

UNIVERSITY NAME

DOCTORAL THESIS

---

# Thesis Title

---

*Author:*  
John SMITH

*Supervisor:*  
Dr. James SMITH

*A thesis submitted in fulfillment of the requirements  
for the degree of Doctor of Philosophy  
in the*

Research Group Name  
Department or School Name

September 7, 2017



# Contents

<b>I</b>	<b>Context and State of the Art</b>	<b>1</b>
<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Introduction to Cryptography . . . . .	3
1.1.1	Secret-Key Cryptography . . . . .	4
1.1.2	Public-Key Cryptography . . . . .	4
1.2	Embedded Cryptography and Secure Hardware . . . . .	4
1.2.1	The Example of the Smart Card . . . . .	4
1.2.2	Certification of a Secure Hardware . . . . .	4
1.2.3	Modern More Complex Devices to Certify . . . . .	4
1.2.4	Embedded Cryptography Vulnerabilities . . . . .	4
<b>2</b>	<b>Introduction to Side-Channel Attacks</b>	<b>5</b>
2.1	Introduction to Side-Channel Attacks . . . . .	5
2.1.1	Historical Overview . . . . .	5
2.1.2	Terminology and Generalities . . . . .	5
	Target and Leakage Model . . . . .	5
	Points of Interest . . . . .	5
	Simple vs Advanced SCAs . . . . .	5
	Vertical vs Horizontal SCAs . . . . .	5
	Profiled vs Non-Profiled SCAs . . . . .	5
	Side-Channel Algebraic Attacks . . . . .	5
	Distinguishers . . . . .	5
	SCA Metrics . . . . .	5
2.2	Main Side-Channel Countermeasures . . . . .	5
2.2.1	Random Delays and Jitter . . . . .	5
2.2.2	Shuffling . . . . .	5
2.2.3	Masking . . . . .	5
2.3	Higher-Order Attacks . . . . .	5
2.3.1	Higher-Order Moments Analysis and Combining Functions . . . . .	5
2.3.2	Profiling Higher-Order Attacks . . . . .	5
	Profiling with Masks Knowledge . . . . .	5
	Profiling without Masks Knowledge . . . . .	5
<b>3</b>	<b>Introduction to Machine Learning</b>	<b>7</b>
3.1	Basic Concepts of Machine Learning . . . . .	7
3.1.1	The Task, the Experience and the Performance . . . . .	7
3.1.2	Supervised, Semi-Supervised, Unsupervised Learning . . . . .	7
3.1.3	Training, Validation and Test Sets . . . . .	7
3.1.4	Underfitting, Overfitting and Regularization . . . . .	7
3.1.5	Data Augmentation . . . . .	7
3.1.6	No Free Lunch Theorem . . . . .	7

3.2	Machine Learning Applications in Side-Channel Context . . . . .	7
3.2.1	Profiled Attack as a Classification Problem . . . . .	7
	Support Vector Machine . . . . .	7
	Random Forest . . . . .	7
	Neural Networks . . . . .	7
<b>II</b>	<b>Contributions</b>	<b>9</b>
<b>4</b>	<b>Introduction</b>	<b>11</b>
4.1	Foreword of this Thesis: Research of Points of Interest . . . . .	11
4.2	Dimensionality Reduction Approach . . . . .	11
4.2.1	Linear Methods for First-Order Attacks . . . . .	11
4.2.2	Kernel Methods for Higher-Order Attacks . . . . .	11
4.3	Neural Network Approach . . . . .	11
4.3.1	Toward Getting Rid of Information-Loosing Preprocessing . . . . .	11
<b>5</b>	<b>Points of Interest</b>	<b>13</b>
5.1	Motivations . . . . .	13
5.1.1	The Curse of Dimensionality . . . . .	13
5.2	Selection on Points of Interest: Classical Statistics . . . . .	13
5.3	Related Issues: Leakage Detection and Leakage Assessment . . . . .	13
5.4	Generalized SNR for Multi-Variate Attacks . . . . .	13
5.5	Observations Leading to Take a Dimensionality Reduction Approach . . . . .	13
<b>6</b>	<b>Linear Dimensionality Reduction</b>	<b>15</b>
6.1	Introduction . . . . .	15
6.1.1	Principal Component Analysis . . . . .	15
6.1.2	Linear Discriminant Analysis . . . . .	15
6.1.3	Projection Pursuits . . . . .	15
6.2	Principal Component Analysis . . . . .	15
6.2.1	Statistical Point of View . . . . .	15
6.2.2	Geometrical Point of View . . . . .	15
6.3	Application of PCA in SCAs . . . . .	15
6.3.1	Original vs Class-Oriented PCA . . . . .	15
6.3.2	The Choice of the Principal Components . . . . .	15
6.3.3	The Explained Local Variance Selection Method . . . . .	17
6.4	Linear Discriminant Analysis . . . . .	18
6.4.1	Statistical Point of View . . . . .	18
6.4.2	Geometrical Point of View . . . . .	18
6.5	Application of LDA in SCAs . . . . .	18
6.5.1	The Small Sample Size problem . . . . .	18
6.6	Experimental Results . . . . .	18
	The testing adversary. . . . .	19
	Scenario 1. . . . .	19
	Scenario 2. . . . .	19
6.7	Misaligning Effects . . . . .	20

<b>7</b>	<b>Kernel Dimensionality Reduction</b>	<b>23</b>
7.1	Motivation	23
7.1.1	Higher-Order Attacks	23
	Higher-Order Version of Projection Pursuits	23
7.2	Kernel Function and Kernel Trick	23
7.2.1	Local Kernel Functions as Similarity Metrics	23
7.3	Kernel Discriminant Analysis	23
7.4	Experiments over Atmega328P	23
7.4.1	The Regularization Problem	23
7.4.2	The Multi-Class Trade-Off	23
7.4.3	Multi-Class vs 2-class Approach	23
7.4.4	Asymmetric Preprocessing/Attack Approach	23
	Comparison with Projection Pursuits	23
7.5	Drawbacks of Kernel Methods	23
	Misalignment Effects	23
	Memory Complexity and Actual Number of Parameters	23
	Two-Phases Approach: Preprocessing-Templates	23
<b>8</b>	<b>Convolutional Neural Networks against Jitter-Based Countermeasures</b>	<b>25</b>
8.1	Moving from Kernel Machines to Neural Networks	25
8.2	Misalignment of Side-Channel Traces	25
8.2.1	The Necessity and the Risks of Applying Realignment Techniques	25
8.2.2	Analogy with Image Recognition Issues	25
8.3	Convolutional Layers to Impose Shift-Invariance	25
8.4	Data Augmentation for Misaligned Side-Channel Traces	25
8.5	Experiments against Software Countermeasures	25
8.6	Experiments against Artificial Hardware Countermeasures	25
8.7	Experiments against Real-Case Hardware Countermeasures	25
<b>9</b>	<b>KDA vs Neural Networks Approach for HO-Attacks</b>	<b>27</b>
9.1	Simulated Experiment for Profiled HO-Attacks	27
9.1.1	The Simulations	27
9.1.2	Comparison between KDA and MLP	27
9.2	Real-Case Experiments over ARM Cortex-M4	27
<b>10</b>	<b>Siamese Neural Networks for Collision Attacks</b>	<b>29</b>
10.1	Introduction	29
10.2	Siamese Neural Networks	29
10.2.1	Distances and Loss Functions	29
10.2.2	Relation with Kernel Machines	29
10.3	Collision Attacks with Siamese NNs	29
10.3.1	Experimental Results	29
<b>11</b>	<b>Conclusions and Perspectives</b>	<b>31</b>
11.1	Summary	31
11.2	Strengthen Embedded Security: the Main Challenge for Machine Learning Applications	31
<b>A</b>	<b>Scenario 3 and 4 of CARDIS '15 paper</b>	<b>33</b>
	Scenario 3.	33
	Scenario 4.	33



# List of Figures

6.1	First and sixth PCs in DPA contest v4 trace set (between time samples 198001 and 199000) . . . . .	16
6.2	The first six PCs. Acquisition campaign on an 8-bit AVR Atmega328P (see Sec. 6.6). . . . .	17
6.3	Guessing Entropy as function of the number of attack traces . . . . .	19
6.4	Guessing Entropy as function of the number of profiling traces. . . . .	21
6.5	Degradation of linear-reduction-based template attacks in presence of misalignment. . . . .	22





# List of Tables



# List of Abbreviations

SCA Side Channel Attack



# List of Symbols



## **Part I**

# **Context and State of the Art**





## Chapter 1

# Introduction

### 1.1 Introduction to Cryptography

The terms *Cryptography* (from the Greek *kryptòs* (secret) and *graphein* (writing)) and *Cryptanalysis*, denotes two branches of a science named *Cryptology*, or *science of secret*. Cryptography initially refers to the art of *encrypting* messages, which means writing meaningful messages in such a way to appear nonsense to anyone unaware of the encryption process. In general, cryptography aims to construct protocols to secure communication, while cryptanalysis studies the resistance of cryptographic techniques, developing *attacks* to break the cryptosystems' security claims. These two complementary domains evolve in parallels, since the evolution of attack techniques allows conceiving more resistant cryptographic algorithms, and inversely the resistance of such algorithms requires the conception of more sophisticated attacks.

The art of cryptography is very ancient, probably as ancient as the language, but only the development of information technology made cryptology take the shape of a proper science, sometimes referred to as *Modern cryptology*. The last be seen as a branch of different disciplines, such as applied mathematics, computer science, electrical engineering, and communication science. Modern cryptosystems exploit algorithms based on mathematical tools and are implemented as computer programs, or electronic circuits. Their goal is to provide security functionality for communications that use *insecure channels*, for example the internet. In particular, modern cryptosystems are designed in order to ensure at least one of the four following information security properties:

- a. *confidentiality*: the transmitted message must be readable only by a chosen pool of authorized entities;
- b. *authenticity*: the receiver can verify the identity of the sender of a message;
- c. *non-repudiation*: the sender of a message cannot deny having sent the message afterwards;
- d. *data integrity*: the receiver can be convinced that the message has not been corrupted during the transmission.

Two branches of cryptography may be distinguished: the *symmetric cryptography* and the *asymmetric cryptography*. The first one historically appeared before and is based on the hypothesis that the two communicating entities share a common secret, or private key; for this reason this is also called *secret key cryptography*. The second one, introduced around 1970, allow any entity to encrypt a message in such a way that only a unique chosen other entity could decrypt it; this is also called *public key cryptography*. In the following part of this section the two standard cryptographic

primitives, *i.e.* building block algorithms used to build cryptographic protocols, that will be used in this thesis; a symmetric one, the *Advanced Encryption Standard*, and an asymmetric one, the RSA.

### **1.1.1 Secret-Key Cryptography**

### **1.1.2 Public-Key Cryptography**

## **1.2 Embedded Cryptography and Secure Hardware**

### **1.2.1 The Example of the Smart Card**

### **1.2.2 Certification of a Secure Hardware**

### **1.2.3 Modern More Complex Devices to Certify**

### **1.2.4 Embedded Cryptography Vulnerabilities**

## Chapter 2

# Introduction to Side-Channel Attacks

### 2.1 Introduction to Side-Channel Attacks

#### 2.1.1 Historical Overview

#### 2.1.2 Terminology and Generalities

Target and Leakage Model

Points of Interest

Simple vs Advanced SCAs

Vertical vs Horizontal SCAs

Profiled vs Non-Profiled SCAs

Side-Channel Algebraic Attacks

Distinguishers

SCA Metrics

### 2.2 Main Side-Channel Countermeasures

#### 2.2.1 Random Delays and Jitter

#### 2.2.2 Shuffling

#### 2.2.3 Masking

### 2.3 Higher-Order Attacks

#### 2.3.1 Higher-Order Moments Analysis and Combining Functions

#### 2.3.2 Profiling Higher-Order Attacks

Profiling with Masks Knowledge

Profiling without Masks Knowledge



## Chapter 3

# Introduction to Machine Learning

### 3.1 Basic Concepts of Machine Learning

#### 3.1.1 The Task, the Experience and the Performance

#### 3.1.2 Supervised, Semi-Supervised, Unsupervised Learning

#### 3.1.3 Training, Validation and Test Sets

#### 3.1.4 Underfitting, Overfitting and Regularization

#### 3.1.5 Data Augmentation

#### 3.1.6 No Free Lunch Theorem

### 3.2 Machine Learning Applications in Side-Channel Context

#### 3.2.1 Profiled Attack as a Classification Problem

Support Vector Machine

Random Forest

Neural Networks



## **Part II**

# **Contributions**





## **Chapter 4**

# **Introduction**

### **4.1 Foreword of this Thesis: Research of Points of Interest**

### **4.2 Dimensionality Reduction Approach**

#### **4.2.1 Linear Methods for First-Order Attacks**

#### **4.2.2 Kernel Methods for Higher-Order Attacks**

### **4.3 Neural Network Approach**

#### **4.3.1 Toward Getting Rid of Information-Loosing Preprocessing**



## **Chapter 5**

# **Points of Interest**

### **5.1 Motivations**

#### **5.1.1 The Curse of Dimensionality**

### **5.2 Selection on Points of Interest: Classical Statistics**

### **5.3 Related Issues: Leakage Detection and Leakage Assessment**

### **5.4 Generalized SNR for Multi-Variate Attacks**

### **5.5 Observations Leading to Take a Dimensionality Reduction Approach**



## Chapter 6

# Linear Dimensionality Reduction

### 6.1 Introduction

#### 6.1.1 Principal Component Analysis

#### 6.1.2 Linear Discriminant Analysis

#### 6.1.3 Projection Pursuits

### 6.2 Principal Component Analysis

#### 6.2.1 Statistical Point of View

#### 6.2.2 Geometrical Point of View

### 6.3 Application of PCA in SCAs

#### 6.3.1 Original vs Class-Oriented PCA

*Remark.* Stacked Auto-Encoders...

#### 6.3.2 The Choice of the Principal Components

The introduction of the PCA method in SCA context (either in its classical or class-oriented version) has raised some important questions: *how many* principal components and *which ones* are sufficient/necessary to reduce the trace size (and thus the attack processing complexity) without losing important discriminative information?

Until now, an answer to the questions above has been given in [choudary2014efficient], linked to the concept of *explained variance* (or *explained global variance*, EGV for short) of a PC  $\alpha_i$ :

$$\text{EGV}(\alpha_i) = \frac{\lambda_i}{\sum_{k=1}^r \lambda_k}, \quad (6.1)$$

where  $r$  is the rank of the covariance matrix  $\mathbf{S}$ , and  $\lambda_j$  is the eigenvalue associated to the  $j$ -th PC  $\alpha_j$ .  $\text{EGV}(\alpha_i)$  is the variance of the data projected over the  $i$ -th PC (which equals  $\lambda_i$ ) divided by the total variance of the original data (given by the trace of the covariance matrix  $\mathbf{S}$ , *i.e.* by the sum of all its non-zero eigenvalues). By definition of EGV, the sum of all the EGV values is equal to 1; that is why this quantity is often multiplied by 100 and expressed as percentage. Exploiting the EGV to choose among the PCs consists in fixing a wished *cumulative explained variance*  $\beta$  and in keeping  $C$  different PCs, where  $C$  is the minimum integer such that

$$\text{EGV}(\alpha_1) + \text{EGV}(\alpha_2) + \dots + \text{EGV}(\alpha_C) \geq \beta. \quad (6.2)$$

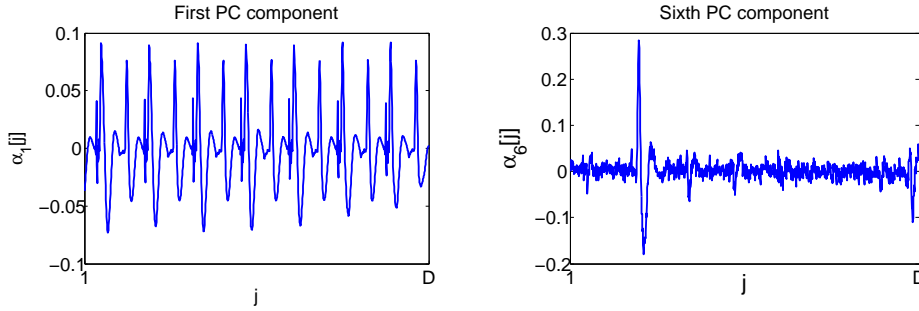


FIGURE 6.1: First and sixth PCs in DPA contest v4 trace set (between time samples 198001 and 199000)

However, if the adversary has a constraint for the reduced dimension  $C$ , the EGV notion simply suggests to keep the first  $C$  components, taking for granted that the optimal way to choose PCs is in their natural order. This assumption is not always confirmed in SCA context: in some works, researchers have already remarked that the first components sometimes contain more noise than information [Batina2012; specht] and it is worth discarding them. For the sake of providing a first example of this behaviour on publicly accessible traces, we applied a class-oriented PCA on 3000 traces from the DPA contest v4 [DPAcontest]; we focused over a small 1000-dimensional window in which, in complete knowledge about masks and other countermeasures, information about the first Sbox processing leaks (during the first round). In Fig. 6.1 the first and the sixth PCs are plotted. It may be noticed that the first component indicates that one can attend a high variance by exploiting the regularity of the traces, given by the clock signal, while the sixth one has high coefficients localised in a small time interval, very likely to signalize the instants in which the target sensitive variable leaks.

To the best of our knowledge, a single method adapted to SCA context has been proposed until now to automatically choose PCs [SCAclassProbl] while dealing with the issue raised in Fig. 6.1. It is based on the following assumption:

*Assumption 1.* The leaking side-channel information is localised in few points of the acquired trace.

In the rest of the paper, we conduct our own analyses under Assumption 1 that we think to be reasonable in SCA contexts where the goal of the security developers is to minimize the number of leaking points. Under this assumption, the authors of [SCAclassProbl] use for side-channel attack purposes the *Inverse Participation Ratio* (IPR), a measure widely exploited in Quantum Mechanics domain (see for example [guhr1998random]). They propose to use such a score to evaluate the eigenvectors *localization*. It is defined as follows:

$$\text{IPR}(\alpha_i) = \sum_{j=1}^D \alpha_i[j]^4. \quad (6.3)$$

The authors of [SCAclassProbl] suggest to collect the PCs in decreasing order with respect to the IPR score.

The selection methods provided by the evaluation of the EGV and of the IPR are somehow complementary: the former is based only on the eigenvalues associated to the PCs and does not consider the form of the PCs themselves; the latter completely

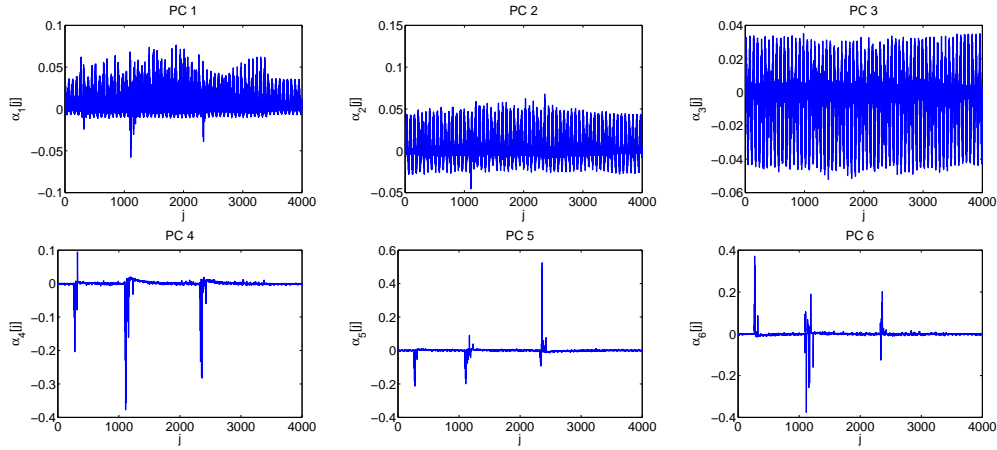


FIGURE 6.2: The first six PCs. Acquisition campaign on an 8-bit AVR Atmega328P (see Sec. 6.6).

discards the information given by the eigenvalues of the PCs, considering only the distribution of their coefficients. One of the contributions of the present paper is to propose a new selection method, that builds a bridge between the EGV and the IPR approaches. As we will argue, our method, based on the so-called *explained local variance*, does not only lead to the construction of a new selection criterion, but also permits to modify the PCs, choosing individually the coefficients to keep and those to discard.

### 6.3.3 The Explained Local Variance Selection Method

#### Riprendere le notazioni e mettere apposto i newcommand

*Definition 1.* The *Explained Local Variance* of a PC  $\alpha_i$  in a sample  $j$ , is defined by

$$\text{ELV}(\alpha_i, j) = \frac{\lambda_i \alpha_i[j]^2}{\sum_{k=1}^r \lambda_k} = \text{EGV}(\alpha_i) \alpha_i[j]^2. \quad (6.4)$$

Let  $\mathcal{J} = \{j_1^i, j_2^i, \dots, j_D^i\} \subset \{1, 2, \dots, D\}$  be a set of indexes sorted such that  $\text{ELV}(\alpha_i, j_1^i) \geq \text{ELV}(\alpha_i, j_2^i) \geq \dots \geq \text{ELV}(\alpha_i, j_D^i)$ . It may be observed that the sum over all the  $\text{ELV}(\alpha_i, j)$ , for  $j \in [1, \dots, D]$ , equals  $\text{EGV}(\alpha_i)$ . If we operate such a sum in a cumulative way following the order provided by the sorted set  $\mathcal{J}$ , we obtain a complete description of the trend followed by the component  $\alpha_i$  to achieve its EGV. As we can see in Fig. ??, where such cumulative ELVs are represented, the first 3 components are much slower in achieving their final EGV, while the 4<sup>th</sup>, the 5<sup>th</sup> and the 6<sup>th</sup> achieve a large part of their final EGVs very quickly (*i.e.* by adding the ELV contributions of much less time samples). For instance, for  $i = 4$ , the sum of the  $\text{ELV}(\alpha_4, j_k^4)$ , with  $k \in [1, \dots, 30]$ , almost equals  $\text{EGV}(\alpha_4)$ , whereas the same sum for  $i = 1$  only achieves about the 15% of  $\text{EGV}(\alpha_1)$ . Actually, the EGV of the 4<sup>th</sup>, the 5<sup>th</sup> and the 6<sup>th</sup> component only essentially depends on a very few time samples. This observation, combined with Assumption 1, suggests that they are more suitable for SCA than the three first ones. To validate this statement, it suffices to look at the form of such components (Fig. 6.2): the leading ones are very influenced by the clock, while the latest ones are well localised over the leaking points.

Operating a selection of components *via* ELV, in analogy with the EGV, requires to fix the reduced space dimension  $C$ , or a threshold  $\beta$  for the cumulative ELV. In the

first case, the maximal ELVs of each PC are compared, and the  $C$  components achieving the highest values of such ELVs are chosen. In the second case, all pairs (PC, time sample) are sorted in decreasing order with respect to their ELV, and summed until the threshold  $\beta$  is achieved. Then only PCs contributing in this sum are selected.

We remark that the ELV is a score associated not only to the whole components, but to each of their coefficients. This interesting property can be exploited to further remove, within a selected PC, the non-significant points, *i.e.* those with a low ELV. In practice this is done by setting these points to zero. That is a natural way to exploit the ELV score in order to operate a kind of *denoising* for the reduced data, making them only depend on the significant time samples. In Sec. 6.6 (scenario 4) we test the performances of an attack varying the number of time samples involved in the computation of the reduced data, and showing that such a denoising processing might impact significantly.

## 6.4 Linear Discriminant Analysis

### 6.4.1 Statistical Point of View

### 6.4.2 Geometrical Point of View

## 6.5 Application of LDA in SCAs

### 6.5.1 The Small Sample Size problem

## 6.6 Experimental Results

In this section we compare the different extractors provided by the PCA and the LDA in association with the different techniques of components selection. Defining an universal criterion to compare the different extractors would not make sense since the latter one should encompass a lot of parameters, sometimes opposite, that vary according to the context (amount of noise, specificity of the information leakage, nature of the side channel, etc.). For this reason we choose to split our comparisons into four scenarios. Each scenario has a single varying parameter that, depending on the attacker context, may wish to be minimized. Hereafter the definition of the four scenario. In the following only results of the two first is reported, the interested reader might refer to Appendix A for results of in the two other scenarios.

Scenario 1 varying parameter: number of attack traces  $N_a$ ,

Scenario 2 varying parameter: number of profiling traces  $N_p$ ,

Scenario 3 varying parameter: number of projecting components selected  $C$ ,

Scenario 4 varying parameter: number of original time samples implied into the trace preprocessing computation  $\#PoI$ .

For scenarios in which  $nbProfilingTraces$  is fixed, the value of  $N_p$  is chosen high enough to avoid the SSS problem, and the extensions of LDA presented in Sec. ?? are not evaluated. This choice of  $N_p$  will imply that the LDA is always performed in a favourable situation, which makes expect the LDA to be particularly efficient for these experiments. Consequently, for the scenarios in which  $N_p$  is high, our goal is to study whether the



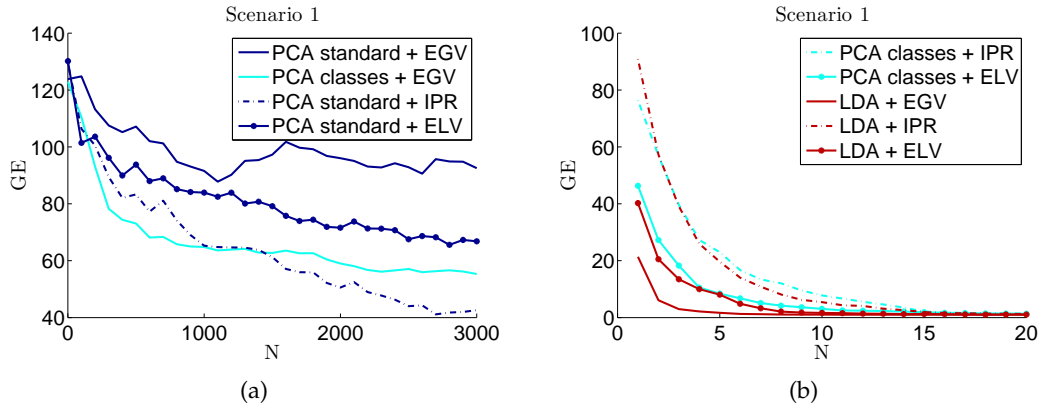


FIGURE 6.3: Guessing Entropy as function of the number of attack traces for different extraction methods. All Guessing Entropies are estimated as the average rank of the right key over 100 independent experiments.

PCA can be made almost as efficient as the LDA thanks to the component selection methods discussed in Sec. 6.3.2.

**This part will maybe be useless: somewhere I will have described all trace sets**

**The testing adversary.**

**Scenario 1.**

**cos'e**  $N_z$  To analyse the dependence between the extraction methods presented in Sections ?? and ?? and the number of attack traces  $N_a$  needed to achieve a given GE, we fixed the other parameters as follows:  $N_z = 50$  ( $N_p = 50 \times 256$ ),  $C = 3$  and  $\sharp\text{PoI} = 3996$  (all points are allowed to participate in the building of the PCs and of the LDCs). The experimental results, depicted in Fig. 6.3(a)-(b), show that the PCA standard method has very bad performances in SCA, while the LDA outperforms the others. Concerning the class-oriented PCA, we observe that its performance is close to that of LDA when combined with the selection methods ELV (which performs best) or IPR.

**Scenario 2.**

Now we test the behaviour of the extraction methods as function of the number  $N_z$  of available profiling traces per class. The number of components  $C$  is still fixed to 3,  $\sharp\text{PoI} = 3996$  again and the number of attack traces is  $N_a = 100$ . This scenario has to be divided into two parts: if  $N_z \leq 15$ , then  $N_p < D$  and the SSS problem occurs. Thus, in this case we test the four extensions of LDA presented in Sec. ??, associated to either the standard selection, to which we abusively refer as EGV,<sup>1</sup> or to the IPR selection. We compare them to the class-oriented PCA associated to EGV, IPR or ELV. The ELV selection is not performed for the techniques extending LDA, since for some of them the projecting LDCs are not associated to some eigenvalues in a meaningful way. On the contrary, if  $N_z \geq 16$  there is no need to approximate the LDA technique, so the classical one is performed. Results for this scenario are shown in Fig. 6.4. It may be noticed that the combinations class-oriented PCA + ELV/IPR select exactly the same components, for our data, see Fig. 6.4(e) and do

<sup>1</sup>It consists in keeping the  $C$  first LDCs (the  $C$  last for the Direct LDA)

not suffer from the lack of profiling traces. They are slightly outperformed by the  $S_W$  Null Space method associated with the EGV, see Fig.6.4(d). The Direct LDA (Fig. 6.4(b)) method also provides a good alternative, while the other tested methods do not show a stable behaviour. The results in absence of the SSS problem (Fig.6.4(f)) confirm that the standard PCA is not adapted to SCA, even when provided with more profiling traces. It also shows that among class-oriented PCA and LDA, the class-oriented PCA converges faster.

## 6.7 Misaligning Effects

**give parameters: 6 4 citare Choudary, Template Attacks over different devices**

In this section we experimentally show how the approach based on linear dimensionality reduction described in this chapter is affected by traces misalignment. To this aim, we simply take the same data and parameters exploited for Scenario 1 in Sec. 6.6, and artificially misalign them through the technique proposed in Appendix A with parameters **parameters here**. Then we tried to pre-process attack them through the 9 methodology tested in Scenario 1. It may be noticed in Fig. 6.5 that none of the 9 techniques is still efficient, included the optimal LDA+EGV that lead to null guessing entropy the synchronized traces using less 7 attack traces. In this case it cannot lead to successful attack in less than 3000 traces.

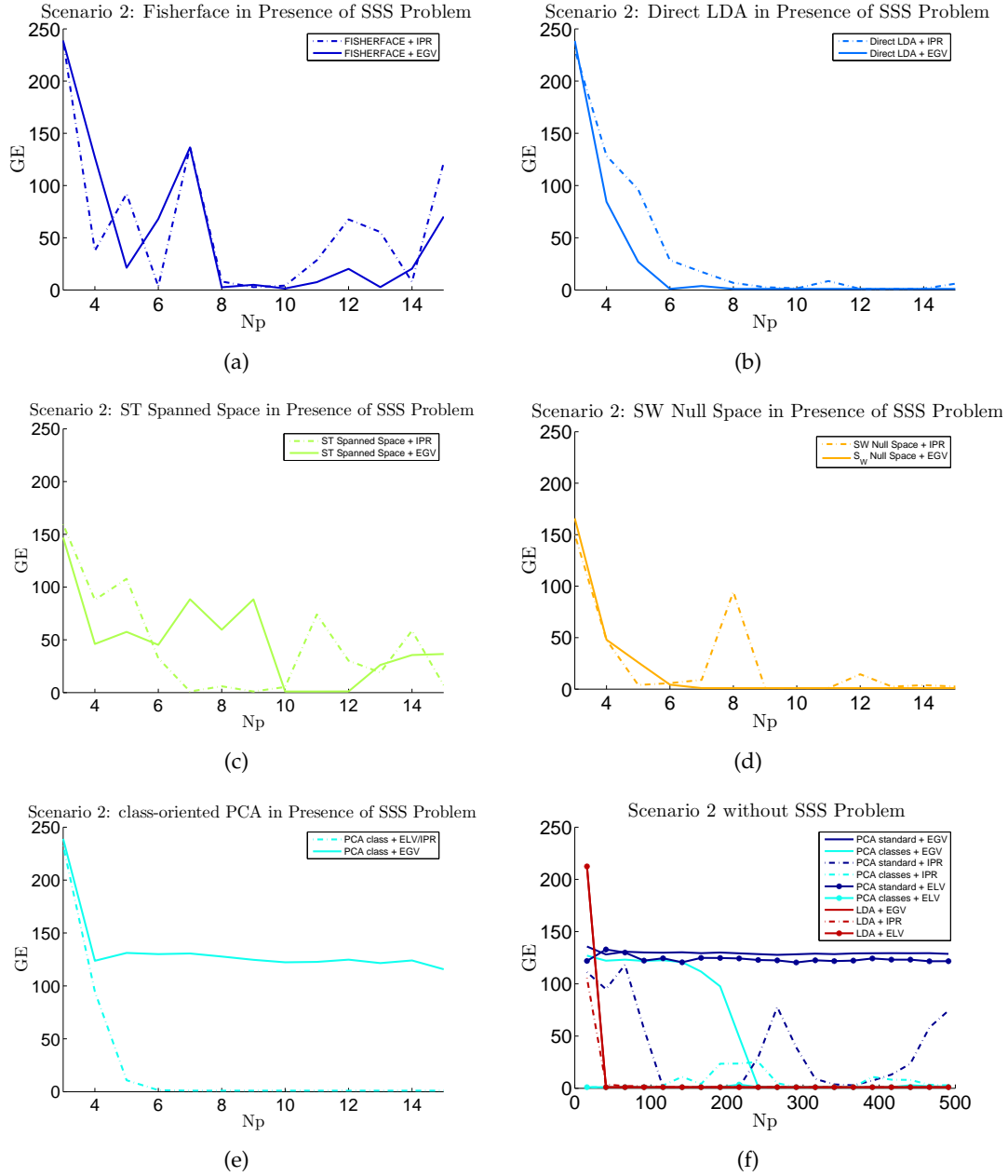


FIGURE 6.4: Guessing Entropy as function of the number of profiling traces. Figures (a)-(d): methods extending the LDA in presence of SSS problem; Figure (e): class-oriented PCA in presence of the SSS problem; Figure (f): number of profiling traces high enough to avoid the SSS problem.

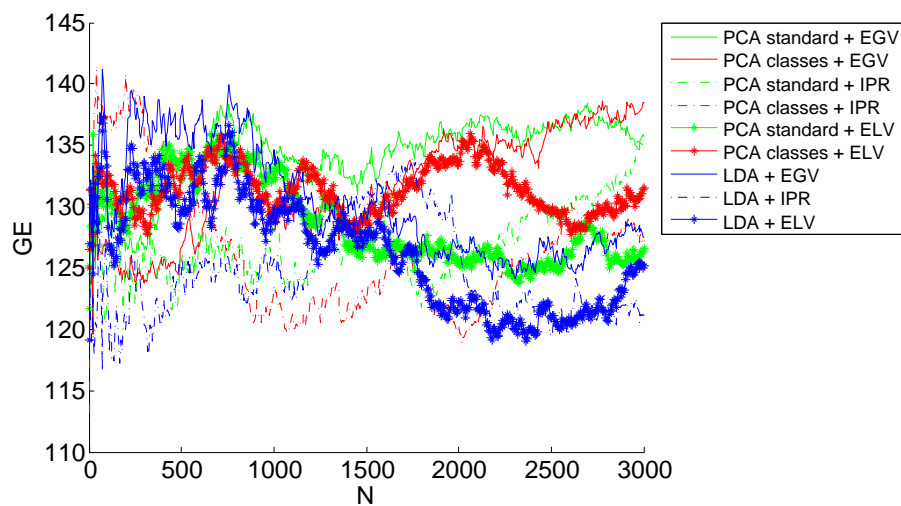


FIGURE 6.5: Degradation of linear-reduction-based template attacks in presence of misalignment.

## Chapter 7

# Kernel Dimensionality Reduction

### 7.1 Motivation

#### 7.1.1 Higher-Order Attacks

Higher-Order Version of Projection Pursuits

### 7.2 Kernel Function and Kernel Trick

#### 7.2.1 Local Kernel Functions as Similarity Metrics

### 7.3 Kernel Discriminant Analysis

### 7.4 Experiments over Atmega328P

#### 7.4.1 The Regularization Problem

#### 7.4.2 The Multi-Class Trade-Off

#### 7.4.3 Multi-Class vs 2-class Approach

#### 7.4.4 Asymmetric Preprocessing/Attack Approach

Comparison with Projection Pursuits

### 7.5 Drawbacks of Kernel Methods

Misalignment Effects

Memory Complexity and Actual Number of Parameters

Two-Phases Approach: Preprocessing-Templates



## Chapter 8

# Convolutional Neural Networks against Jitter-Based Countermeasures

### 8.1 Moving from Kernel Machines to Neural Networks

### 8.2 Misalignment of Side-Channel Traces

#### 8.2.1 The Necessity and the Risks of Applying Realignment Techniques

#### 8.2.2 Analogy with Image Recognition Issues

### 8.3 Convolutional Layers to Impose Shift-Invariance

### 8.4 Data Augmentation for Misaligned Side-Channel Traces

### 8.5 Experiments against Software Countermeasures

### 8.6 Experiments against Artificial Hardware Countermeasures

### 8.7 Experiments against Real-Case Hardware Countermeasures





## **Chapter 9**

# **KDA vs Neural Networks Approach for HO-Attacks**

### **9.1 Simulated Experiment for Profiled HO-Attacks**

#### **9.1.1 The Simulations**

#### **9.1.2 Comparison between KDA and MLP**

### **9.2 Real-Case Experiments over ARM Cortex-M4**



## **Chapter 10**

# **Siamese Neural Networks for Collision Attacks**

### **10.1 Introduction**

### **10.2 Siamese Neural Networks**

#### **10.2.1 Distances and Loss Functions**

#### **10.2.2 Relation with Kernel Machines**

### **10.3 Collision Attacks with Siamese NNs**

#### **10.3.1 Experimental Results**



## **Chapter 11**

# **Conclusions and Perspectives**

### **11.1 Summary**

### **11.2 Strengthen Embedded Security: the Main Challenge for Machine Learning Applications**



## Appendix A

# Scenario 3 and 4 of CARDIS '15 paper

### Scenario 3.

Let  $C$  be now variable and let the other parameters be fixed as follows:  $N_a = 100$ ,  $N_z = 200$ ,  $\#PoI = 3996$ . Looking at Fig. ??, we might observe that the standard PCA might actually well perform in SCA context if provided with a larger number of kept components; on the contrary, a little number of components suffices to the LDA. Finally, keeping more of the necessary components does not worsen the efficiency of the attack, which allows the attacker to choose  $C$  as the maximum value supported by his computational means.

*Remark.* In our experiments the ELV selection method only slightly outperforms the IPR. Nevertheless, since it relies on more sound and more general observations, *i.e.* the maximization of explained variance concentrated into few points, it is likely to be more robust and less case-specific. For example, in Fig. 6.4(f) it can be remarked that while the class-oriented PCA + ELV line keeps constant on the value 0 of guessing entropy, the class-oriented PCA + IPR is sometimes higher than 0.

**Is the table with results overview interesting?**

### Scenario 4.

This is the single scenario in which we allow the ELV selection method to not only select the components to keep but also to modify them, keeping only some coefficients within each component, setting the other ones to zero. We select pairs (*component, time sample*) in decreasing order of the ELV values, allowing the presence of only  $C = 3$  components and  $\#PoI$  different times samples: *i.e.*, we impose that the matrix  $A$  defining the extractor (see (??)) has  $C = 3$  rows (storing the 3 chosen components, transposed) and exactly  $\#PoI$  non-zero columns. Looking at Fig. ?? we might observe that the LDA allows to achieve the maximal guessing entropy with only 1 PoI in each of the 3 selected components. Actually, adding PoIs worsen its performances, which is coherent with the assumption that the vulnerable information leaks in only a few points. Such points are excellently detected by the LDA. Adding contribution from other points raises the noise, which is then compensated by the contributions of further noisy points, in a very delicate balance. Such a behaviour is clearly visible in standard PCA case: the first 10 points considered raise the level of noise, that is then balanced by the last 1000 points.





## **Artificially Simulate Jitter**