

# Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal<sup>1</sup>, M.R. Oliveira<sup>1</sup>, R. Valadas<sup>2</sup>,  
P. Filzmoser<sup>3</sup>, P. Salvador<sup>4</sup> and A. Pacheco<sup>1</sup>

<sup>1</sup>CEMAT - IST - UTL, <sup>2</sup>IT - IST - UTL, <sup>3</sup>TU Wien, <sup>4</sup>IT - UA

March 29, 2012

31<sup>th</sup> Annual International Conference on Computer Communications  
**IEEE INFOCOM 2012**

March 25 -30, 2012  
Orlando, Florida USA

IEEE COMMUNICATIONS SOCIETY

The banner features a tropical sunset background with palm trees. The IEEE logo is at the bottom left, and the IEEE Communications Society logo is at the bottom right.



# Outline

## Introduction

## Outlier Detection PCA

## Feature Selection

Robust mutual information  
Automatic method to select features

## Results

Objects and features  
Performance evaluation

## Conclusions

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

Conclusions



# Introduction

In our work, “**Robust**” means “**Robust**” in the **Statistical** sense

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

3 Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

Conclusions

40

# Introduction

In our work, “**Robust**” means “**Robust**” in the **Statistical** sense

- ▶ A **branch of Statistics** created by Huber and Hampel in the 80's

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

3

## Introduction

Outlier Detection  
PCA

## Feature Selection

Robust mutual information  
Automatic method to select  
features

## Results

Objects and features  
Performance evaluation

## Conclusions

40

# Introduction

In our work, “**Robust**” means “**Robust**” in the **Statistical** sense

- ▶ A **branch of Statistics** created by Huber and Hampel in the 80's
  - ▶ is an **alternative** approach to **Classical Statistics**

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

3

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

Conclusions

40

# Introduction

In our work, “**Robust**” means “**Robust**” in the **Statistical** sense

- ▶ A **branch of Statistics** created by Huber and Hampel in the 80's
  - ▶ is an **alternative** approach to **Classical Statistics**
  - ▶ provides estimators that are **not unduly affected** by **outliers** or **small departures** from model assumptions

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

3

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

Conclusions

40

# Introduction

In our work, “**Robust**” means “**Robust**” in the **Statistical** sense

- ▶ A **branch of Statistics** created by Huber and Hampel in the 80's
  - ▶ is an **alternative** approach to **Classical Statistics**
  - ▶ provides estimators that are **not unduly affected** by **outliers** or **small departures** from model assumptions

## Outlier detection

Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

3

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select features

Results

Objects and features  
Performance evaluation

Conclusions

40

# Introduction

In our work, “**Robust**” means “**Robust**” in the **statistical** sense

- ▶ A **branch of Statistics** created by Huber and Hampel in the 80's
  - ▶ is an **alternative** approach to **Classical Statistics**
  - ▶ provides estimators that are **not unduly affected** by **outliers** or **small departures** from model assumptions

## Feature selection      Outlier detection

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

4

Introduction

Outlier Detection  
PCA

Feature Selection  
Robust mutual information  
Automatic method to select  
features

Results  
Objects and features  
Performance evaluation

Conclusions

40

# Introduction

In our work, “**Robust**” means “**Robust**” in the **statistical** sense

- ▶ A **branch of Statistics** created by Huber and Hampel in the 80's
  - ▶ is an **alternative** approach to **Classical Statistics**
  - ▶ provides estimators that are **not unduly affected** by **outliers** or **small departures** from model assumptions

**Robust**

**Robust**

**Feature selection**

**Outlier detection**

Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

5

Introduction

Outlier Detection  
PCA

Feature Selection  
Robust mutual information  
Automatic method to select features

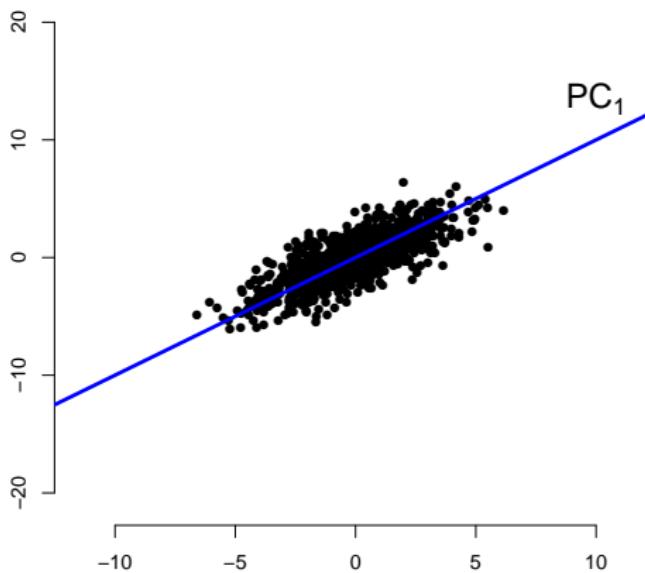
Results  
Objects and features  
Performance evaluation

Conclusions

# Outlier Detection

## PCA

**PC<sub>1</sub>**- Linear combination of the features, with maximum variance



Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

Outlier Detection

PCA

Feature Selection

Robust mutual information  
Automatic method to select features

Results

Objects and features  
Performance evaluation

Conclusions

6

40

# Outlier Detection

## PCA

Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select features

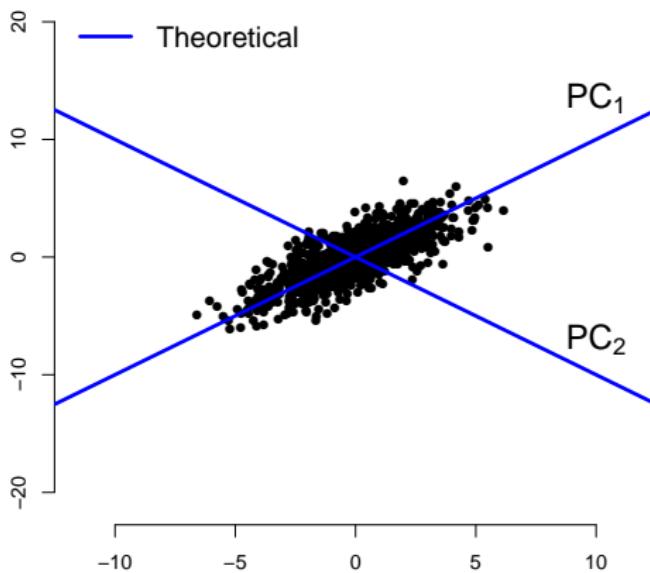
Results

Objects and features  
Performance evaluation

Conclusions

7

**PC<sub>2</sub>** - Orthogonal to **PC<sub>1</sub>**, with maximum variance



40

# Outlier Detection

## PCA

Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

Outlier Detection

PCA

Feature Selection

Robust mutual information  
Automatic method to select features

Results

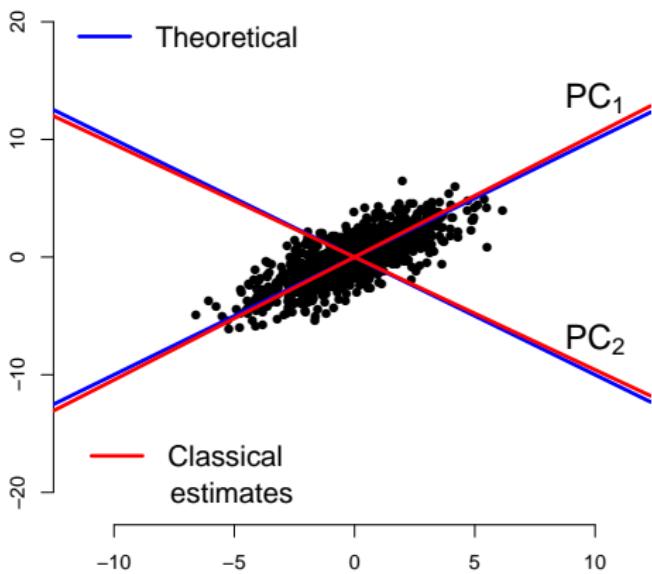
Objects and features  
Performance evaluation

Conclusions

8

40

**Classical estimator:** very sensitive to outliers



# Outlier Detection

## PCA

Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select features

Results

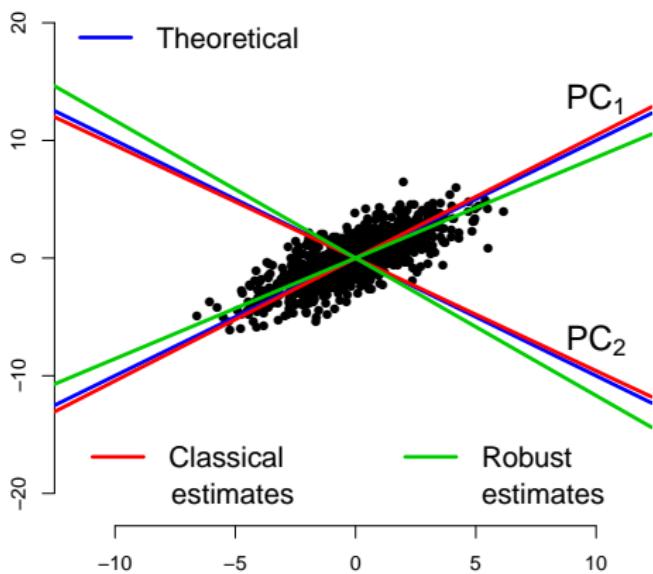
Objects and features  
Performance evaluation

Conclusions

9

40

**Robust** estimator: **insensitive to a small amount of outliers**





# Outlier Detection

## PCA

### Effect of a single outlier:

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

Conclusions

10

40



# Outlier Detection

## PCA

**Toy example: How does PCA detect outliers?**

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

Conclusions

11

40

# Outlier Detection

## PCA

**Main idea:** An **outlier** is a point **very far** from the majority of the data, **projected** into a proper space

Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select features

Results

Objects and features  
Performance evaluation

Conclusions

12

40

# Outlier Detection

PCA

**Main idea:** An **outlier** is a point **very far** from the majority of the data, **projected** into a proper space

## Algorithm

Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

Outlier Detection

PCA

Feature Selection

Robust mutual information  
Automatic method to select features

Results

Objects and features  
Performance evaluation

Conclusions

12

40

# Outlier Detection

## PCA

**Main idea:** An **outlier** is a point **very far** from the majority of the data, **projected** into a proper space

## Algorithm

1. **Estimate** the first PCs and **project** the data into the main directions

Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select features

Results

Objects and features  
Performance evaluation

Conclusions

12

40

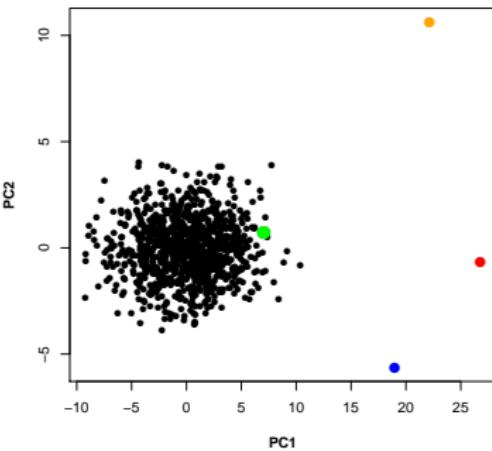
# Outlier Detection

## PCA

**Main idea:** An **outlier** is a point **very far** from the majority of the data, **projected** into a proper space

## Algorithm

1. **Estimate** the first PCs and **project** the data into the main directions



Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select features

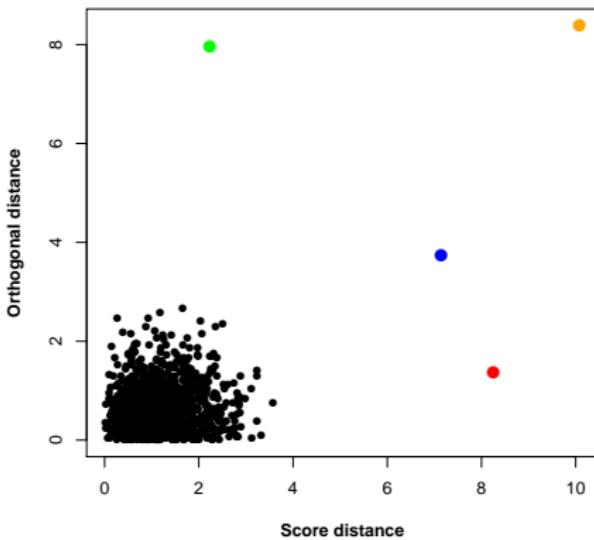
Results

Objects and features  
Performance evaluation

Conclusions

12

40



# Outlier Detection

## PCA

Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

Outlier Detection

PCA

Feature Selection

Robust mutual information  
Automatic method to select features

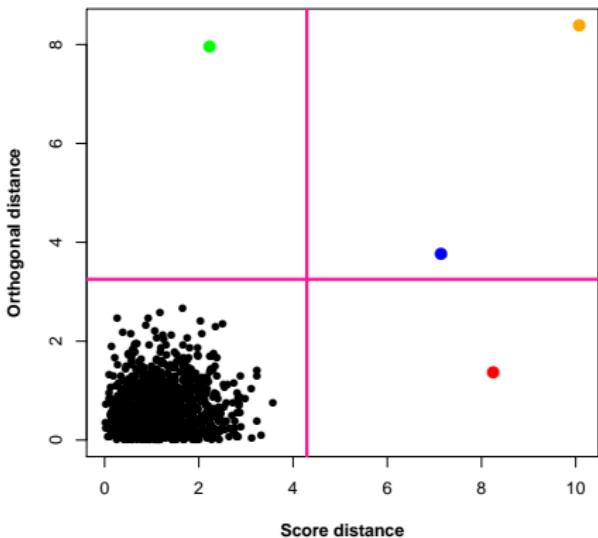
Results

Objects and features  
Performance evaluation

Conclusions

14

40



# Outlier Detection

## PCA

Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

Outlier Detection

PCA

Feature Selection

Robust mutual information  
Automatic method to select features

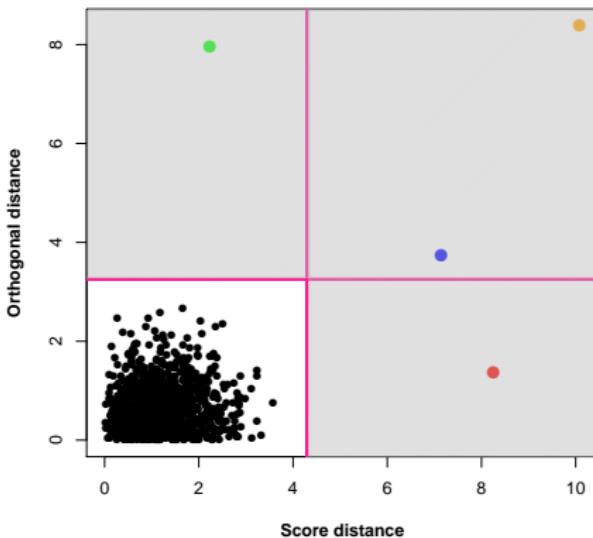
Results

Objects and features  
Performance evaluation

Conclusions

15

40





# Feature Selection

## Filter method

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

### 16 Feature Selection

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

Conclusions



# Feature Selection

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

## 16 Feature Selection

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

Conclusions

## Filter method

- ▶ Selects the **most informative features** to detect anomalies

# Feature Selection

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

## 16 Feature Selection

Robust mutual information  
Automatic method to select  
features

Results  
Objects and features  
Performance evaluation

Conclusions

## Filter method

- ▶ Selects the **most informative features** to detect anomalies
- ▶ Based on a measure of **association between features and classes (anomalous or regular)**

# Feature Selection

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

## 16 Feature Selection

Robust mutual information  
Automatic method to select  
features

Results  
Objects and features  
Performance evaluation

Conclusions

## Filter method

- ▶ Selects the **most informative features** to detect anomalies
- ▶ Based on a measure of **association between features and classes** (**anomalous** or **regular**)
- ▶ We use **Mutual Information**



# Feature Selection

## Mutual Information (MI)

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

17 **Feature Selection**

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

Conclusions

# Feature Selection

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

## 17 Feature Selection

Robust mutual information  
Automatic method to select  
features

Results  
Objects and features  
Performance evaluation

Conclusions

- ▶ Index that measures **linear** and **non-linear dependencies** between features

# Feature Selection

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

17 **Feature Selection**

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

Conclusions

## Mutual Information (MI)

- ▶ Index that measures **linear** and **non-linear dependencies** between features
- ▶ The **MI** between  $X$  and  $Y$  is

# Feature Selection

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

## 17 Feature Selection

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

Conclusions

- ▶ Index that measures **linear** and **non-linear dependencies** between features
- ▶ The **MI** between  $X$  and  $Y$  is

$$MI(X, Y) = \iint f_{XY}(x, y) \ln \left( \frac{f_{XY}(x, y)}{f_X(x) f_Y(y)} \right) dx dy$$

# Feature Selection

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

17 Feature Selection

Robust mutual information  
Automatic method to select  
features

Results  
Objects and features  
Performance evaluation  
Conclusions

## Mutual Information (MI)

- ▶ Index that measures **linear** and **non-linear dependencies** between features
- ▶ The **MI** between  $X$  and  $Y$  is

$$MI(X, Y) = \iint f_{XY}(x, y) \ln \left( \frac{f_{XY}(x, y)}{f_X(x) f_Y(y)} \right) dx dy$$

## Estimation

- ▶ When the joint probability density function  $f_{XY}(x, y)$  is unknown...

# Feature Selection

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

17 Feature Selection

Robust mutual information  
Automatic method to select  
features

Results  
Objects and features  
Performance evaluation  
Conclusions

## Mutual Information (MI)

- ▶ Index that measures **linear** and **non-linear dependencies** between features
- ▶ The **MI** between  $X$  and  $Y$  is

$$MI(X, Y) = \iint f_{XY}(x, y) \ln \left( \frac{f_{XY}(x, y)}{f_X(x) f_Y(y)} \right) dx dy$$

## Estimation

- ▶ When the joint probability density function  $f_{XY}(x, y)$  is unknown...several **estimators** were **proposed!**

# Feature Selection

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

17 Feature Selection

Robust mutual information  
Automatic method to select  
features

Results  
Objects and features  
Performance evaluation  
Conclusions

## Mutual Information (MI)

- ▶ Index that measures **linear** and **non-linear dependencies** between features
- ▶ The **MI** between  $X$  and  $Y$  is

$$MI(X, Y) = \iint f_{XY}(x, y) \ln \left( \frac{f_{XY}(x, y)}{f_X(x) f_Y(y)} \right) dx dy$$

## Estimation

- ▶ When the joint probability density function  $f_{XY}(x, y)$  is unknown...several **estimators** were **proposed!**
- ▶ There is **no single best estimator!**

# Feature Selection

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

17 Feature Selection

Robust mutual information  
Automatic method to select  
features

Results  
Objects and features  
Performance evaluation  
Conclusions

## Mutual Information (MI)

- ▶ Index that measures **linear** and **non-linear dependencies** between features
- ▶ The **MI** between  $X$  and  $Y$  is

$$MI(X, Y) = \iint f_{XY}(x, y) \ln \left( \frac{f_{XY}(x, y)}{f_X(x) f_Y(y)} \right) dx dy$$

## Estimation

- ▶ When the joint probability density function  $f_{XY}(x, y)$  is unknown...several **estimators** were **proposed!**
- ▶ There is **no single best estimator!**
- ▶ **Sensitive**

# Feature Selection

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

17 Feature Selection

Robust mutual information  
Automatic method to select  
features

Results  
Objects and features  
Performance evaluation  
Conclusions

## Mutual Information (MI)

- ▶ Index that measures **linear** and **non-linear dependencies** between features
- ▶ The **MI** between  $X$  and  $Y$  is

$$MI(X, Y) = \iint f_{XY}(x, y) \ln \left( \frac{f_{XY}(x, y)}{f_X(x) f_Y(y)} \right) dx dy$$

## Estimation

- ▶ When the joint probability density function  $f_{XY}(x, y)$  is unknown...several **estimators** were **proposed!**
- ▶ There is **no single best estimator!**
- ▶ **Sensitive**
  - ▶ to **atypical observations** in each **class**
  - ▶ to **mislabeling errors**



# Feature Selection

## Robustification of the MI estimator

## Robustification of the MI estimator

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

### Feature Selection

18

Robust mutual information  
Automatic method to select  
features

### Results

Objects and features  
Performance evaluation

### Conclusions

## Robustification of the MI estimator

**Main idea:** **Exclude** outliers from the data used to estimate **MI**

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

18

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

Conclusions



# Feature Selection

Robustification of the MI estimator

## Robustification of the MI estimator

**Main idea:** **Exclude** outliers from the data used to estimate **MI**

## Algorithm

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

18  
Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

Conclusions

## Robustification of the MI estimator

**Main idea:** **Exclude** outliers from the data used to estimate **MI**

## Algorithm

1. Given a **labeled data set**,  
**separate** in 2 subsets:

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select  
features

Results  
Objects and features  
Performance evaluation

Conclusions

18

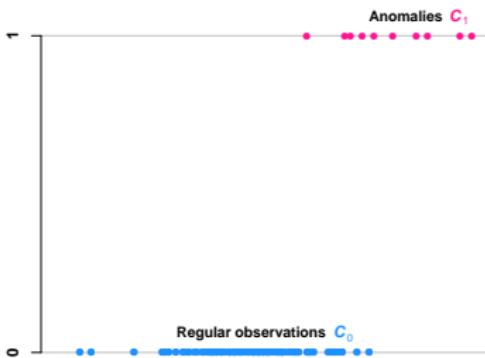
40

## Robustification of the MI estimator

Main idea: **Exclude** outliers from the data used to estimate **MI**

## Algorithm

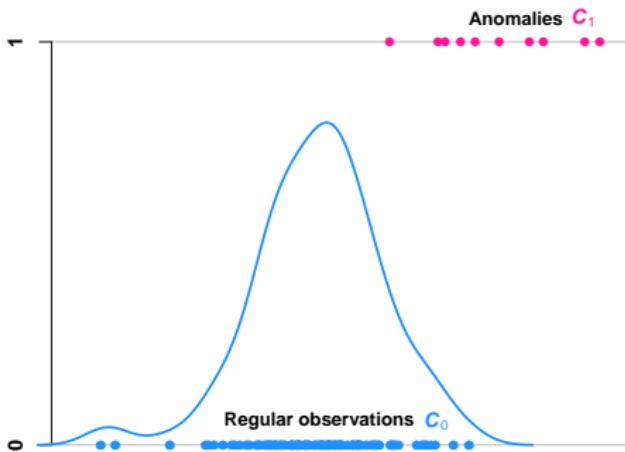
- Given a **labeled data set**,  
**separate** in 2 subsets:
  - $\mathcal{C}_0$  - regular observations (**label 0**)
  - $\mathcal{C}_1$  - anomalies (**label 1**)



# Feature Selection

## Robustification of the MI estimator

2. Assume that **regular observations** have **approximately** normal distribution



Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

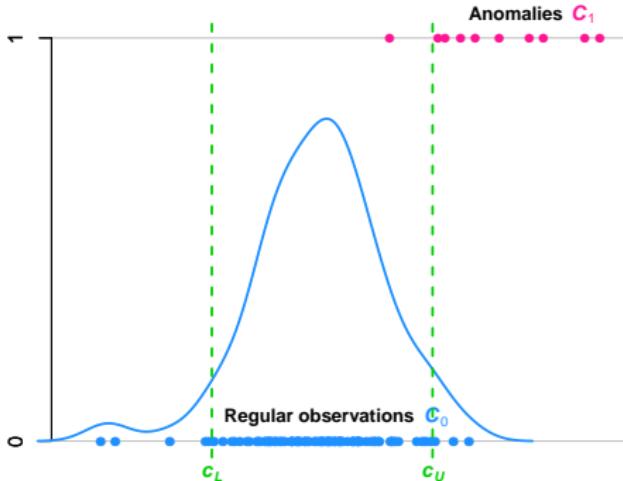
Outlier Detection  
PCA

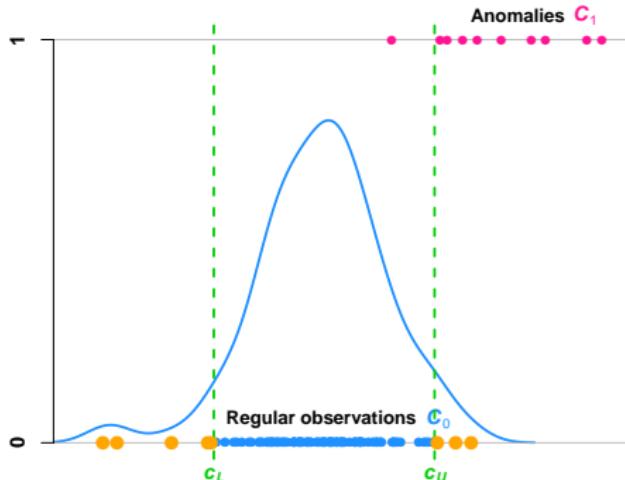
Feature Selection

19  
Robust mutual information  
Automatic method to select features

Results  
Objects and features  
Performance evaluation

Conclusions





# Feature Selection

## Robustification of the MI estimator

Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

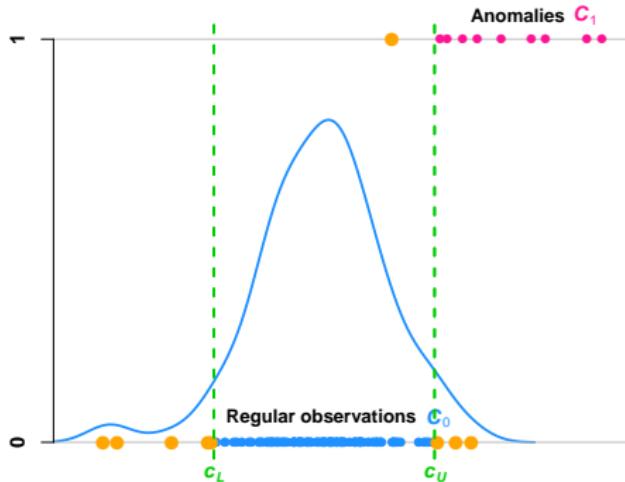
Outlier Detection  
PCA

Feature Selection

22 Robust mutual information  
Automatic method to select features

Results  
Objects and features  
Performance evaluation

Conclusions



# Feature Selection

## Robustification of the MI estimator

Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

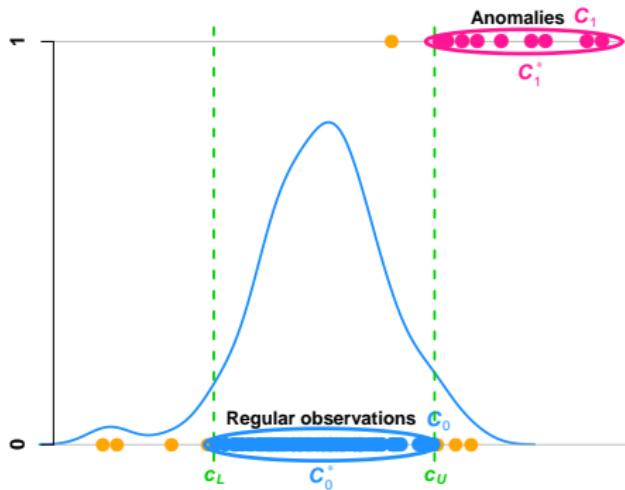
Outlier Detection  
PCA

Feature Selection

23  
Robust mutual information  
Automatic method to select features

Results  
Objects and features  
Performance evaluation

Conclusions





# Feature Selection

Automatic method to select features

## Proposed methods

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

### Feature Selection

Robust mutual information  
Automatic method to select  
features

### Results

Objects and features  
Performance evaluation

### Conclusions

24

40



# Feature Selection

Automatic method to select features

## Proposed methods

- ▶ Number of features **fixed in advance**

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

Conclusions

24

# Feature Selection

Automatic method to select features

## Proposed methods

- ▶ Number of features **fixed in advance**
- ▶ **User-defined threshold** (fixed in advance) for **MI**

24

Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select features

Results

Objects and features  
Performance evaluation

Conclusions

40

# Feature Selection

Automatic method to select features

## Proposed methods

- ▶ Number of features **fixed in advance**
- ▶ **User-defined threshold** (fixed in advance) for **MI**
- ▶ **PCA**, ...

24

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

Conclusions

40

# Feature Selection

Automatic method to select features

## Proposed methods

- ▶ Number of features **fixed in advance**
- ▶ **User-defined threshold** (fixed in advance) for **MI**
- ▶ **PCA**, ...

Automatic  
method  
to  
select  
relevant  
features

24

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

Conclusions

40



# Feature Selection

Automatic method to select features

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

Conclusions

25

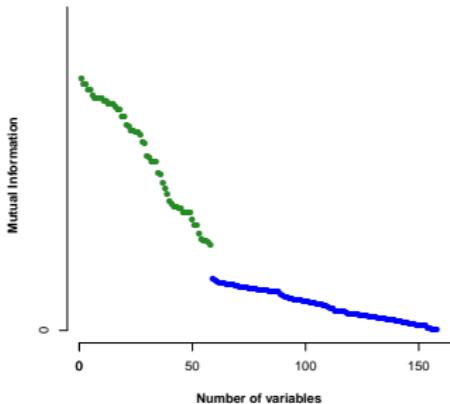
40

# Feature Selection

Automatic method to select features

## Automatic method to select features

- ▶ **Partition** the ordered features in **2 subsets**:
  - ▶ **Relevant features**: highest MI
  - ▶ **Non-relevant features**: others
- ▶ Search for the partition with **highest separation** between subsets (**sharp drop** on the ordered MI)



Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select features

Results

Objects and features  
Performance evaluation

Conclusions

25

40



# Feature Selection

Automatic method to select features

## Algorithm

1. Order the MI estimates,  $mi = (mi_{(1)}, \dots, mi_{(p)})$

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select  
features

26

Results

Objects and features  
Performance evaluation

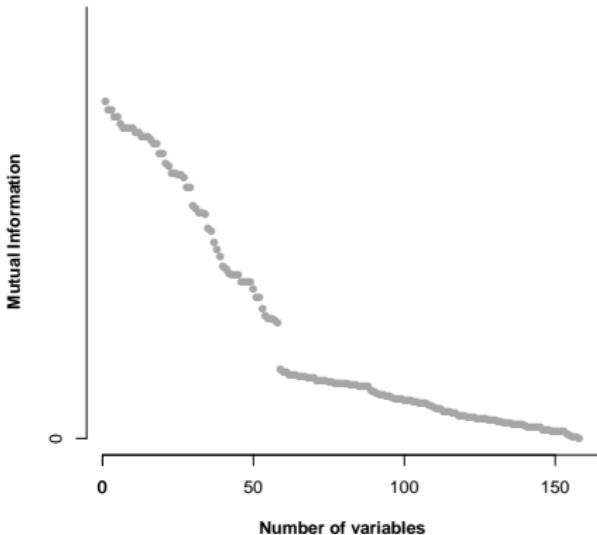
Conclusions

# Feature Selection

Automatic method to select features

## Algorithm

1. Order the MI estimates,  $mi = (mi_{(1)}, \dots, mi_{(p)})$



Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select features

Results  
Objects and features  
Performance evaluation

Conclusions

26

40

# Feature Selection

Automatic method to select features

2. Separate  $mi$  in 2 groups:  $mi_1^{(k)} = (mi_{(1)}, \dots, mi_{(k)})$  and  
 $mi_2^{(k)} = (mi_{(k+1)}, \dots, mi_{(p)})$

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select  
features

27

Results

Objects and features  
Performance evaluation

Conclusions

# Feature Selection

Automatic method to select features

Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

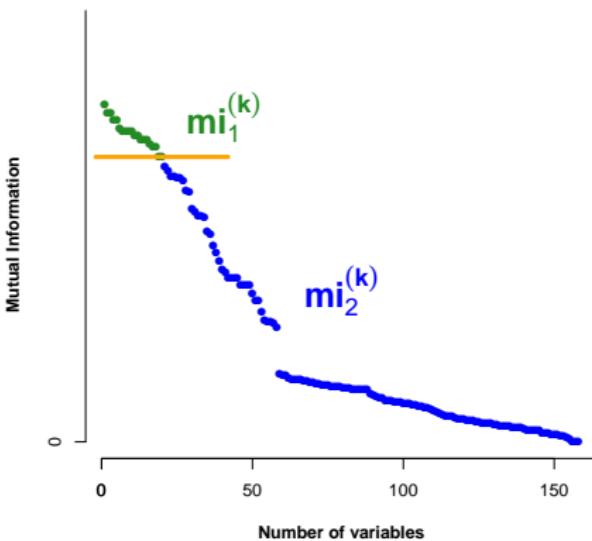
Feature Selection

Robust mutual information  
Automatic method to select features

Results  
Objects and features  
Performance evaluation

Conclusions

27



40

# Feature Selection

Automatic method to select features

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

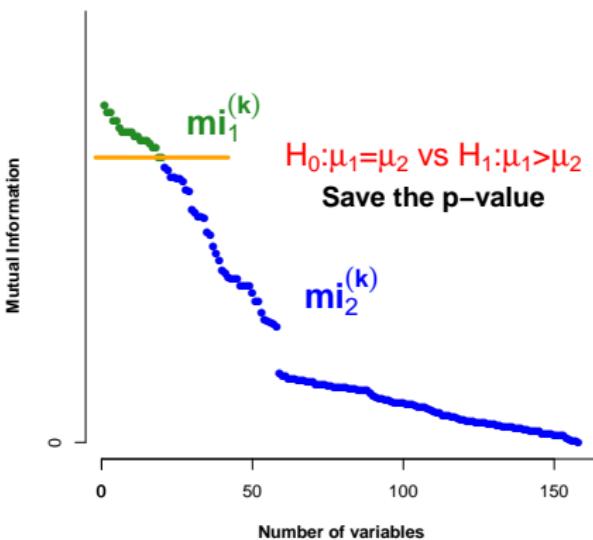
Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

Conclusions

28



40

# Feature Selection

Automatic method to select features

4. Update  $mi_1^{(k+1)}$  and  $mi_2^{(k+1)}$  and repeat step 2., 3. and 4.

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

Conclusions

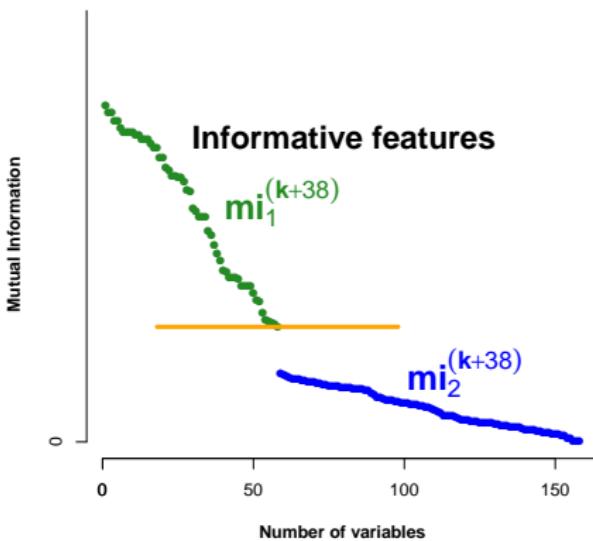
29

40

# Feature Selection

Automatic method to select features

5. Choose  $k$  corresponding to the lowest p-value, highest separation



Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select features

Results  
Objects and features  
Performance evaluation

Conclusions

30

40

# Results

## Objects and features

- ▶ To **evaluate** our method we need a data set with **all traffic objects correctly labeled - perfect ground-truth**

Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select features

Results

31  
Objects and features  
Performance evaluation

Conclusions

# Results

## Objects and features

- ▶ To **evaluate** our method we need a data set with **all traffic objects correctly labeled - perfect ground-truth**
- ▶ **Impossible** to obtain with real data

Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select features

Results

Objects and features  
Performance evaluation

Conclusions

31

40

# Results

## Objects and features

- ▶ To **evaluate** our method we need a data set with **all traffic objects correctly labeled - perfect ground-truth**
- ▶ **Impossible** to obtain with real data



we arranged a **small private laboratory network** with **highly-protected** Internet access  
(CISCO ASA 5510)

Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select features

Results  
Objects and features  
Performance evaluation

Conclusions

31

40



# Results

Objects and features

## Traffic types:

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

Conclusions

32

40



# Results

Objects and features

## Traffic types:

Licit Traffic

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

Conclusions

32

40

# Results

Objects and features

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

Conclusions

32

40

## Traffic types:

### Licit Traffic

- ▶ **10 users** invited to **generate traffic** from inside lab network
- ▶ **Mix** of predominant Internet **applications**:
  - ▶ File sharing (**BitTorrent**)
  - ▶ Video **Streaming**
  - ▶ Web browsing (**HTTP**)

# Results

Objects and features

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

Conclusions

32

40

## Traffic types:

### Licit Traffic

- ▶ **10 users** invited to **generate traffic** from inside lab network
- ▶ **Mix** of predominant Internet **applications**:
  - ▶ File sharing (**BitTorrent**)
  - ▶ Video **Streaming**
  - ▶ Web browsing (**HTTP**)

### Attacks

CEMAT

Dept. of Mathematics,  
Tech. Univ. of Lisbon,  
Portugal

# Results

Objects and features

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

Conclusions

32

40

## Traffic types:

### Licit Traffic

- ▶ **10 users** invited to **generate traffic** from inside lab network
- ▶ **Mix** of predominant Internet **applications**:
  - ▶ File sharing (**BitTorrent**)
  - ▶ Video **Streaming**
  - ▶ Web browsing (**HTTP**)

### Attacks

- ▶ **Produced** only within the lab network
- ▶ **2 broad classes** of **attacks**:
  - ▶ **Port-scan** (NMAP)
  - ▶ **Snapshots** (emulated using FTP)

CEMAT

Dept. of Mathematics,  
Tech. Univ. of Lisbon,  
Portugal



# Results

Dataset, traffic object definition and features

**Dataset:** captured **March 2, 2011, for 8 hours**

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

Conclusions

33

40

# Results

Dataset, traffic object definition and features

**Dataset:** captured **March 2, 2011**, for **8 hours**

## Traffic objects

Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select features

Results

Objects and features  
Performance evaluation

Conclusions

33

40

# Results

Dataset, traffic object definition and features

**Dataset:** captured **March 2, 2011**, for **8 hours**

## Traffic objects

- ▶ **Aggregates** all traffic of one application that **enters or leaves** one machine in 5 minutes (*datastream*)
  - ▶ The **same IP** source address
  - ▶ **One** of the **TCP port** numbers **equal**

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

Conclusions

33

40

# Results

Dataset, traffic object definition and features

**Dataset:** captured **March 2, 2011**, for **8 hours**

## Traffic objects

- ▶ **Aggregates** all traffic of one application that **enters or leaves** one machine in 5 minutes (*datastream*)
  - ▶ The **same IP** source address
  - ▶ **One** of the **TCP port** numbers **equal**

## Features

Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

Outlier Detection

PCA

Feature Selection

Robust mutual information  
Automatic method to select features

Results

Objects and features  
Performance evaluation

Conclusions

33

40

# Results

Dataset, traffic object definition and features

**Dataset:** captured **March 2, 2011**, for **8 hours**

## Traffic objects

- ▶ **Aggregates** all traffic of one application that **enters or leaves** one machine in 5 minutes (*datastream*)
  - ▶ The **same IP** source address
  - ▶ **One** of the **TCP port** numbers **equal**

## Features

- ▶ **5 traffic characteristics** (computed in 0.1 seconds interval):
  - ▶ Nr **packets**: upstream (PUp) and downstream (PDw)
  - ▶ Nr **bytes**: upstream (BUp) downstream (BDw)
  - ▶ Nr **active TCP sessions** (Ses)
- ▶ **8 summary statistics**: min, Q<sub>1</sub>, med, mean, Q<sub>3</sub>, max, sd, MAD

Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

Outlier Detection

PCA

Feature Selection

Robust mutual information  
Automatic method to select features

Results

Objects and features  
Performance evaluation

Conclusions

33

40

# Results

## Performance evaluation

Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

Outlier Detection

PCA

Feature Selection

Robust mutual information  
Automatic method to select features

Results

Objects and features

Performance evaluation

Conclusions

34

Scenario	Detector	Nr Ftrs	Performance measures		
			Recall	FPR	Precision
<b>B1</b>	<b>Ø-NR</b>	-	1	0.063	0.500
	<b>Ø-R</b>	-	1	0.063	0.500
	<b>NR-NR</b>	17	1	0	1
	<b>NR-R</b>	17	1	0	1
	<b>R-NR</b>	5	1	0	1
<b>B2</b>	<b>R-R</b>	5	1	0	1
	<b>Ø-NR</b>	-	0.167	0	1
	<b>Ø-R</b>	-	1	0.679	0.240
	<b>NR-NR</b>	13	0.167	0	1
	<b>NR-R</b>	13	1	0	1
<b>B3</b>	<b>R-NR</b>	7	0.167	0	1
	<b>R-R</b>	7	1	0	1
	<b>Ø-NR</b>	-	0.273	0	1
	<b>Ø-R</b>	-	0.818	0.783	0.333
	<b>NR-NR</b>	18	0.273	0	1
	<b>NR-R</b>	18	0.454	0.522	0.294
	<b>R-NR</b>	7	0.273	0.044	0.750
	<b>R-R</b>	7	1	0	1

# Results

## Performance evaluation

Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

Outlier Detection

PCA

Feature Selection

Robust mutual information  
Automatic method to select features

Results

Objects and features

Performance evaluation

Conclusions

35

Scenario	Detector	Nr Ftrs	Performance measures		
			Recall	FPR	Precision
<b>B1</b>	<b>Ø-NR</b>	-	1	0.063	0.500
	<b>Ø-R</b>	-	1	0.063	0.500
	<b>NR-NR</b>	17	1	0	1
	<b>NR-R</b>	17	1	0	1
	<b>R-NR</b>	5	1	0	1
<b>B2</b>	<b>R-R</b>	5	<b>1</b>	<b>0</b>	<b>1</b>
	<b>Ø-NR</b>	-	0.167	0	1
	<b>Ø-R</b>	-	1	0.679	0.240
	<b>NR-NR</b>	13	0.167	0	1
	<b>NR-R</b>	13	1	0	1
<b>B3</b>	<b>R-NR</b>	7	0.167	0	1
	<b>R-R</b>	7	<b>1</b>	<b>0</b>	<b>1</b>
	<b>Ø-NR</b>	-	0.273	0	1
	<b>Ø-R</b>	-	0.818	0.783	0.333
	<b>NR-NR</b>	18	0.273	0	1
	<b>NR-R</b>	18	0.454	0.522	0.294
	<b>R-NR</b>	7	<b>0.273</b>	<b>0.044</b>	<b>0.750</b>
	<b>R-R</b>	7	<b>1</b>	<b>0</b>	<b>1</b>

# Results

## Performance evaluation

Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

Outlier Detection

PCA

Feature Selection

Robust mutual information  
Automatic method to select features

Results

Objects and features

Performance evaluation

Conclusions

36

Scenario	Detector	Nr Ftrs	Performance measures		
			Recall	FPR	Precision
<b>R1</b>	<b>Ø-NR</b>	-	1	0.033	0.600
	<b>Ø-R</b>	-	1	0.100	0.333
	<b>NR-NR</b>	15	1	0.067	0.429
	<b>NR-R</b>	15	1	0.100	0.333
	<b>R-NR</b>	10	<b>1</b>	<b>0</b>	<b>1</b>
	<b>R-R</b>	10	1	0.017	0.750
<b>R2</b>	<b>Ø-NR</b>	-	0.417	0	1
	<b>Ø-R</b>	-	1	0.600	0.286
	<b>NR-NR</b>	5	1	0.020	0.923
	<b>NR-R</b>	5	1	0.040	0.857
	<b>R-NR</b>	10	<b>1</b>	<b>0</b>	<b>1</b>
	<b>R-R</b>	10	1	0.040	0.857
<b>R3</b>	<b>Ø-NR</b>	-	0.125	0	1
	<b>Ø-R</b>	-	0.875	0.786	0.298
	<b>NR-NR</b>	6	0.125	0	1
	<b>NR-R</b>	6	0.688	0.476	0.355
	<b>R-NR</b>	7	0.063	0	1
	<b>R-R</b>	7	<b>1</b>	<b>0</b>	<b>1</b>

# Results

## Performance evaluation

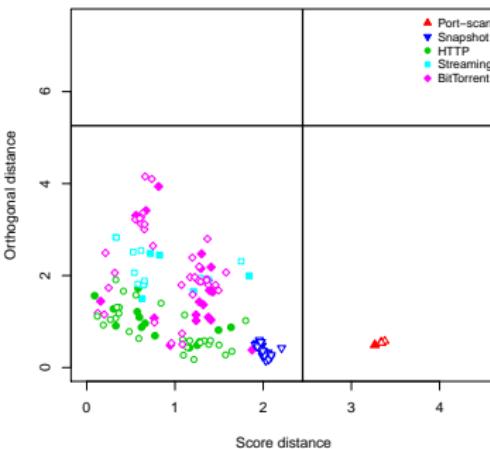
R3

### Detector: R-NR

#### Features: 7

- ▶ **Ses**:  $Q_3$ , max, sd, MAD
- ▶ **BUp**: max, sd
- ▶ **BDw**: mean

Recall	False Positive	Precision
0.125 (0.125)	0 (0)	1 (1)



Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select features

Results  
Objects and features  
Performance evaluation

Conclusions

37

# Results

## Performance evaluation

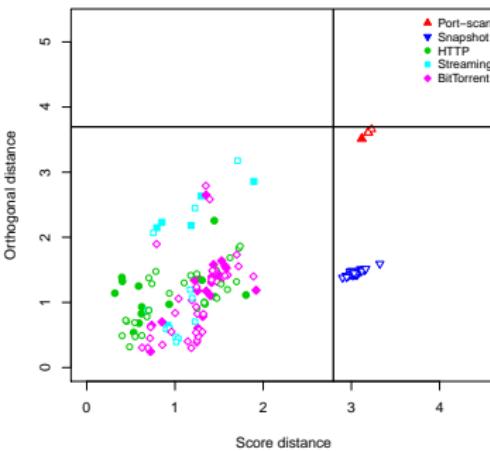
R3

**Detector:** R-R

**Features:** 7

- ▶ **Ses:**  $Q_3$ , max, sd, MAD
- ▶ **BUp:** max, sd
- ▶ **BDw:** mean

Recall	False Positive	Precision
1 (0.875)	0 (0.786)	1 (0.298)



Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select features

Results  
Objects and features  
Performance evaluation

Conclusions

38

40

# Conclusions

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

39 Conclusions

# Conclusions

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

39

Conclusions

# Conclusions

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

- ▶ Feature selection is found to be an **important preprocessing step**, based on **robust MI**, which leads to **better results**
- ▶ The combination of **robust feature selection** and **robust PCA outlier detection**:
  - ▶ achieves **very high performance**
  - ▶ **adaptive** to different traffic conditions

39

Conclusions

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

40

# Conclusions

Robust Feature Selection and Robust PCA for Internet Traffic Anomaly Detection

C. Pascoal

- ▶ Feature selection is found to be an **important preprocessing step**, based on **robust MI**, which leads to **better results**
- ▶ The combination of **robust feature selection** and **robust PCA outlier detection**:
  - ▶ achieves **very high performance**
  - ▶ **adaptive** to different traffic conditions
- ▶ **Robust statistics** brings effective **tools** for anomaly detection
  - ▶ can be **helpful** in addressing other related network problems

39

Conclusions

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select features

Results

Objects and features  
Performance evaluation

40

Robust Feature  
Selection and Robust  
PCA for Internet Traffic  
Anomaly Detection

C. Pascoal

Introduction

Outlier Detection  
PCA

Feature Selection

Robust mutual information  
Automatic method to select  
features

Results

Objects and features  
Performance evaluation

40 Conclusions

# Thank you Obrigada