xi_X(x)} \right]. \quad (3.18)$$

Proof. Starting from Lemma 6:

$$\begin{aligned} \text{Rep}_{\text{ML}}(y, w) &= \arg \max_{s \in \mathcal{S}} \xi_{Y|X}(y|\tilde{x}_s(w)) \xi_S(s) \quad \boxed{\text{Using } \xi_{Y|S,W}(y|s, w) = \xi_{Y|X}(y|\tilde{x}_s(w)).} \\ &= \tilde{s} \left[\arg \max_{x \in \mathcal{X}_w} \xi_{Y|X}(y|x) \xi_S(\tilde{s}(x)) \right] \\ &= \tilde{s} \left[\arg \max_{x \in \mathcal{X}_w} \xi_{X,Y}(x, y) \frac{\xi_S(\tilde{s}(x))}{\xi_X(x)} \right]. \end{aligned}$$

□

Theorems 3 and 4 together suggest the following theorem.

Theorem 5. *If the zero leakage requirement holds and \tilde{w} is fully invertible on each \mathcal{X}_s ,*

$$|\det J(x)| = \frac{\xi_X(x)}{\xi_S(\tilde{s}(x))\xi_W(\tilde{w}(x))}. \quad (3.19)$$

Proof. Analogous to the operand in Lemma 4 and Theorem 2:

$$\begin{aligned} \xi_{Y,S,W}(y, s, w) &= \int_{\mathcal{X}} dx \xi_{X,Y,S,W}(x, y, s, w) \\ &= \int_{\mathcal{X}} dx \xi_{X,Y}(x, y) \xi_{S|X,Y}(s|x, y) \xi_{W|X,Y,S}(w|x, y, s) \\ &= \int_{\mathcal{X}} dx \xi_{X,Y}(x, y) \delta_{\tilde{s}(x), s} \delta(\tilde{w}(x) - w) \\ &= \int_{\mathcal{X}_s} dx \xi_{X,Y}(x, y) \delta(\tilde{w}(x) - w) \\ &= \frac{\xi_{X,Y}(x, y)}{|\det J(\tilde{x}_s(w))|}. \end{aligned} \quad (3.20)$$

Analogous to the operand in Lemma 6:

$$\begin{aligned} \xi_{Y,S,W}(y, s, w) &= \xi_{S,W}(s, w) \xi_{Y|S,W}(y|s, w) \\ &= \xi_S(s) \xi_W(w) \xi_{Y|X}(y|\tilde{x}_s(w)) \\ &= \frac{\xi_S(s) \xi_W(w)}{\xi_X(\tilde{x}_s(w))} \xi_{X,Y}(\tilde{x}_s(w), y). \end{aligned} \quad (3.21)$$

Because $\tilde{x}_s(w)$ spans the entire \mathcal{X} space, Equations (3.20) and (3.21) can be combined:

$$|\det J(\tilde{x}_s(w))| = \frac{\xi_X(\tilde{x}_s(w))}{\xi_S(s) \xi_W(w)} \quad \forall s, w \in \mathcal{S} \times \mathcal{W} \quad (3.22)$$

\Leftrightarrow

$$|\det J(x)| = \frac{\xi_X(x)}{\xi_S(\tilde{s}(x)) \xi_W(\tilde{w}(x))} \quad \forall x \in \mathcal{X}. \quad (3.23)$$

□

Theorem 6. *Let the dimensional components of X be independent and both the secret and helper data be determined for each dimensional component independently, i.e., $J(x)$ is diagonal and $\tilde{s}(x) = (\tilde{s}_1(x_1), \tilde{s}_2(x_2), \dots, \tilde{s}_{\dim \mathcal{X}}(x_{\dim \mathcal{X}}))$, then*

$$\frac{\xi_{X_d}(x_d)}{\xi_{S_d}(\tilde{s}_d(x_d)) \xi_{W_d}(\tilde{w}_d(x_d))} = \left| \frac{d\tilde{w}_d(x_d)}{dx_d} \right| \quad \forall d \in \mathbb{N}_{\dim \mathcal{X}}. \quad (3.24)$$

Proof. Starting from Theorem 5:

$$\begin{aligned} \frac{\xi_X(x)}{\xi_S(\tilde{s}(x)) \xi_W(\tilde{w}(x))} &= |\det J(x)| \\ &= \\ \prod_{d \in \mathbb{N}_D} \frac{\xi_{X_d}(x_d)}{\xi_{S_d}(\tilde{s}_d(x_d)) \xi_{W_d}(\tilde{w}_d(x_d))} &= \left| \prod_{d \in \mathbb{N}_D} \frac{d\tilde{w}_d(x_d)}{dx_d} \right| \\ &\Leftrightarrow \\ \frac{\xi_{X_d}(x_d)}{\xi_{S_d}(\tilde{s}_d(x_d)) \xi_{W_d}(\tilde{w}_d(x_d))} &= \left| \frac{d\tilde{w}_d(x_d)}{dx_d} \right| \quad \forall d \in \mathbb{N}_D. \end{aligned} \quad (3.25)$$

□

Choosing S and W to be uniformly distributed leads to

$$\xi_{X_d}(x_d)|\mathcal{S}||\mathcal{W}| = \left| \frac{d\tilde{w}_d(x_d)}{dx_d} \right| \quad \forall d \in \mathbb{N}_D, \quad (3.26)$$

where $|\mathcal{S}|$ is the cardinality of \mathcal{S} and $|\mathcal{W}|$ is the *volume* of \mathcal{W} . Equation (3.26) is the underlying relation which motivates the choice of the PTS of Section 2.2.10. The PTS is not the only possible system satisfying all requirements, but it is conjectured to be the optimal system for symmetrical fading noise, i.e., achieving the lowest error rates out of all possible systems.

3.3 General Zero Leakage Scheme

Theorem 6 suggests that zero leakage helper data can be derived for any given quantization function by integration of Equation (3.24). Because zero leakage requires that every helper data value occurs for every secret, the integration has to be performed on each of the partitions separately. Choosing $\mathcal{W} = [0, 1)$ and \tilde{w} to be both uniformly distributed and monotonically increasing, except for discontinuities at partition boundaries, gives

$$d\tilde{w}(x) = \frac{\xi_X(x)}{\xi_S(\tilde{s}(x))} dx, \quad (3.27)$$

where $\xi_S(s) > 0$ is required for all s . Integration of Equation (3.27) gives

$$\tilde{w}(x) = \frac{\int_{-\infty}^x dx' \xi_X(x') \delta_{\tilde{s}(x), \tilde{s}(x')}}{\xi_S(\tilde{s}(x))}. \quad (3.28)$$

This scheme is not necessarily optimal, as is demonstrated by the example in Figure 3.1.

Corollary 1. *If all secret partitions are contiguous and occur in increasing order with x , then Equation (3.28) can be written as:*

$$\tilde{w}(x) = \frac{\Xi_X(x) - \sum_{s=0}^{\tilde{s}(x)-1} \xi_S(\tilde{s}(x))}{\xi_S(\tilde{s}(x))}. \quad (3.29)$$

Corollary 2. *If Corollary 1 holds and S is uniformly distributed, then Equation (3.28) can be written as:*

$$\tilde{w}(x) = n\Xi_X(x) - \tilde{s}(x). \quad (3.30)$$

Corollaries 1 and 2 are published in [2].

The Partitioning Template System presented in Section 2.2.10 is a specific form of this general scheme in which all secrets carry equal weight and the secret partitions are contiguous, corresponding to Corollary 2. The general form was not presented earlier, because there is no guarantee of optimality, whereas the PTS is conjectured to be optimal under certain conditions.

3.4 Partitioning Template System

The maximum likelihood reproducer, and therefore the corresponding matcher, for the PTS can now easily be derived.

Theorem 7. *Under the noise model of Section 2.2.8, the maximum likelihood reproducer for the PTS is given by*

$$\text{Rep}_{\text{ML}}(y, w) = \tilde{s} \left[\arg \min_{x \in \mathcal{X}_w} |y - \lambda x| \right]. \quad (3.31)$$

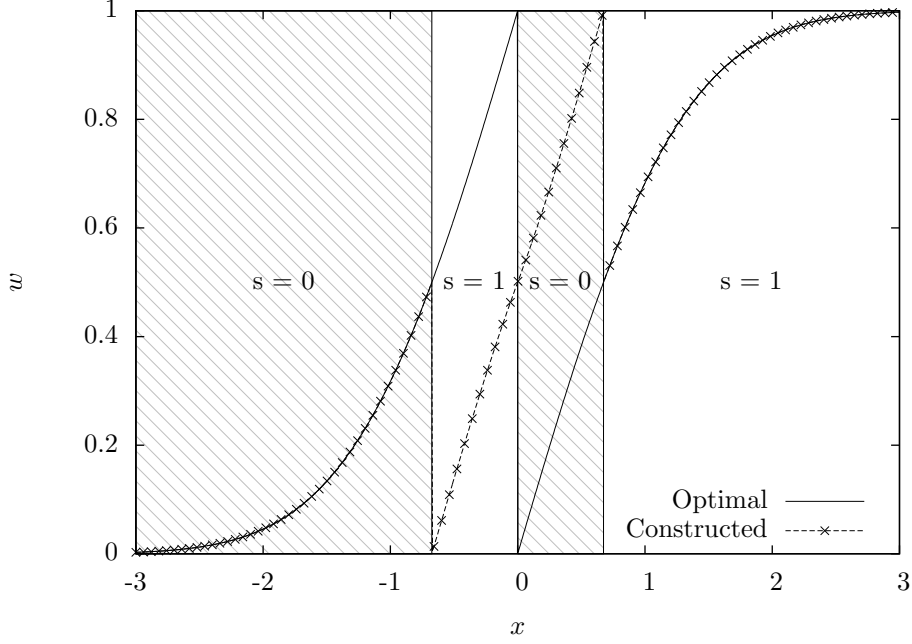


Figure 3.1: The helper data w as a function of the biometric sample x for both the constructed helper data function, Equation (3.28) and the optimal function for this case. $s = 0$ in the highlighted region and $s = 1$ otherwise.

Proof. Starting from Lemma 6:

$$\begin{aligned}
\text{Rep}_{\text{ML}}(y, w) &= \arg \max_{s \in \mathcal{S}} \xi_{Y|S,W}(y|s, w) \xi_S(s) && \boxed{\text{Using } \xi_{Y|S,W}(y|s, w) = \xi_{Y|X}(y|\tilde{x}_s(w)).} \\
&= \arg \max_{s \in \mathcal{S}} \xi_{Y|X}(y|\tilde{x}_s(w)) \xi_S(s) && \boxed{S \text{ is uniformly distributed.}} \\
&= \arg \max_{s \in \mathcal{S}} \xi_{Y|X}(y|\tilde{x}_s(w)) \\
&= \tilde{s} \left[\arg \max_{x \in \mathcal{X}_w} \xi_{Y|X}(y|x) \right] && (3.32) \\
&= \tilde{s} \left[\arg \max_{x \in \mathcal{X}_w} \xi_N(|y - \lambda x|) \right] && \boxed{N \text{ has a fading distribution (Section 2.2.8).}} \\
&= \tilde{s} \left[\arg \min_{x \in \mathcal{X}_w} |y - \lambda x| \right].
\end{aligned}$$

□

Theorem 7 states that the reproduced secret of a measured y will be that of the sibling point in \mathcal{X}_w that (after correction λ) is nearest to y . This means that the region on \mathcal{Y} that will match a given x is bounded from below by

$$\frac{\lambda}{2} [x + \max\{x' \in \mathcal{X}_{\tilde{w}(x)} | x' < x\}] = \frac{\lambda}{2} [x + \tilde{x}_{\tilde{s}(x)-1}(\tilde{w}(x))] \quad (3.33)$$

and from above by

$$\frac{\lambda}{2} [x + \min\{x' \in \mathcal{X}_{\tilde{w}(x)} | x' > x\}] = \frac{\lambda}{2} [x + \tilde{x}_{\tilde{s}(x)+1}(\tilde{w}(x))], \quad (3.34)$$

if these exist.

Chapter 4

Properties of Partitioning Template Systems

In this chapter, the properties of a PTS will be analyzed for one-dimensional data. First, the error characteristics will be investigated. This is followed by analysis of the running time improvements, for which numerical results obtained using the Gaussian noise model will be presented. The chapter will conclude with notes on how these results apply to multidimensional data.

4.1 Error Characteristics

For a QHDS, the error indicator random variable can be expressed as

$$E = \begin{cases} 0 & \text{if } \tilde{s}(x) = \text{Rep}(y, \tilde{w}(x)) \\ 1 & \text{otherwise} \end{cases}. \quad (4.1)$$

Then $\rho_{E|X,Y}(0|x, y) = \delta_{\tilde{s}(x), \text{Rep}(y, \tilde{w}(x))}$.

For the PTS with independent, symmetrically fading noise, this becomes

$$\rho_{E|X,Y}(0|x, y) = \begin{cases} 1 & \text{if } \tau_{\tilde{s}(x)}(\tilde{w}(x)) \leq y < \tau_{\tilde{s}(x)+1}(\tilde{w}(x)) \\ 0 & \text{otherwise} \end{cases}. \quad (4.2)$$

Then, for $n \geq 2$,

$$\begin{aligned} \rho_{E|X}(0|x) &= \int_{\tau_{\tilde{s}(x)}(\tilde{w}(x))}^{\tau_{\tilde{s}(x)+1}(\tilde{w}(x))} dy \rho_{Y|X}(y|x) \\ &= \int_{\tau_{\tilde{s}(x)}(\tilde{w}(x))}^{\tau_{\tilde{s}(x)+1}(\tilde{w}(x))} dy \rho_N(y - \lambda x) \\ &= P_N(\tau_{\tilde{s}(x)+1}(\tilde{w}(x)) - \lambda x) - P_N(\tau_{\tilde{s}(x)}(\tilde{w}(x)) - \lambda x) \\ &= \begin{cases} P_N\left(\frac{\lambda}{2} \left(P_X^{-1}\left(P_X(x) + \frac{1}{n}\right) - x\right)\right) & \text{if } P_X(x) < \frac{1}{n} \\ P_N\left(\frac{\lambda}{2} \left(P_X^{-1}\left(P_X(x) + \frac{t}{n}\right) - x\right)\right)\Big|_{t=-1}^1 & \text{if } \frac{1}{n} \leq P_X(x) < \frac{n-1}{n} \\ 1 - P_N\left(\frac{\lambda}{2} \left(P_X^{-1}\left(P_X(x) - \frac{1}{n}\right) - x\right)\right) & \text{if } \frac{n-1}{n} \leq P_X(x) \end{cases}. \end{aligned} \quad (4.3)$$

The expected error probability conditioned on X for the Gaussian noise model is displayed in Figure 4.1 for three different values of λ . This demonstrates the error probabilities decrease when

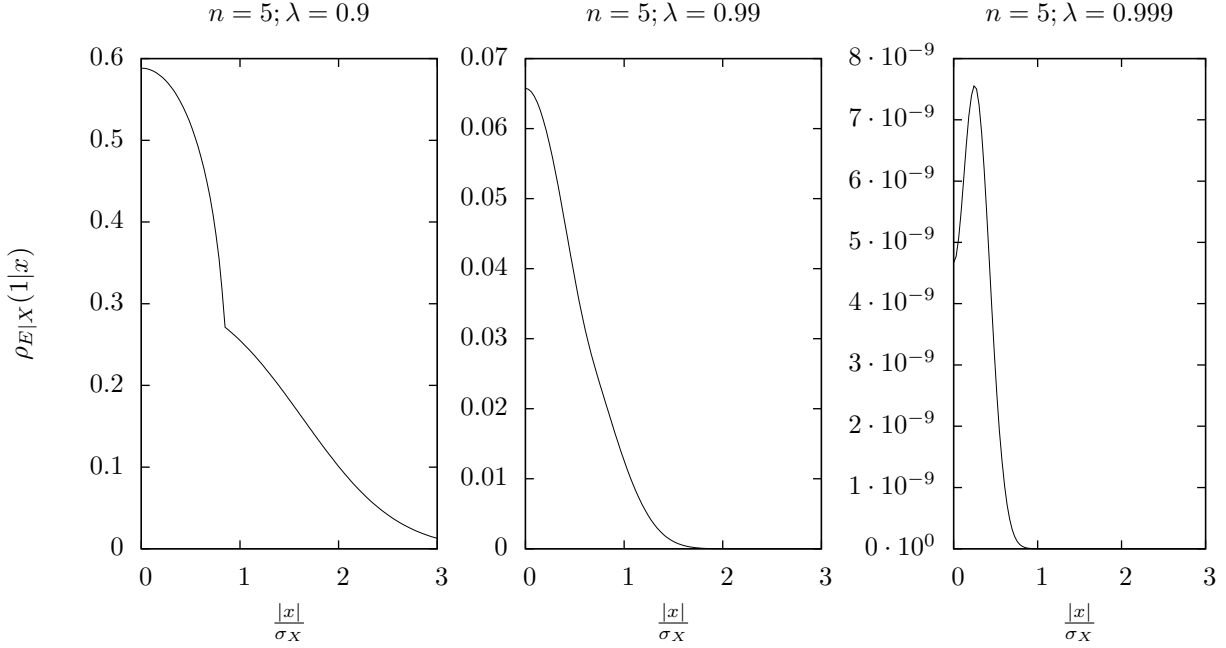


Figure 4.1: The expected error probability as a function of x using a Gaussian noise model and $n = 5$. From left to right, the values of λ are 0.9, 0.99 and 0.999.

λ approaches 1. When λ is small, the expected error probability shows a discontinuity in its derivative when transitioning between the three cases in Equation (4.3) as displayed for $\lambda = 0.9$. This gradually smooths out when λ increases, as displayed for $\lambda = 0.99$. Finally, when λ increases even further, the maximum expected error probability no longer occurs at $x = 0$, but near the point x such that $P_X^{-1}(P_X(x \pm \frac{1}{n})) = -x$. This behavior occurs generally for symmetrically fading distributions. While $\rho_{E|X}(1|x)$ is symmetrical about $x = 0$, only the positive range is displayed to show more detail.

The expected error probability is given by

$$\begin{aligned}
\rho_E(1) &= 1 - \rho_E(0) \\
&= 1 - \int_{-\infty}^{\infty} dx \rho_{E|X}(0|x) \rho_X(x) \\
&= 1 - \int_0^1 d(P_X(x)) \rho_{E|X}[0|P_X^{-1}(P_X(x))] \\
&= 1 - \int_0^1 du \rho_{E|X}(0|P_X^{-1}(u)) \\
&= 1 - \left\{ \int_0^{1-\frac{1}{n}} du P_N \left[\frac{\lambda}{2} \left(P_X^{-1} \left(u + \frac{1}{n} \right) - P_X^{-1}(u) \right) \right] + \int_{1-\frac{1}{n}}^1 du 1 \right. \\
&\quad \left. - \int_0^{\frac{1}{n}} du 0 - \int_{\frac{1}{n}}^1 du P_N \left[\frac{\lambda}{2} \left(P_X^{-1} \left(u - \frac{1}{n} \right) - P_X^{-1}(u) \right) \right] \right\}
\end{aligned} \tag{4.4}$$

(Continued)

$$\begin{aligned}
&= 1 - \left\{ \frac{1}{n} + \int_{\frac{1}{n}}^1 du \left[1 - P_N \left(\frac{\lambda}{2} \left(P_X^{-1} \left(u - \frac{1}{n} \right) - P_X^{-1}(u) \right) \right) \right] \right. \\
&\quad \left. - \int_{\frac{1}{n}}^1 du \left[\frac{\lambda}{2} \left(P_X^{-1} \left(u - \frac{1}{n} \right) - P_X^{-1}(u) \right) \right] \right\} \\
&= 2 \int_{\frac{1}{n}}^1 du P_N \left[\frac{\lambda}{2} \left(P_X^{-1} \left(u - \frac{1}{n} \right) - P_X^{-1}(u) \right) \right].
\end{aligned}$$

This integral cannot in general be evaluated and is not analytical in the case of the Gaussian noise model. To estimate the expected error probability it can be bounded from above using the maximum expected error probability. Let $\varepsilon_{\max} = \max_x \rho_{E|X}(1|x) = 1 - \min_x \rho_{E|X}(0|x)$.

If X is distributed symmetrically fading around zero, then a local extremum in the reconstruction error probability will occur for $X = 0$. Without proof, it is stated that this extremum will be the global maximum, unless X and Y are highly correlated and consequently the expected error probability is exceedingly small. This case can be disregarded, as it defeats the purpose of this thesis.

Then, for $n \geq 3$,

$$\begin{aligned}
\varepsilon_{\max} &= 1 - \left\{ P_N \left[\frac{\lambda}{2} \left(P_X^{-1} \left(\frac{1}{2} + \frac{1}{n} \right) \right) \right] - P_N \left[\frac{\lambda}{2} \left(P_X^{-1} \left(\frac{1}{2} - \frac{1}{n} \right) \right) \right] \right\} \\
&= 2P_N \left[\frac{\lambda}{2} P_X^{-1} \left(\frac{1}{2} - \frac{1}{n} \right) \right].
\end{aligned} \tag{4.5}$$

This relation between ε_{\max} , λ , and n can be rewritten as follows:

$$\begin{aligned}
\varepsilon_{\max} &= 2P_N \left[\frac{\lambda}{2} P_X^{-1} \left(\frac{1}{2} - \frac{1}{n} \right) \right] \\
&\Leftrightarrow \\
P_N^{-1} \left(\frac{\varepsilon}{2} \right) &= \frac{\lambda}{2} P_X^{-1} \left(\frac{1}{2} - \frac{1}{n} \right) \\
&\Leftrightarrow \\
P_X \left(\frac{2}{\lambda} P_N^{-1} \left(\frac{\varepsilon}{2} \right) \right) &= \frac{1}{2} - \frac{1}{n} \\
&\Leftrightarrow \\
\frac{1}{n} &= \frac{1}{2} - P_X \left(\frac{2}{\lambda} P_N^{-1} \left(\frac{\varepsilon}{2} \right) \right) \\
&\Leftrightarrow \\
n &= \frac{1}{\frac{1}{2} - P_X \left(\frac{2}{\lambda} P_N^{-1} \left(\frac{\varepsilon}{2} \right) \right)}.
\end{aligned} \tag{4.6}$$

This relation can be used to determine the maximum possible value for n , given a constraint on the maximum error probability. Several examples of this are given in Figures 4.2 and 4.3

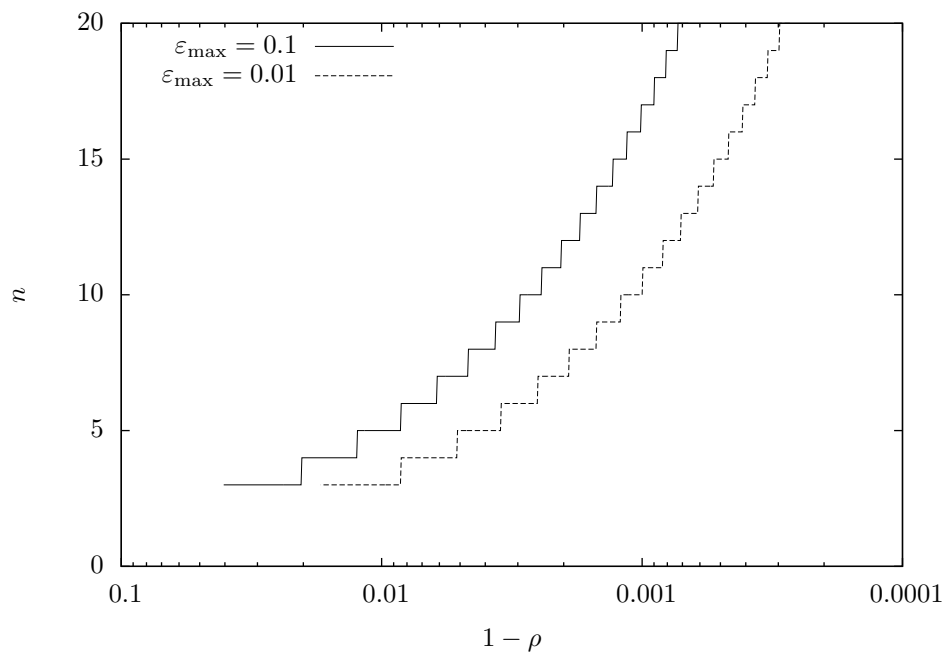


Figure 4.2: The maximum possible value for n versus $1 - \rho$ for two given values of ϵ_{\max} . A logarithmic scale is used for ρ approaching 1.

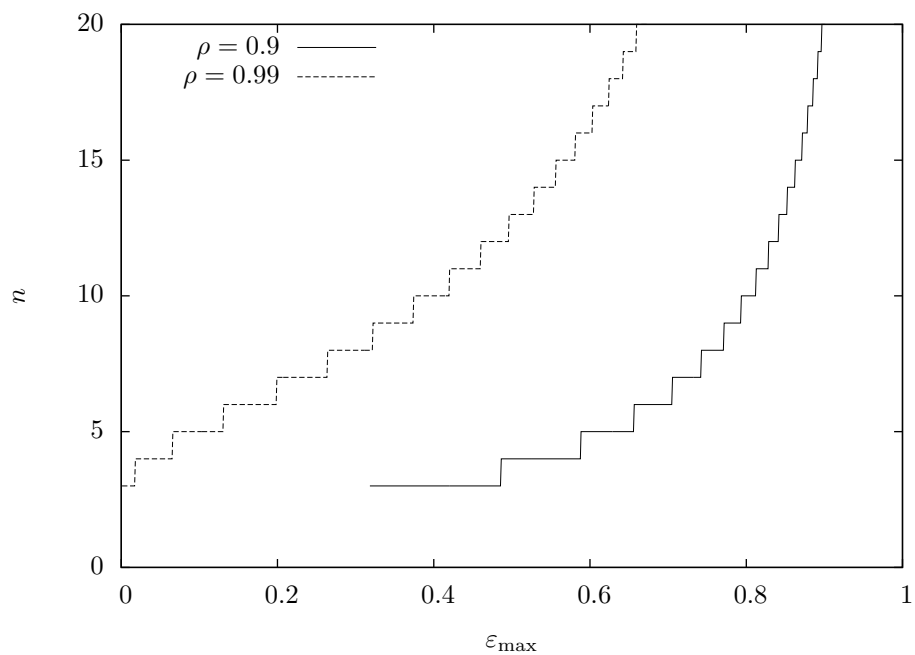


Figure 4.3: The maximum possible value for n versus ϵ_{\max} for two given values of ρ .

4.2 Running Time Improvement

The expected running time for the two stage priority selection identification process is (Equation (2.16))

$$\begin{aligned}
T_{2\text{CFM}}(v, N) &= \tau_{\text{Pre}} + \frac{1}{2} \varepsilon_{0,0} \omega_{\text{Pre}} N t_{\text{Auth}} \\
&\quad + \frac{1}{2} \varepsilon_{1,0} (1 + \omega_{\text{Pre}}) N t_{\text{Auth}} \\
&\quad + \varepsilon_{\text{Auth}} N t_{\text{Auth}}.
\end{aligned} \tag{4.7}$$

Ideally, the optimal priority selection can be made by selecting from \mathcal{W} those values that have the highest probability given y . However, this approach is not suitable for analysis. Instead, the values from \mathcal{X} that have the highest probability given y are chosen and mapped to \mathcal{W} , after which a selection can be made from the database. If the overall error rate is small enough, the difference between these two processes is negligible.

Since $\rho_{X|Y}(x|y)$ is symmetrically fading around λy , the optimal selection region in \mathcal{X} also symmetrical and can be described as $(\lambda y - \Delta, \lambda y + \Delta)$ for some $\Delta > 0$.

To minimize $T_{2\text{CFM}}$ its derivative with respect to Δ is determined. Before that the second and third term of Equation (4.7) are recombined using $\varepsilon_{0,0} + \varepsilon_{1,0} = 1 - \varepsilon_{\text{Auth}}$ to simplify the result:

$$\begin{aligned}
T_{2\text{CFM}}(v, N) &= \tau_{\text{Pre}} + \frac{1}{2} (1 - \varepsilon_{\text{Auth}}) (1 + \omega_{\text{Pre}}) N t_{\text{Auth}} \\
&\quad - \frac{1}{2} \varepsilon_{0,0} N t_{\text{Auth}} \\
&\quad + \varepsilon_{\text{Auth}} N t_{\text{Auth}}.
\end{aligned} \tag{4.8}$$

Only ω_{Pre} and $\varepsilon_{0,0}$ depend on Δ , so to optimize Δ , the following equation has to be solved:

$$\frac{\partial T_{2\text{CFM}}(y, N)}{\partial \Delta} = \left((1 - \varepsilon_{\text{Auth}}) \frac{d\omega_{\text{Pre}}}{d\Delta} - \frac{d\varepsilon_{0,0}}{d\Delta} \right) \frac{1}{2} N t_{\text{Auth}} = 0 \tag{4.9}$$

\Leftrightarrow

$$(1 - \varepsilon_{\text{Auth}}) \frac{d\omega_{\text{Pre}}}{d\Delta} = \frac{d\varepsilon_{0,0}}{d\Delta}. \tag{4.10}$$

Substituting the definitions of $\varepsilon_{0,0}$ and ω_{Auth} yields

$$\begin{aligned}
\frac{d\varepsilon_{0,0}(y)}{d\Delta} &= \frac{d}{d\Delta} \int_{\bar{\lambda}y - \Delta}^{\bar{\lambda}y + \Delta} dx \rho_{X|Y}(x|y) \\
&= \rho_{X|Y}(\bar{\lambda}y + \Delta|y) + \rho_{X|Y}(\bar{\lambda}y - \Delta|y)
\end{aligned} \tag{4.11}$$

and

$$\begin{aligned}
\frac{d\omega_{\text{Pre}}(y)}{d\Delta} &= n \frac{d}{d\Delta} (P(\bar{\lambda}y + \Delta) - P(\bar{\lambda}y - \Delta)) \\
&= n (\rho_X(\bar{\lambda}y + \Delta) + \rho_X(\bar{\lambda}y - \Delta)).
\end{aligned} \tag{4.12}$$

Furthermore,

$$\begin{aligned}
\varepsilon_{\text{Auth}} &= \int_{x_{\perp}(y)}^{x_{\top}(y)} dx \rho_{X|Y}(x|y) \\
&= \int_{x_{\perp}(y)}^{x_{\top}(y)} dx \frac{\phi\left(\frac{x - \bar{\lambda}y}{\sqrt{1 - \rho^2} \sigma_X}\right)}{\sqrt{1 - \rho^2} \sigma_X},
\end{aligned} \tag{4.13}$$

where $x_{\perp}(y)$ and $x_{\top}(y)$ are the solutions to

$$\frac{\lambda}{2} \left(x_{\perp} + P_X^{-1} \left(P_X(x_{\perp}) + \frac{1}{n} \right) \right) = y \quad (4.14)$$

$$\frac{\lambda}{2} \left(x_{\top} + P_X^{-1} \left(P_X(x_{\top}) - \frac{1}{n} \right) \right) = y. \quad (4.15)$$

No further analytical evaluation can be done without assumptions on the probability distributions.

4.2.1 Gaussian Distribution

When X and Y are jointly Gaussian distributed:

$$\begin{aligned} \varepsilon_{0,0}(y) &= \int_{\bar{\lambda}y-\Delta}^{\bar{\lambda}y+\Delta} dx \rho_{X|Y}(x|y) \\ &= P_{X|Y}(\bar{\lambda}y + \Delta|y) - P_{X|Y}(\bar{\lambda}y - \Delta|y) \\ &= \Phi \left(\frac{\Delta}{\sqrt{1-\rho^2}\sigma_X} \right) - \Phi \left(\frac{-\Delta}{\sqrt{1-\rho^2}\sigma_X} \right) \end{aligned} \quad (4.16)$$

and

$$\begin{aligned} \omega_{\text{Pre}}(y) &= n (P_X(\bar{\lambda}y + \Delta) - P_X(\bar{\lambda}y - \Delta)) \\ &= n (\Phi(\bar{\lambda}y + \Delta) - \Phi(\bar{\lambda}y - \Delta)), \end{aligned} \quad (4.17)$$

making Equations (4.11) and (4.12)

$$\begin{aligned} \frac{d\varepsilon_{0,0}(y)}{d\Delta} &= \frac{\phi \left(\frac{\Delta}{\sqrt{1-\rho^2}\sigma_X} \right)}{\sqrt{1-\rho^2}\sigma_X} + \frac{\phi \left(\frac{-\Delta}{\sqrt{1-\rho^2}\sigma_X} \right)}{\sqrt{1-\rho^2}\sigma_X} \\ &= \frac{2}{\sqrt{1-\rho^2}\sigma_X} \phi \left(\frac{\Delta}{\sqrt{1-\rho^2}\sigma_X} \right) \end{aligned} \quad (4.18)$$

and

$$\frac{d\omega_{\text{Pre}}(y)}{d\Delta} = n (\phi(\bar{\lambda}y + \Delta) + \phi(\bar{\lambda}y - \Delta)). \quad (4.19)$$

At this point it is not possible to obtain solutions analytically. Further evaluation has to be done numerically.

4.2.2 Numerical Results

At this point in the analysis, the equations may not be analytically soluble. Therefore, for the remainder of this chapter, numerical solutions to the above equations will be presented.

To obtain these results some standard parameters are used.

- $\sigma_{N_X} = \sigma_{N_Y}$, i.e., the first special case of the noise model of Section 2.2.8, meaning $\rho = \lambda = \bar{\lambda}$.
- Where appropriate, results will be displayed for two values of the correlation coefficient: a moderate value, $\rho = 0.99$ and a high value $\rho = 0.9999$.

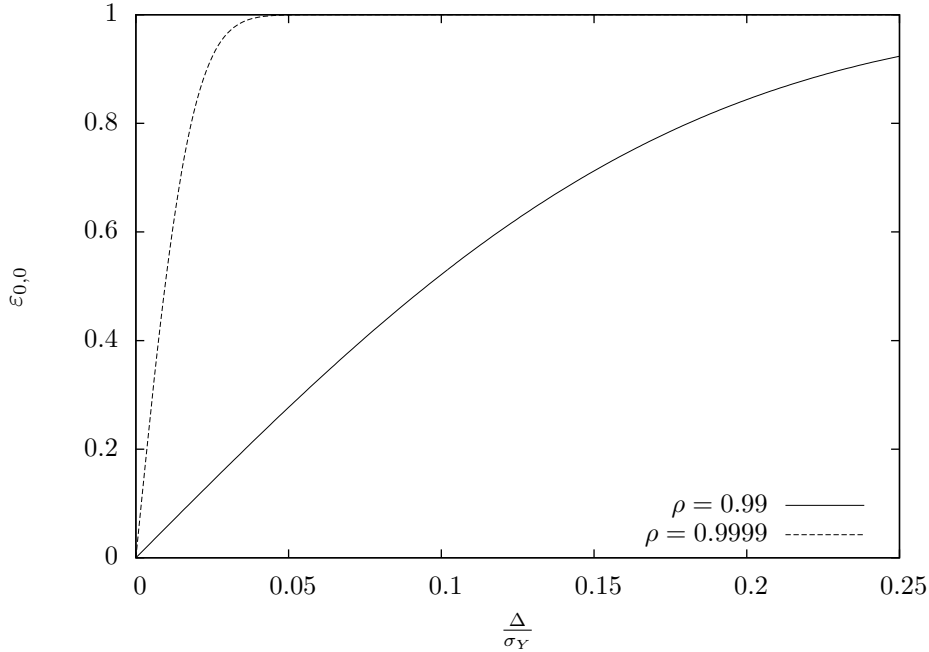


Figure 4.4: $\varepsilon_{0,0}$ as a function of Δ using the standard parameters.

- For the number of quantization intervals, again, two values will be used: $n = 3$ and $n = 5$.
- Similarly, the example values $y = 0$ and $|y| = \sigma_Y$ will be used.

In the figures presented below the results are symmetrical in y about 0. To improve readability, only the positive range will be displayed and labeled with $|y|$.

First of all, the behaviors of $\varepsilon_{0,0}$ and ω are presented as a function of Δ in Figures 4.4 and 4.5 respectively.

In Figure 4.4, when $\Delta = 0$, the priority selection region is empty, and it is, therefore, impossible to make the correct identification during this step, i.e., $\varepsilon_{0,0} = 0$. As Δ increases, so will $\varepsilon_{0,0}$ as it asymptotically approaches 1. For the greater correlation coefficient, $\varepsilon_{0,0}$ will increase more strongly with Δ than for the smaller one, as a large correlation coefficient indicates a high probability of finding the correct entry within the priority selection region.

Because both standard correlation coefficients are quite close to 1, the behavior of ω as a function of Δ , as displayed in Figure 4.5 is nearly linear. Similar to $\varepsilon_{0,0}$ when $\Delta = 0$, the priority selection region is empty and no work will be done during this phase. For $y = 0$ and $n > 2$, when Δ reaches the threshold point of $x = 0$, then $\omega = 1$. Beyond this point, the numerical results no longer hold, as they no longer conform to the defined system. For $y \neq 0$, this point occurs just before ω reaches 1, due to the asymmetry in threshold points, while the priority selection region remains symmetrical.

To inspect the actual running time gains, T_{2CFM} can be compared to T_{Auth} . Since all terms except τ_{Pre} scale linearly with T_{Auth} , the quantity $\frac{T_{2CFM} - \tau_{Pre}}{T_{Auth}}$ is suitable for making the comparison.

Figure 4.6 displays the normalized running time as a function of Δ . The solutions for all possible values of the standard parameters are shown for completeness. Most importantly, this figure shows the existence of the running time minima and, for greater correlation coefficients, the steepness of the dependence of the running time on Δ . The dependence of the optimal expected running time on the other parameters will be studied in more detail below.

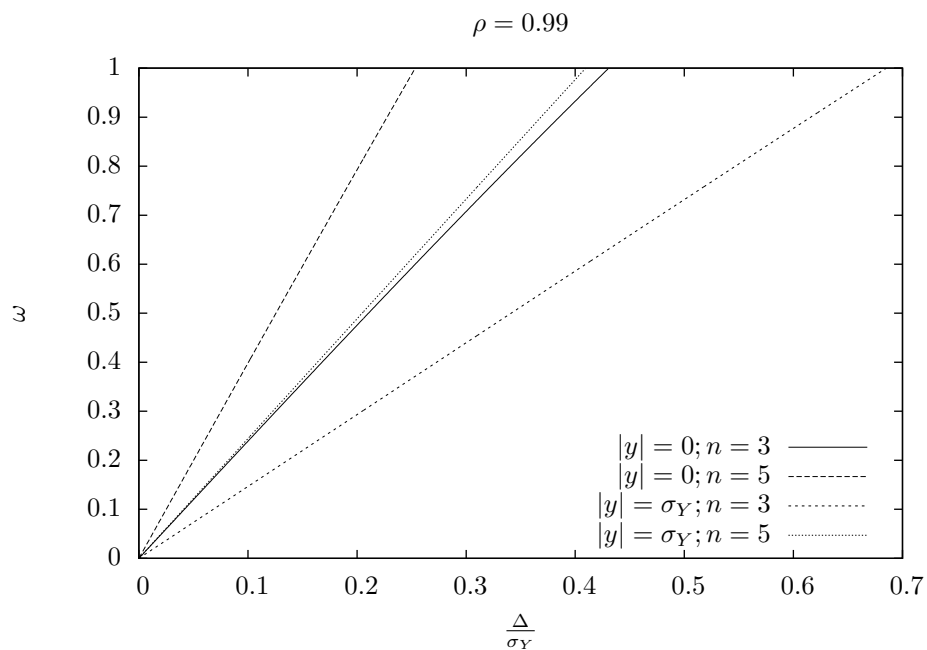


Figure 4.5: ω as a function of Δ using the standard parameters.

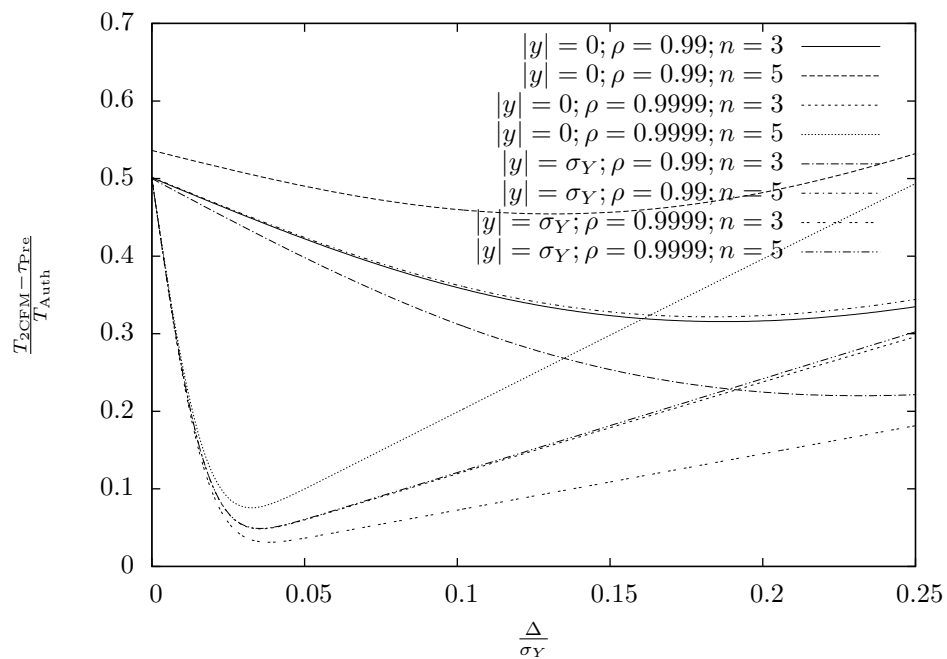


Figure 4.6: The normalized running time as a function of Δ , using the standard parameters.

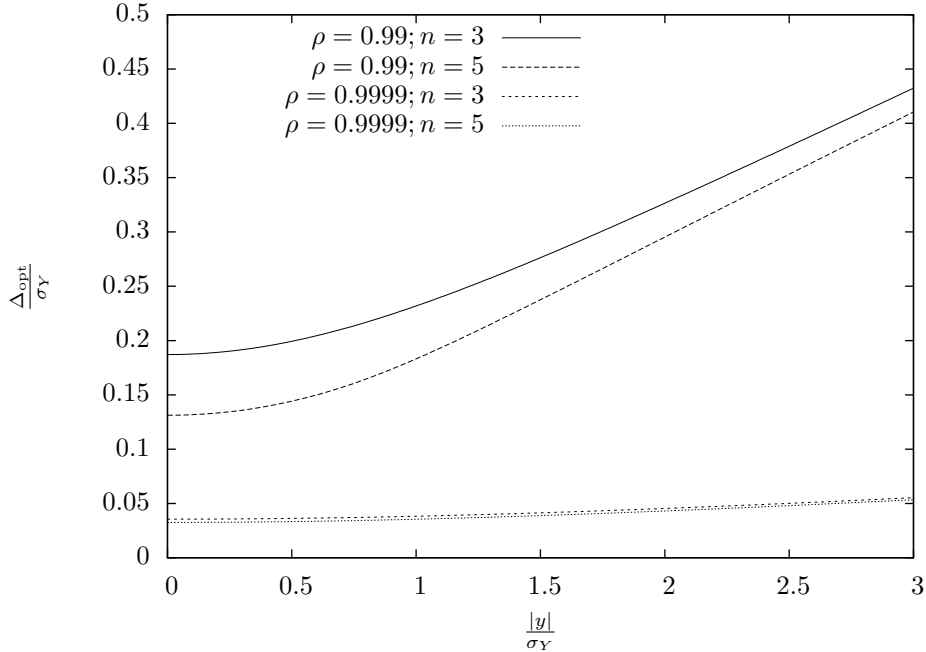


Figure 4.7: The value of Δ yielding the minimal T_{2CFM} as a function of y using the standard parameters.

Equations (4.13), (4.18) and (4.19) can be used to numerically solve Equation (4.9) and find the value of Δ yielding the optimal running time T_{2CFM} . The solutions using the standard parameters are displayed in Figure 4.7.

The optimal value of Δ increases with $|y|$, because, as can be seen in Figures 4.4 and 4.5, ω decreases with y , while $\varepsilon_{0,0}$ is independent. This means that when y is increased, the priority selection region can be widened, resulting in a higher probability of correctly identifying within that region, while still reducing the amount of work in that region. Because $\varepsilon_{0,0}$ is much greater for the higher correlation coefficient, only a small value of Δ is required to perform priority selection.

This optimal normalized running time is displayed in Figure 4.8 as a function of y , using the standard parameters. The optimal running time decreases with $|y|$. This is because the distribution of Y has most mass near 0, and therefore even for small values of Δ , relatively many entries will be included in the priority selection, whereas in the tails of the distribution, far fewer entries are considered for priority selection, even when larger values of Δ are used.

Figures 4.9 and 4.10 show the normalized optimal running time as a function of ρ and n respectively. The running time decreases with increasing ρ , because greater correlation allows to find the correct entry within a small region around λy with greater probability. The increase in running time with increasing n is due to another reason: for greater n , the same priority selection region in the \mathcal{X} space will translate to a larger region in the \mathcal{W} space than for smaller n . This means more unrelated entries will have to be processed on average before a match is found. In Figure 4.10 only solutions using the high value of the correlation coefficient, $\rho = 0.9999$, are shown, because it is unrealistic to use large values of n for $\rho = 0.99$.

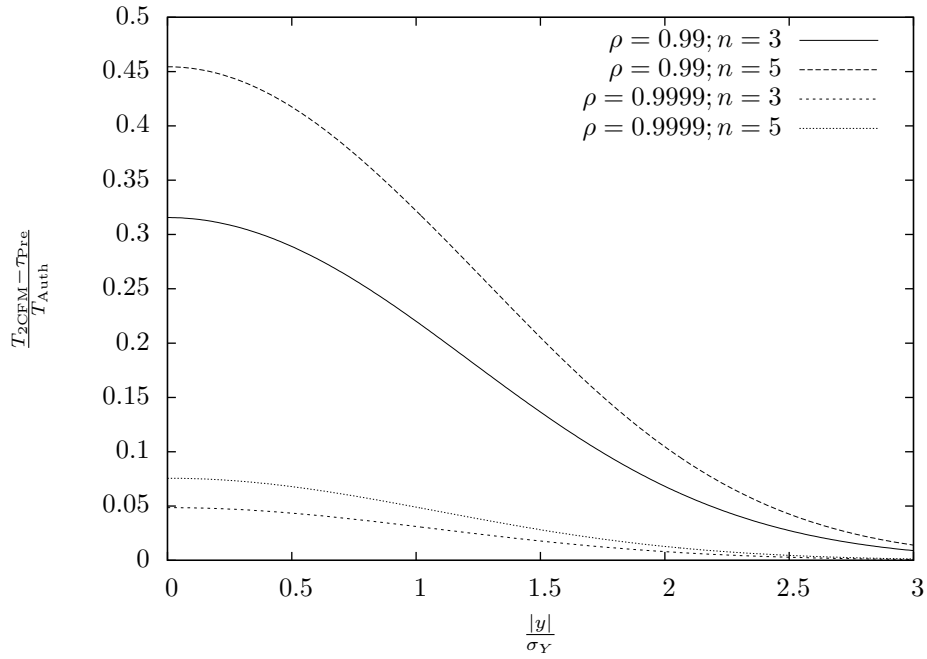


Figure 4.8: The optimal normalized running time as a function of y , using the standard parameters.

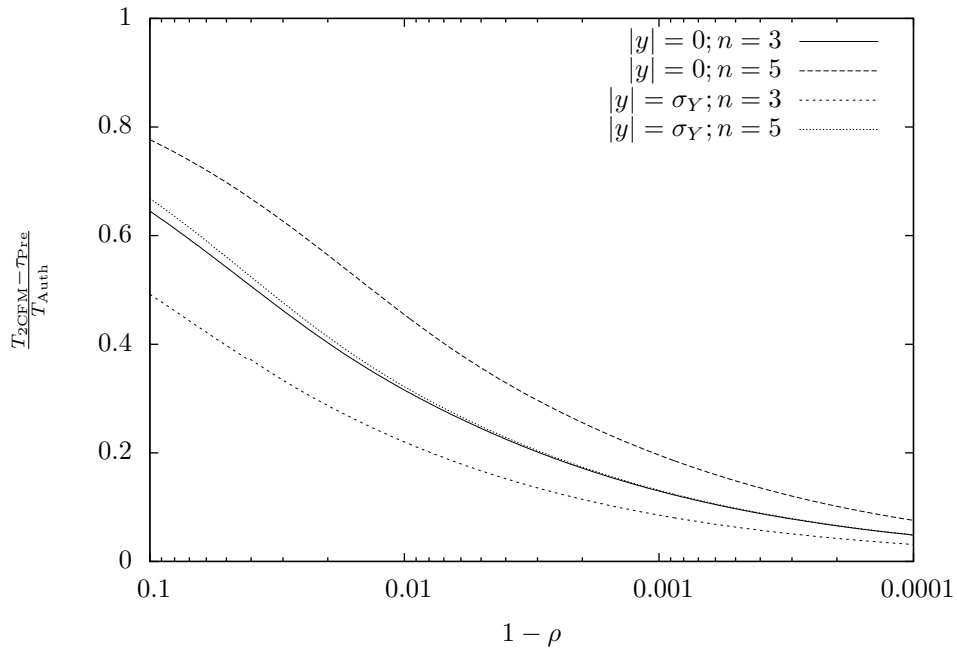


Figure 4.9: The optimal normalized running time as a function of ρ , using the standard parameters.

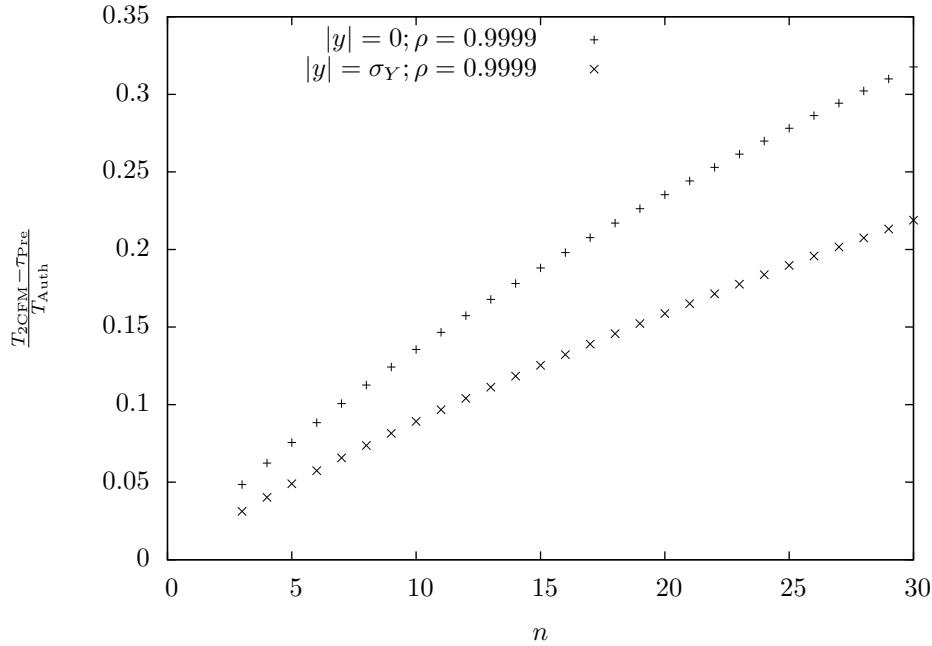


Figure 4.10: The optimal normalized running time as a function of n , using the standard values of y . Only the solution for $\rho = 0.9999$ is displayed, as large values of n are unrealistic for the lower value of the correlation coefficient.

4.3 Multidimensional Data

The analysis and results presented in this chapter apply only to the one-dimensional case, but the intent is for each dimensional component of multidimensional data to be processed independently, and for each of the results to be aggregated. This raises the following issues.

When using multidimensional data for which each dimensional component can be considered independently, the priority selection region can be determined separately for each dimensional component, but the decision whether an entry matches a presented biometric can only be made after aggregation of all one-dimensional reproduction steps. The most straightforward way of obtaining the multidimensional priority selection region is by simply combining the one-dimensional priority selection regions. The resulting work load and probability of finding a matching entry within the priority selection region will be the product of all one-dimensional ω and $\varepsilon_{0,0}$ values. Therefore, using the optimal one-dimensional priority selection regions will result in only a very small amount of work when the matching entry will be found in the priority selection region. However, the probability of doing so will also decrease, and consequently the cost of failure during this step will be relatively high. For this reason, the optimal priority selection region will have to be determined for the entire multidimensional system. It is possible that a different measure than the expected running time will have to be used to obtain the desired system performance. With decreasing probability of finding the correct entry within the priority selection region and its increasing relative cost, using more than two selection stages may become increasingly beneficial.

Figures 4.1 and 4.8 show that those components that have values near the center of the distribution are both likely to produce an error and incur a high processing cost. For a typical composition, it is expected that both outliers and points near the distribution's center will occur. Therefore, this observation suggests the use of error correcting codes as part of the reproduction process. Components with a value near the distribution's center, are at enrollment and verification time known to have poor reliability. This information can be used to choose an appropriate code. Apart

from correcting errors in reproducing particular components, a more resilient code can be used which would tolerate some components not being reproduced at all, or only when needed.

These issues are not considered within the scope of this thesis, but the questions they raise may form the basis for future work.

Chapter 5

Information Theoretic Approach

In Chapters 3 and 4, the properties of an identification system using multi-stage selection have been studied using an analytical approach. In this chapter, the relation between the biometric sample for identification, Y , and the helper data, W , will be studied information theoretically. The information theoretic work has been performed earlier on in the project, but this line of research has been abandoned in favor of the analytical approach detailed in Chapters 3 and 4.

In order to perform a fast biometric database lookup using multi stage selection, the first selection that has to be made cannot involve costly computations. With feature reconstruction considered costly, this means that it must be possible to reconstruct some part of the public data from the sample without any using any additional information.

The work in this chapter is based on the idea of assigning an entry to one or more clusters of entries during assignment, and attempting to perform ignorant reconstruction, i.e., without using any helper data, of at least one of these clusters correctly during verification [19].

The usual information theoretic approach to determining the error probabilities involved in such a reconstruction process is to apply Fano's inequality. This is not directly applicable in this case, however, because the data are only partially reconstructed.

First a generalization of Fano's inequality will be presented, that may or may not have been useful in further studying this issue. The chapter will conclude with a few numerical results regarding the mutual information contained between Y and W for the zero leakage helper data system.

5.1 Fano's Inequality

For two correlated random variables A and B , Fano's inequality relates the probability of correctly guessing A given B to the conditional entropy $H(A|B)$ [1].

Theorem 8 (Fano's inequality). *Let $\hat{A} = g(B)$ and $\varepsilon = \Pr[A \neq \hat{A}]$, then*

$$H(A|B) \leq h(\varepsilon) + \varepsilon \log(|\mathcal{A}| - 1). \quad (5.1)$$

Proof. Let E be an error indicator random variable such that $E = 0$ if $A = \hat{A}$ and $E = 1$ otherwise. The conditional entropy $H(E, A|B)$ can be expanded in two ways:

$$H(E, A|B) = H(A|B) + \underbrace{H(E|A, B)}_{=0} \quad (5.2)$$

and

$$\begin{aligned}
H(E, A|B) &= H(E|B) + H(A|B, E) \\
&= \underbrace{H(E|B)}_{\leq h(\varepsilon)} + \underbrace{(1 - \varepsilon)H(A|B, E = 0)}_{=0} + \underbrace{\varepsilon H(A|B, E = 1)}_{\leq \log(|\mathcal{A}|-1)}.
\end{aligned} \tag{5.3}$$

□

5.2 Multi-stage Selection

Fano's inequality is not directly applicable to analyze the properties of an identification system using multi-stage selection, because during preselection, no exact match is required. Therefore, in order to analyze this matter, Fano's inequality has to be modified by relaxing the goal of correctly guessing A from B exactly. Instead, the goal is to guess a value sufficiently similar to A . There are several ways to formulate such a requirement, one of which is based on set intersection. This is the approach used in this thesis.

Consider $A \subseteq \mathcal{C}$ for some set \mathcal{C} , i.e., $\mathcal{A} = \wp(\mathcal{C})$, and a guess to be acceptable if $A \cap B \neq \emptyset$. An inequality similar to Equation (5.1) can be derived using the same two expansions as in Equation (5.3), however, in this case, $H(A|B, E = 0) = 0$ would no longer hold, so it's necessary to inspect its behavior.

5.2.1 General Case

Theorem 9. *Let $\varepsilon = \Pr[A \cap g(B) \neq \emptyset]$. Without any further assumptions, this generalization of Fano's inequality for subset intersection becomes*

$$H(A|B) \leq h(\varepsilon) + |\mathcal{C}| - \varepsilon |g(B)|. \tag{5.4}$$

Proof. In the general case, there are $2^{|\mathcal{C}|-|g(B)|}$ subsets of \mathcal{C} such that the intersection of one subset and $g(B)$ is empty. Because there are $2^{|\mathcal{C}|}$ subsets of $|\mathcal{C}|$ in total, the number of subsets that do intersect with $g(B)$ is equal to $2^{|\mathcal{C}|} - 2^{|\mathcal{C}|-|g(B)|}$. Without making any assumptions as to the distributions of A and B , the conditional entropies can be bounded as:

$$H(A|B, E = 0) \leq \log \left(2^{|\mathcal{C}|} - 2^{|\mathcal{C}|-|g(B)|} \right) \tag{5.5}$$

$$H(A|B, E = 1) \leq \log \left(2^{|\mathcal{C}|} \right), \tag{5.6}$$

giving

$$\begin{aligned}
H(A|B, E) &\leq (1 - \varepsilon) \log \left(2^{|\mathcal{C}|} - 2^{|\mathcal{C}|-|g(B)|} \right) + \varepsilon \log \left(2^{|\mathcal{C}|-|g(B)|} \right) \\
&\leq (1 - \varepsilon)|\mathcal{C}| + \varepsilon(|\mathcal{C}| - |g(B)|) \\
&= |\mathcal{C}| - \varepsilon |g(B)|
\end{aligned} \tag{5.7}$$

□

Because the purpose of preselection is to reduce the number of enrolled entries to be fully processed, $|g(B)|$ must be much smaller than $|\mathcal{C}|$. Therefore, Equation (5.4) provides a very weak upper bound. However, it is an indication that high entropy may be achievable, which could be beneficial for privacy.

5.2.2 Fixed Cardinality

For the second case, the cardinalities $|A|$ and $|g(B)|$ are fixed. Furthermore, $|A| + |g(B)| \leq |\mathcal{C}|$, as otherwise the two sets would always intersect.

Theorem 10. *Let $\varepsilon = \Pr[A \cap g(B) \neq \emptyset]$ and both $|A|$ and $|g(B)|$ be fixed. In this case the generalization of Fano's inequality for subset intersection becomes*

$$H(A|B) \leq h(\varepsilon) + (1 - \varepsilon) \log \left(\binom{|\mathcal{C}|}{|A|} - \binom{|\mathcal{C}| - |g(B)|}{|A|} \right) + \varepsilon \log \binom{|\mathcal{C}| - |g(B)|}{|A|}. \quad (5.8)$$

Proof. The number of subsets of $|\mathcal{C}|$ having cardinality $|A|$ is given by the binomial coefficient $\binom{|\mathcal{C}|}{|A|}$. The subsets of \mathcal{C} that do not intersect with $g(B)$ must be subsets of $\mathcal{C} \setminus g(B)$, therefore, number of such subsets is given by $\binom{|\mathcal{C}| - |g(B)|}{|A|}$. Using these facts, the following conditional entropies can be bounded:

$$H(A|B, E = 0) \leq \log \left(\binom{|\mathcal{C}|}{|A|} - \binom{|\mathcal{C}| - |g(B)|}{|A|} \right) \quad (5.9)$$

$$H(A|B, E = 1) \leq \log \binom{|\mathcal{C}| - |g(B)|}{|A|}, \quad (5.10)$$

□

When $|A| = |g(B)| = 1$, Fano's original inequality emerges. Should either be greater than 1, Stirling's approximation, see for example page 945 of [4], can be applied to expand the binomial coefficients, giving the following corollary.

Corollary 3. *If $|A|$ and $|g(B)|$ are relatively small compared to $|\mathcal{C}|$, then Equation (5.8) has an upper bound of*

$$h(\varepsilon) + \varepsilon \log \left(\frac{|\mathcal{C}|^{|A|}}{|\mathcal{C}| - |g(B)|} \right) > H(A|B) - |A| \log |\mathcal{C}| + \log(|A|!) - 1. \quad (5.11)$$

Proof. Stirling's approximation can be used to estimate large factorials:

$$\sqrt{2\pi n} \left(\frac{n}{e} \right)^n \leq n! \leq e\sqrt{n} \left(\frac{n}{e} \right)^n. \quad (5.12)$$

Since $|A|$ and $|g(B)|$ are much smaller than $|\mathcal{C}|$:

$$\begin{aligned} \binom{|\mathcal{C}|}{|A|} &= \frac{|\mathcal{C}|!}{|A|!(|\mathcal{C}| - |A|)!} \\ &\leq \frac{1}{|A|!} \frac{e}{\sqrt{2\pi}} \frac{|\mathcal{C}|!}{(|\mathcal{C}| - |A|)!} \frac{\left(\frac{|\mathcal{C}|!}{e} \right)^{|\mathcal{C}|}}{\left(\frac{(|\mathcal{C}| - |A|)!}{e} \right)^{(|\mathcal{C}| - |A|)!}} \\ &= \frac{1}{|A|!} \frac{e}{\sqrt{2\pi}} \sqrt{\frac{|\mathcal{C}|}{|\mathcal{C}| - |A|}} \left(\frac{|\mathcal{C}|}{|\mathcal{C}| - |A|} \right)^{|\mathcal{C}| - |A|} \left(\frac{|\mathcal{C}|}{e} \right)^{|A|}. \end{aligned} \quad (5.13)$$

The factor $\left(\frac{|\mathcal{C}|}{|\mathcal{C}| - |A|} \right)^{|\mathcal{C}| - |A|} = \left(1 + \frac{|A|}{|\mathcal{C}| - |A|} \right)^{|\mathcal{C}| - |A|}$ can be estimated further using [20]:

$$e^{\frac{xy}{x+y}} < \left(1 + \frac{x}{y} \right)^y < e^x, \quad (5.14)$$

as

$$\left(\frac{|\mathcal{C}|}{|\mathcal{C}| - |A|}\right)^{|\mathcal{C}| - |A|} < e^{|A|} \quad (5.15)$$

Combining Equations (5.13) and (5.15) gives

$$\binom{|\mathcal{C}|}{|A|} < \frac{1}{|A|!} \frac{e}{\sqrt{2\pi}} \sqrt{\frac{|\mathcal{C}|}{|\mathcal{C}| - |A|}} |\mathcal{C}|^{|A|} \quad (5.16)$$

A lower limit on this binomial expansion by Stirling's approximation is given by multiplying the upper limit with $\frac{2\pi}{e^2}$:

$$\begin{aligned} \binom{|\mathcal{C}| - |g(B)|}{|A|} &\geq \frac{1}{|A|!} \frac{\sqrt{2\pi}}{e} \sqrt{\frac{|\mathcal{C}| - |g(B)|}{|\mathcal{C}| - |g(B)| - |A|}} \\ &\left(\frac{|\mathcal{C}| - |g(B)|}{|\mathcal{C}| - |g(B)| - |A|}\right)^{|\mathcal{C}| - |g(B)| - |A|} \left(\frac{|\mathcal{C}| - |g(B)|}{e}\right)^{|A|}. \end{aligned} \quad (5.17)$$

Using Equation (5.14) a lower bound can be established:

$$\begin{aligned} \left(\frac{|\mathcal{C}| - |g(B)|}{|\mathcal{C}| - |g(B)| - |A|}\right)^{|\mathcal{C}| - |g(B)| - |A|} &> e^{\frac{|A|(|\mathcal{C}| - |g(B)| - |A|)}{|\mathcal{C}| - |g(B)|}} \\ &= e^{|A|} e^{-\frac{|A|^2}{|\mathcal{C}| - |g(B)|}}, \end{aligned} \quad (5.18)$$

which can be combined with Equation (5.17) to give

$$\binom{|\mathcal{C}| - |g(B)|}{|A|} > \frac{1}{|A|!} \frac{\sqrt{2\pi}}{e} \sqrt{\frac{|\mathcal{C}| - |g(B)|}{|\mathcal{C}| - |g(B)| - |A|}} (|\mathcal{C}| - |g(B)|)^{|A|} e^{-\frac{|A|^2}{|\mathcal{C}| - |g(B)|}}. \quad (5.19)$$

These approximations can be substituted into Equation (5.8):

$$\begin{aligned} H(A|B) &< h(\varepsilon) + (1 - \varepsilon) [|A| \log |\mathcal{C}| - \log(|A|!)] \\ &+ (1 - \varepsilon) \left[\log \left(\frac{e}{\sqrt{2\pi}} \sqrt{\frac{|\mathcal{C}|}{|\mathcal{C}| - |A|}} - \frac{\sqrt{2\pi}}{e} \sqrt{\frac{|\mathcal{C}| - |g(B)|}{|\mathcal{C}| - |g(B)| - |A|}} e^{-\frac{|A|^2}{|\mathcal{C}| - |g(B)|}} \right) \right] \\ &+ \varepsilon \left[|A| \log(|\mathcal{C}| - |g(B)|) - \log(|A|!) + \log \left(\frac{e}{\sqrt{2\pi}} \sqrt{\frac{|\mathcal{C}| - |g(B)|}{|\mathcal{C}| - |g(B)| - |A|}} \right) \right]. \end{aligned} \quad (5.20)$$

Because $|A|$ and $|g(B)|$ are much smaller than $|\mathcal{C}|$, the two complex logarithms can be bounded very roughly from above by 1, as their arguments will lie strictly between 0 and 2. This then finally yields the bound for the conditional entropy:

$$H(A|B) < h(\varepsilon) + (1 - \varepsilon) |A| \log |\mathcal{C}| + \varepsilon |A| \log(|\mathcal{C}| - |g(B)|) - \log(|A|!) + 1. \quad (5.21)$$

□

5.3 Numerical Results

The mutual information between the biometric sample for identification and the helper data, $I(W; Y)$ for the zero leakage system is displayed in Figures 5.1 and 5.2 as a function of the correlation coefficient, ρ , and the number of quantization intervals, n , respectively. As expected,

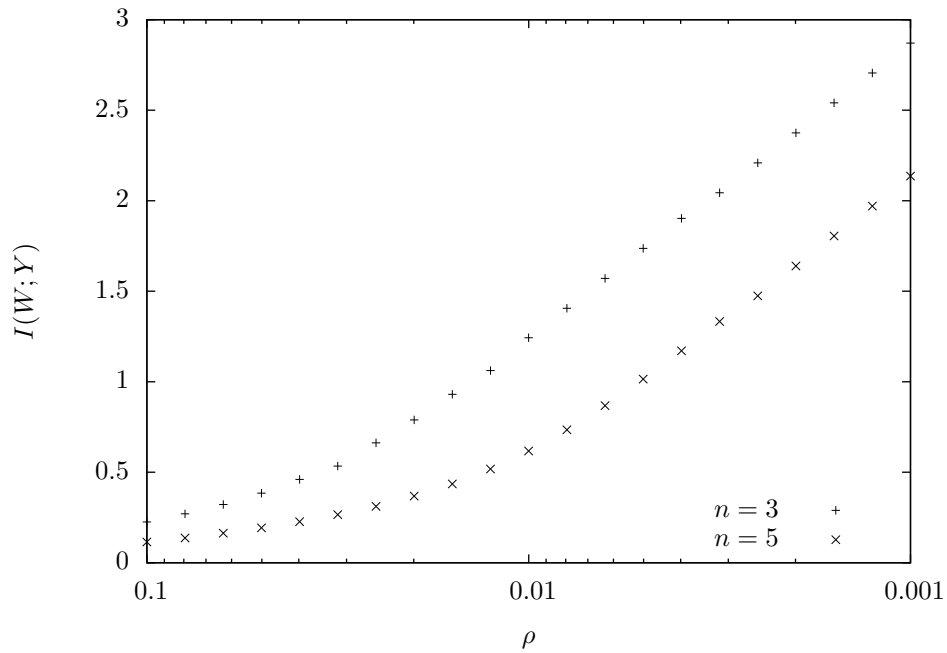


Figure 5.1: The mutual information between the helper data and the biometric sample for identification as a function of the correlation coefficient, using the standard parameters of Chapter 4.

the mutual information increases with increasing correlation coefficient, but decreases with the number of quantization intervals.

These results were obtained using numerical integration. For highly correlated systems, the integration produced invalid results due to numerical instabilities. For this reason, no results were obtained for $\rho = 0.9999$, which was used in Chapter 4.

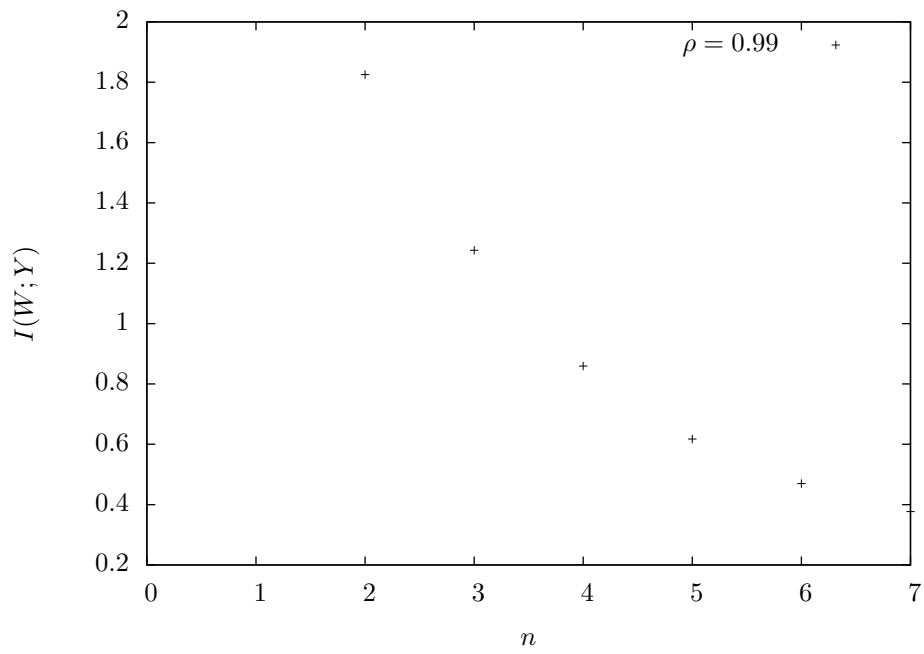


Figure 5.2: The mutual information between the helper data and the biometric sample for identification as a function of the number of quantization intervals, using $\rho = 0.99$.

Chapter 6

Conclusion

This chapter will conclude the thesis by giving a summary of the work presented within and some future work arising from it. Apart from this thesis, theoretical contributions were made to [2].

6.1 Summary

In Chapter 2, besides introducing notation and definitions, the issue of modifying an existing authentication system into an efficient identification system with the same privacy characteristics is investigated.

Chapter 3 deals with the matter of optimally reconstructing an enrolled entry in order to perform authentication. Starting out very general (Section 3.1), the analysis converges onto ZL systems (Section 3.2). Furthermore, a procedure for constructing a ZL authentication system for general continuous data is presented in Section 3.3 along with the system's underlying matching function which determines the system's characteristics in Section 3.4. The results of this section also motivates the PTS of Section 2.2.10, as a special case of the system presented in Section 3.3.

The properties of the PTS for one-dimensional data are further analysed in Chapter 4. Relations between the system's properties are derived analytically where possible and numerically using a Gaussian noise model where required. Results include the gain in running time relative to the naïve search implementation and work towards the trade-off between the system's resilience against errors, the running time and measure of privacy preservation. The chapter concludes with notes on the issues arising from the aggregation of many one-dimensional PTSs to handle multidimensional data.

In Chapter 5, the information theoretical work is presented that was performed for the original research question of how much entropy must be revealed by the helper data to allow for efficient searching. The chapter contains theoretical as well as numerical results, but this line of research was abandoned in favor of the analytical approach centered around the ZL property of Chapters 2 through 4.

6.2 Future Work

During the project, the focus of the research has deviated from what was originally intended. As a result, the original research questions have mostly gone unanswered and remain open for future work.

Contained within the results of Chapter 4 is evidence of a fundamental trade-off between the PTS's resilience against errors, the running time and measure of privacy preservation. To fully quantify this trade-off, the measure of privacy preservation needs to be well defined. This may require the definition of a model of privacy impact of different data contained in biometric samples, though even without such a fully defined model the trade-off may be investigated further.

The analysis and results presented in Chapter 4 apply only to one-dimensional data. The chapter concludes with notes on the issue of creating a system for multidimensional data by the aggregation of many one-dimensional systems. The resulting system's properties may be researched as future work.

Furthermore, the results in Chapter 4 indicate that for a typical composition of many independent dimensions some dimensional components will be problematic to search for and reconstruct, suggesting the use of error correcting codes as an intermediary step between the one-dimensional reconstruction processes and the aggregation thereof. The impact of applying error correcting codes could be a subject of future work.

Bibliography

- [1] Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. John Wiley & Sons, Inc., second edition, 2005.
- [2] Joep de Groot, Boris Skoric, Niels de Vreede, and Jean-Paul Linnartz. Information leakage of continuous-source zero secrecy leakage helper data schemes. Cryptology ePrint Archive, Report 2012/566, 2012. <http://eprint.iacr.org/>.
- [3] Y. Dodis, L. Reyzin, and A. Smith. Fuzzy extractors: How to generate strong keys from biometrics and other noisy data. In *LNCS*. Springer, 2004.
- [4] I. S. Gradshteyn and I. M. Ryzhik. *Table of Integrals, Series, and Products, Fifth Edition*. Academic Press, 5th edition, 1994.
- [5] J. Guajardo, B. Škorić, P. Tuyls, S.S. Kumar, T. Bel, A.H.M. Blom, and G.J. Schrijen. Anti-counterfeiting, key distribution, and key storage in an ambient world via Physical Unclonable Functions. *Information Systems Frontiers*, 11(1):19–41, 2009.
- [6] D.E. Holcomb, W.P. Burleson, and K. Fu. Power-Up SRAM State as an Identifying Fingerprint and Source of True Random Numbers. *Computers, IEEE Transactions on*, 58(9):1198–1210, sept. 2009.
- [7] Tanya Ignatenko and Frans M. J. Willems. Fundamental limits for biometric identification with a database containing protected templates. In *ISITA*, pages 54–59. IEEE, 2010.
- [8] A. Juels and M. Sudan. A fuzzy vault scheme. *Des. Codes Cryptogr.*, 38:237–257, 2006.
- [9] Ari Juels and Martin Wattenberg. A fuzzy commitment scheme. In Juzar Motiwalla and Gene Tsudik, editors, *ACM Conference on Computer and Communications Security*, pages 28–36. ACM, 1999.
- [10] J.-P. Linnartz and P. Tuyls. New shielding functions to enhance privacy and prevent misuse of biometric templates. In *Audio- and Video-Based Biometric Person Authentication*. Springer, 2003.
- [11] Tsutomu Matsumoto, Hiroyuki Matsumoto, Koji Yamada, and Satoshi Hoshino. Impact of artificial ”gummy” fingers on fingerprint systems. *Optical Security and Counterfeit Deterrence Techniques*, 4677:275–289, 2002.
- [12] G. Edward Suh and Srinivas Devadas. Physical unclonable functions for device authentication and secret key generation. In *Proceedings of the 44th annual Design Automation Conference, DAC '07*, pages 9–14, New York, NY, USA, 2007. ACM.
- [13] Pim Tuyls, Geert Jan Schrijen, Boris Škorić, Jan van Geloven, Nynke Verhaegh, and Rob Wolters. Read-proof hardware from protective coatings. In L. Goubin and M. Matsui, editors, *Cryptographic Hardware and Embedded Systems (CHES) 2006*, volume 4249 of *LNCS*, pages 369–383. Springer-Verlag, 2006.

- [14] Pim Tuyls and Boris Škorić. Physical unclonable functions for enhanced security of tokens and tags. In *Highlights of the Information Security Solutions Europe (ISSE) 2006 Conference*, pages 30–37. Vieweg, 2006. Part 1.
- [15] Pim Tuyls, Boris Škorić, and Tom Kevenaar. *Security with Noisy Data: Private Biometrics, Secure Key Storage and Anti-Counterfeiting*. Springer-Verlag New York, Inc., Secaucus, NJ, USA, 2007.
- [16] Ton van der Putte and Jeroen Keuning. Biometrical fingerprint recognition: don't get your fingers burned. In *Proceedings of the fourth working conference on smart card research and advanced applications on Smart card research and advanced applications*, pages 289–303, Norwell, MA, USA, 2001. Kluwer Academic Publishers.
- [17] E. A. Verbitskiy, P. Tuyls, C. Obi, B. Schoenmakers, and B. Škorić. Key extraction from general nondiscrete signals. *Information Forensics and Security, IEEE Transactions on*, 5(2):269–279, June 2010.
- [18] Sviatoslav Voloshynovskiy, Oleksiy Koval, Fokko Beekhof, and Thierry Pun. Unclonable identification and authentication based on reference list decoding. In *Proceedings of the conference on Secure Component and System Identification*, Berlin, Germany, March 17–18 2008.
- [19] Frans M. J. Willems. Searching methods for biometric identification systems: Fundamental limits. In *ISIT*, pages 2241–2245. IEEE, 2009.
- [20] Exponential function: Inequalities. <http://functions.wolfram.com/ElementaryFunctions/Exp/29/>. Retrieved: 14 March 2013.