

# Data Analytics and Models for Insurance

## Presentation of the research chair

Christian ROBERT

**ISFA-COLUMBIA Workshop**  
**Monday June 27, 2016 - Lyon**

---

**2015 - 2020**



CHAIR OF EXCELLENCE

**Data Analytics & Models for Insurance**



**BNP PARIBAS  
CARDIF**



---

**2010 - 2015**

**Management of modelling in  
Life-insurance**



**BNP PARIBAS  
CARDIF**



**Frédéric PLANCHET**

Co-responsable scientifique  
Membre du Laboratoire SAF



Professeur des Universités  
ISFA, Université Lyon 1

**Christian ROBERT**

Co-responsable scientifique  
Directeur du Laboratoire SAF



Professeur des Universités  
ISFA, Université Lyon 1

**Alexis BIENVENUE**

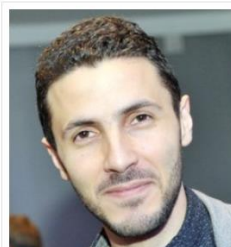
Membre du Laboratoire SAF



Maître de conférences  
ISFA, Université Lyon 1

**Nabil KAZI-TANI**

Membre du Laboratoire SAF



Maître de conférences  
ISFA, Université Lyon 1

**Stéphane LOISEL**

Membre du Laboratoire SAF



Professeur des Universités  
ISFA, Université Lyon 1

**Xavier MILHAUD**

Membre du Laboratoire SAF



Maître de conférences associé  
ISFA, Université Lyon 1

**Didier RULLIERE**

Membre du Laboratoire SAF



Maître de conférences  
ISFA, Université Lyon 1

**Jean-Louis RULLIERE**

Membre du Laboratoire SAF



Professeur des Universités  
ISFA, Université Lyon 1

**Yahia SALHI**

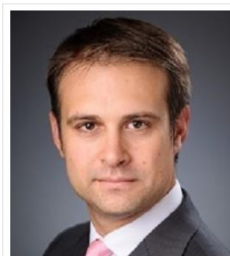
Membre du Laboratoire SAF



Maître de conférences associé  
ISFA, Université Lyon 1

**Pierre THEROND**

Membre du Laboratoire SAF



Enseignant-chercheur associé  
ISFA, Université Lyon 1

**Christophe GEISSLER**

Membre associé



CEO,  
Advestis

**Donatien HAINAUT**

Membre associé



Professeur de Finance  
ESC Rennes

**Bernard BOLLE-REDDAT**

Membre du comité de pilotage



Chief Risks Officer,  
BNP Paribas Cardif

**Michaël DE TOLDI**

Membre du comité de pilotage



Chief Data Officer,  
BNP Paribas Cardif

**Jean-Paul FELIX**

Membre du comité de pilotage



Head of Risk Tools & Processes ,  
BNP Paribas Cardif

**Sébastien CONORT**

Membre associé



Chief Data scientist,  
BNP Paribas Cardif

**Lam DANG**

BNP Paribas Cardif



Data scientist,  
BNP Paribas Cardif


# chaire-dami.fr

Chaire DAMI | BNP Paribas Ca... x +

chaire-dami.fr/fr


Rechercher

OWA ENSAE OWA ISFA ADE Campus Documents ISFA SERVICE VPN UCBL Revues en ligne UCBL Dem compte inv Biblio LSAP Libgen SkyDrive ISFA Travel Planet




CHAIR OF EXCELLENCE

Data Analytics & Models for Insurance



BNP PARIBAS  
CARDIF




LABORATOIRE  
**SAF**  
SCIENCES ACTUARIELLE  
& FINANCIERE

ACCUEIL PROJET ACTEURS EVÉNEMENTS PUBLICATIONS ACTUS Recherche Langue


## Data Analytics & Models for Insurance

**LYON-COLUMBIA  
WORKSHOP**




27 & 28 juin 2016, atelier de  
recherche sur les thématiques de  
la chaire

**RECRUTEMENT  
THESE 2016-2019**



Candidatez !  
« Données incomplètes et  
Apprentissage statistique »

**PETITS DEJEUNERS  
THEMATIQUES**



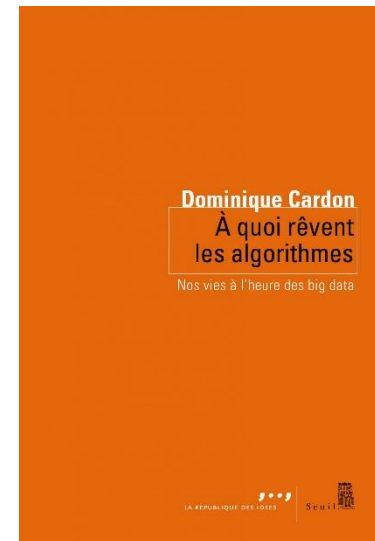
7 juin 2016  
« Market inconsistencies » par N.  
EL KAROUI & J. VEDANI

15:50  
02/06/2016



March 15, 2016

## Seminar – Breakfast – « Politics of algorithms » by Dominique Cardon



June 7, 2016

## Seminar – Breakfast– « Market inconsistencies » by Nicole El Karoui & Julien Védani



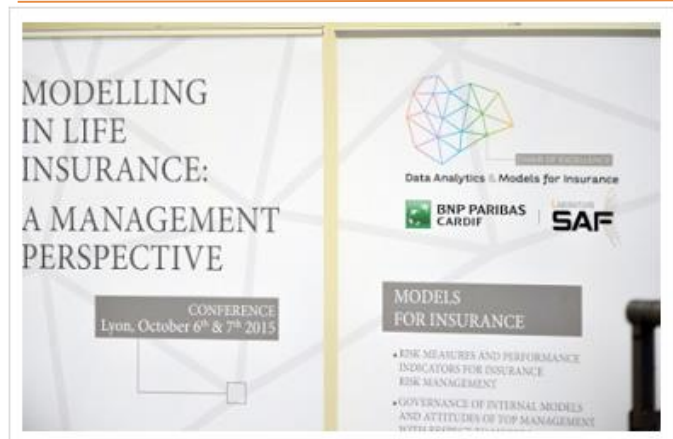
### March 23, 2016 – Topics

- Market inconsistencies of the market-consistent European life insurance economic valuations
- Proxys for SII
- Impact of volatility clustering on equity indexed annuities
- Assessment of beneficiary clauses in free text via Text Mining
- Optimization of treatment of web leads queue with scoring and simulation
- The experiments for observation of human behaviors

### March 25 2015 – Topics

- Credit Losses Impairment
- Agents attitudes towards risk and models: Study of a new analysis and comparison
- Asymmetry & Big Data : which impact for insurance ?
- Working group on the risk-neutral approach
- Longevity risk
- Financial information and Risk in insurance : Change for the better and for worse
- Kaggle AXA competition : methodology of the research lab





**October 6 & 7, 2015**



*David INGRAM (Willis Re) « Bridging the gap between managers and models »*

*Bernard BOLLE-REDDAT (BNP Paribas Cardif) « Management and models »*

*Clément PETIT – Guillaume ALABERGÈRE (ACPR) « Validation in life modelling, a supervisory point of view »*

*Antoon PELSSER (Maastricht University) « The difference between LSMC and replicating portfolio in insurance liability modelling »*

*Michaël SCHMUTZ (FINMA) « Group solvency tests, intragroup transfers and intragroup diversification: A set-valued perspective »*

*Georges DIONNE (HEC Montréal) « Governance of risk management »*

*Thomas BREUER (FHV) « Systemic stress testing and model risk »*

*Andreas TSANAKAS (Cass Business School) « Model risk & culture »*

*Michaël de TOLDI (BNP Paribas Cardif) « Governance for data & analytics in insurance »*

A pair of hands holds a metal caliper, measuring the width of the word "Risk" which is printed in a large, bold, black font. The caliper's jaws are positioned on either side of the word, and the measurement is taken across its full width. The background is plain white.

[illegible]

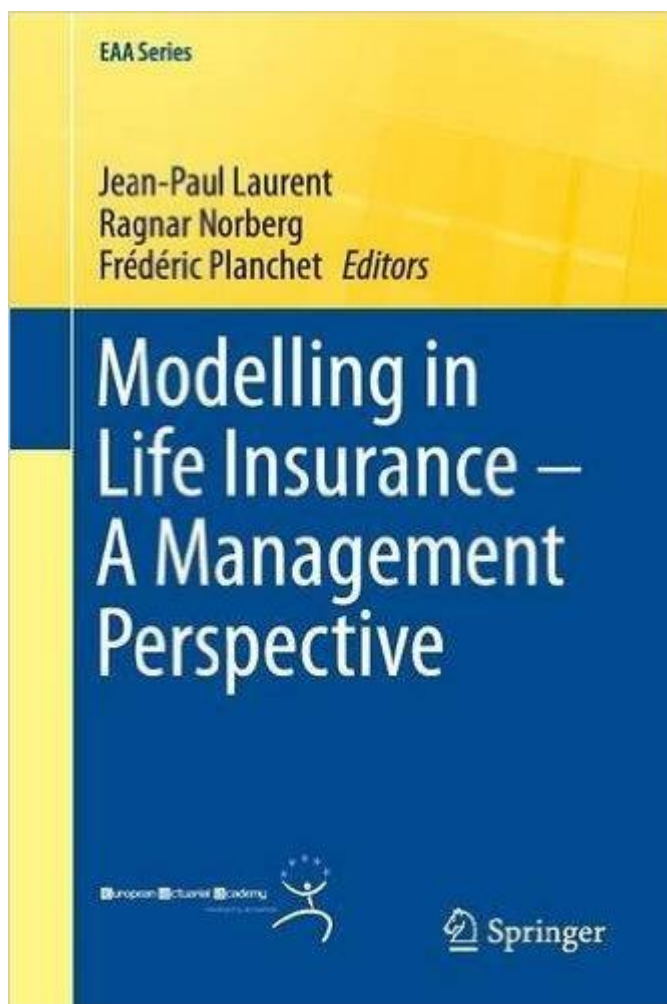
**the impact of the regulatory and accounting environment on  
their development and management**



Figure 1 is a scatter plot illustrating two clusters of data points. The x-axis and y-axis both range from 0.0 to 1.0. The left cluster, represented by blue dots, has a centroid (orange dot) at approximately (0.25, 0.35). The right cluster, represented by red dots, has a centroid (orange dot) at approximately (0.65, 0.75). Lines connect each data point to its respective centroid.

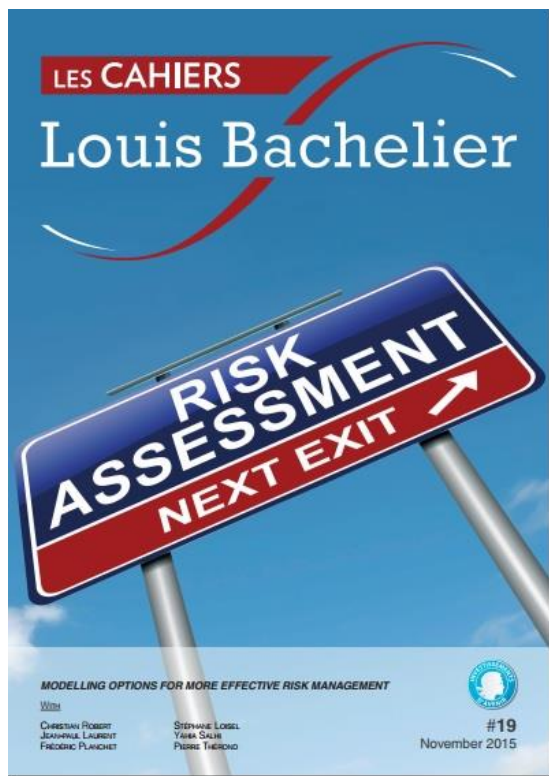
## Proxies, model points and advanced simulation techniques for risk management





## Contents

- 1- Paradigms in life insurance
- 2- About market consistent valuation in insurance
- 3- Cash flow projection models
- 4- Economic scenario generators
- 5- From internal to ORSA models
- 6- Building a model: practical implementation
- 7- Ex-ante model validation and back testing
- 8- The threat of model risk for insurance companies
- 9- Meta-models and consistency issues
- 10- Model feeding & Data Quality
- 11- The role of models in management decision making
- 12- Models and behavior of stakeholders



## Les cahiers de l'ILB – #19 – November 2015

### INDEX

Can ambiguity affect risk reduction?

*Based on an interview with Christian Robert*

Does Basel III succeed in harmonizing the measurement of credit risk?

*Based on an interview with Jean-Paul Laurent*

Valuation of life insurance: how is volatility to be measured?

*Based on the works of Frédéric Planchet*

Risk management: defining an area rather than a threshold

*Based on an interview with Stéphane Loisel*

Insurance: how can sudden changes in the frequency of claims or the intensity of mortality be detected?

*Based on an interview with Yahia Salhi*

IFRS: how are the optimal impairment parameters to be defined?

*Based on an interview with Pierre Thérond*

# Experiments in the lab

Experimental Economics is a branch of economics that focuses on individual behavior in a controlled laboratory setting or out in the field.

Experimental economics helps to prove or disprove economic theories and create predictions and insights about real-world behavior.



## hroot

Hamburg registration and organization online tool

Vous n'avez pas encore de compte? Enregistrez-vous maintenant pour participer aux expérimentations économiques.

[S'enregistrer maintenant](#)

### Vous avez déjà un compte?

Did you already have an account for the registration system? In this case, your data was imported to this system. You can activate your account in the new system here.

[Activez votre compte maintenant !](#)



## Connexion

pour utilisateurs enregistrés

Email

Password

☐ remember me

[Connexion](#)

## WHAT DO WE STUDY ?

- Individual choices  
(choosing under risk, arbitrage, intertemporal choice ...)
- Strategic interactions  
(Negotiation, conflict, contract, incentives, ...)
- Market designs  
(trade efficiency, public good provision, market design ...)

**Privacy concerns, data  
anonymization, open data**



## **Data analytics in insurance**



**Governance for data analytics, new business  
models with big data and analytics**



**Risk-based pricing, predictive  
analytics, machine learning**

# Traineeship: Textual analysis of published and working paper in Machine Learning research

1. Identification of the leading Machine Learning research journals
2. Recovery of titles, abstracts, names of authors and their affiliations
3. Creation of a text-mining tool identifying the key issues and key research center
4. Creation of a visualization tool and mapping of research in Machine Learning in the world
5. Identification of subjects with potential applications for insurance

 Browse Volumes & Issues

International Journal of Machine Learning and Cybernetics

ISSN: 1868-8071 (Print) 1868-808X (Online)

[Browse Volumes & Issues](#)



Latest Articles

 Original Article

Multi-criteria decision-making based on generalized prioritized aggregation operators under simplified neutrosophic uncertain linguistic environment

Zhang-peng Tian, Jing Wang, Hong-yu Zhang... (June 2016)

[» Look Inside](#)

[» Get Access](#)

Available	Volumes
2010 - 2016	7
Issues	Articles
30	548

```
@Article{Tian2016,
author="Tian, Zhang-peng
and Wang, Jing
and Zhang, Hong-yu
and Wang, Jian-qiang",
title="Multi-criteria decision-making based on generalized
prioritized aggregation operators under simplified neutrosophic
uncertain linguistic environment",
journal="International Journal of Machine Learning and
Cybernetics",
year="2016",
pages="1--17",
abstract="A simplified neutrosophic uncertain linguistic set that
integrates quantitative and qualitative evaluation can serve as
an extension of both an uncertain linguistic variable and a
simplified neutrosophic set. It can describe the real preferences
of decision-makers and reflect their uncertainty, incompleteness
and inconsistency. This paper focuses on multi-criteria decision-
making (MCDM) problems in which the criteria occupy different
priority levels and the criteria values take the form of
simplified neutrosophic uncertain linguistic elements. Having
reviewed the relevant literatures, this paper develops some
generalized simplified neutrosophic uncertain linguistic
prioritized weighted aggregation operators and applies them to
solve MCDM problems. Finally, an illustrative example is given,
and two cases of comparison analysis are conducted with other
representative methods to demonstrate the effectiveness and
feasibility of the developed approach.",
issn="1868-808X",
doi="10.1007/s13042-016-0552-9",
url="http://dx.doi.org/10.1007/s13042-016-0552-9"
}
```



# Incomplete data, Machine Learning and Insurance

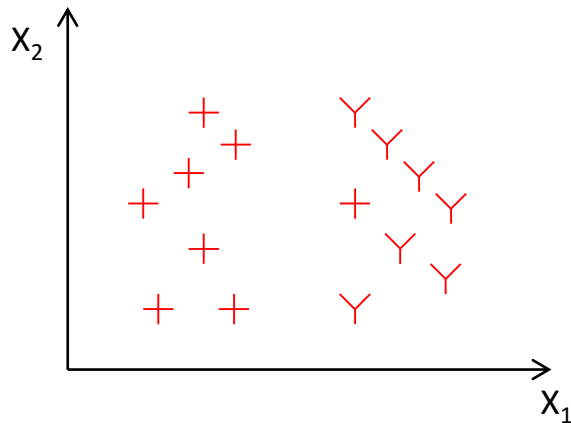
A research project on data science

Christian ROBERT

**ISFA-COLUMBIA Workshop**  
**Monday June 27, 2016 - Lyon**

# Data types

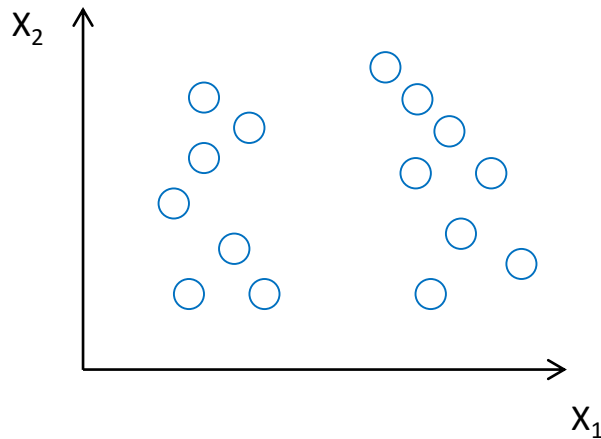
Labeled data



Data: (Y , X)



Unlabeled data



(X)

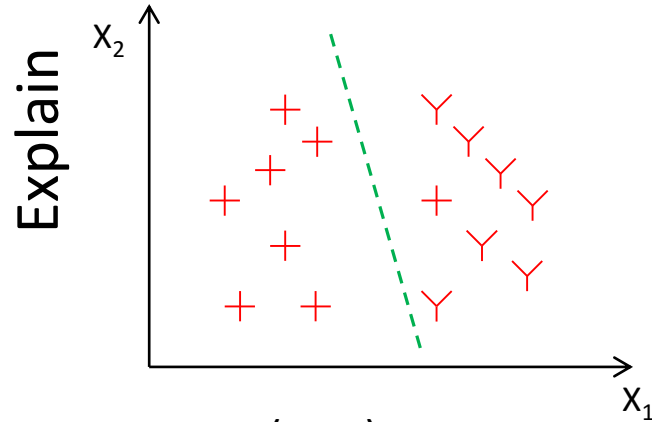


Y : labels = + or Y , response variable, output variable

X : explanatory variables, input variables, covariates, independent variables, control variables, features,...

# Data to be explained and/or to be predicted

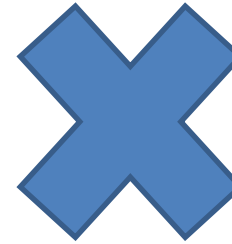
Train



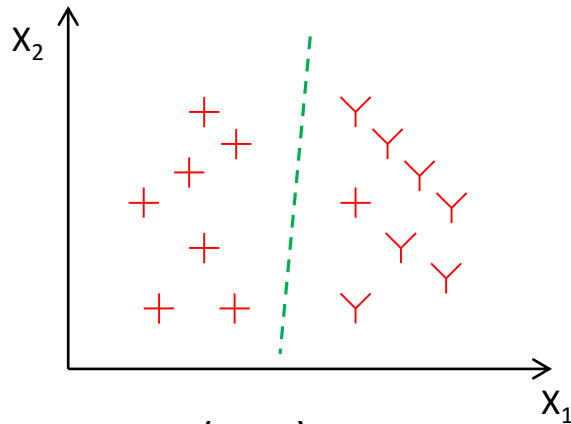
Data:

(Y, X)

Test



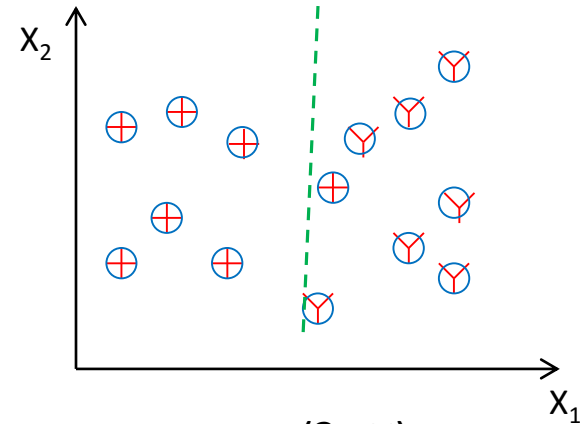
Predict



Data:

(Y, X)

and

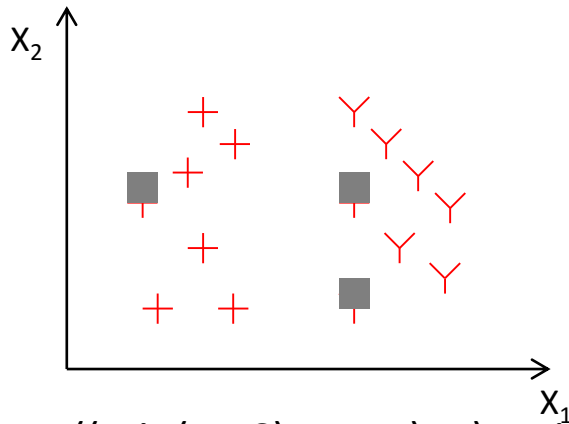


(?, X)



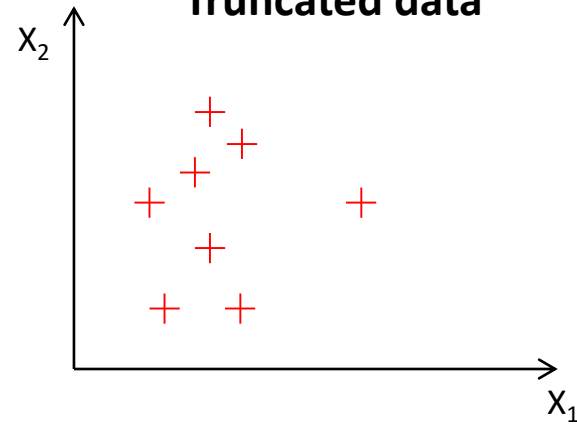
# Imperfect labeled data

**Censored data**



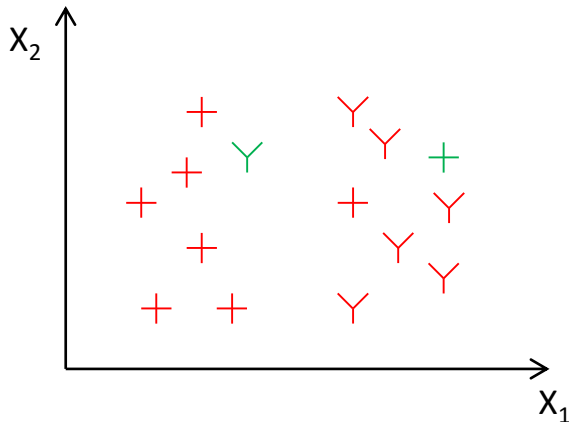
Data:  $((\min(Y, C), 1_{Y > C}), X)$  with  $Y \perp C$

**Truncated data**



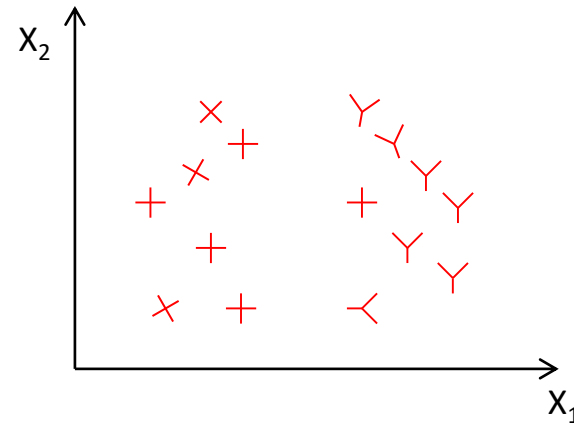
$((Y, C, X) | Y > C)$  with  $Y \perp C$

**Random wrong label**



Data:  $(Y^* = Y 1_{\epsilon=1} + Y^{\wedge} 1_{\epsilon=-1}, X)$  with  $Y \perp Y^{\wedge}$

**Noisy labeled data with endogenous errors**



$(Y^* = Y + \epsilon, X)$  with  $\epsilon \perp X$

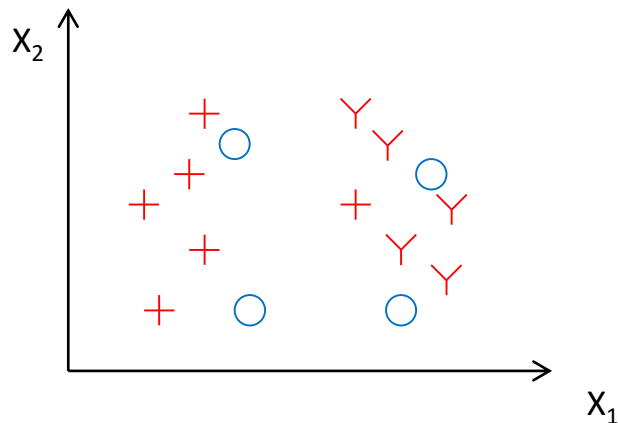


Only probabilistic schemes?

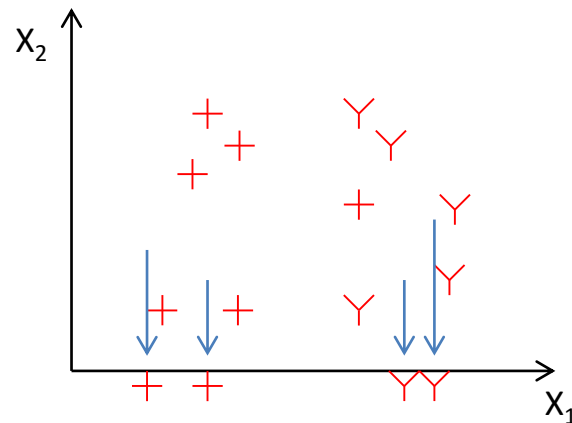


# Labeled with unlabeled data / Missing values

Some labels  $Y$  are not observed



Some components of  $X$  are not observed



**Missing completely at random**

$$(Y^* = Y 1_{\epsilon=1} + \emptyset 1_{\epsilon=-1}, X) \in \perp X$$

$$(Y, X^* = X 1_{\epsilon=1} + \emptyset 1_{\epsilon=-1}) \in \perp X$$

**Missing at random**

$$(Y^* = Y 1_{\epsilon=1} + \emptyset 1_{\epsilon=-1}, X) \in \not\perp X$$

$$(Y, X^* = X 1_{\epsilon=1} + \emptyset 1_{\epsilon=-1}) \in \not\perp X$$

**Missing not a random**

$$(Y^* = Y 1_{Y < c} + \emptyset 1_{Y > c}, X)$$

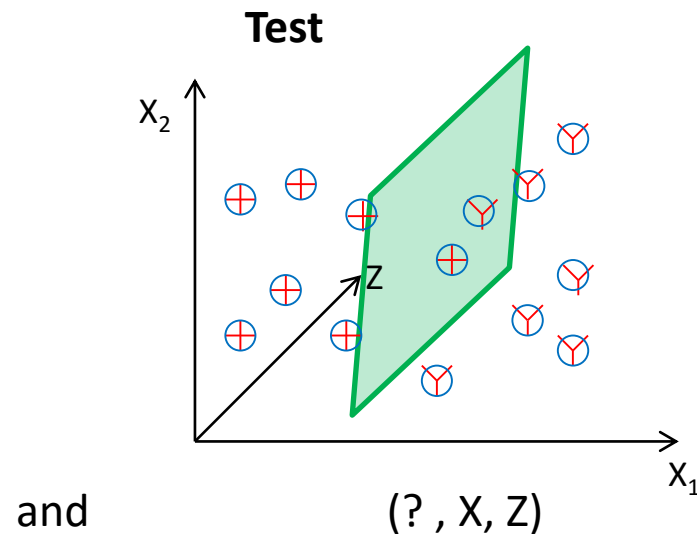
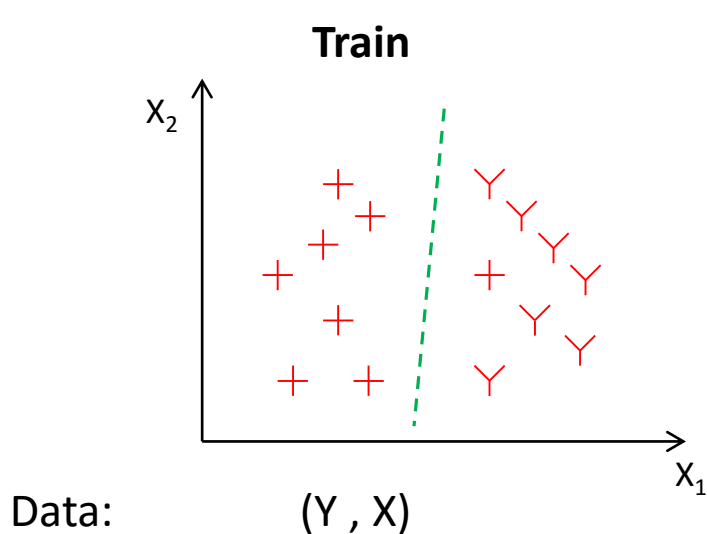
$$(Y, X^* = X 1_{Y < c} + \emptyset 1_{Y > c})$$



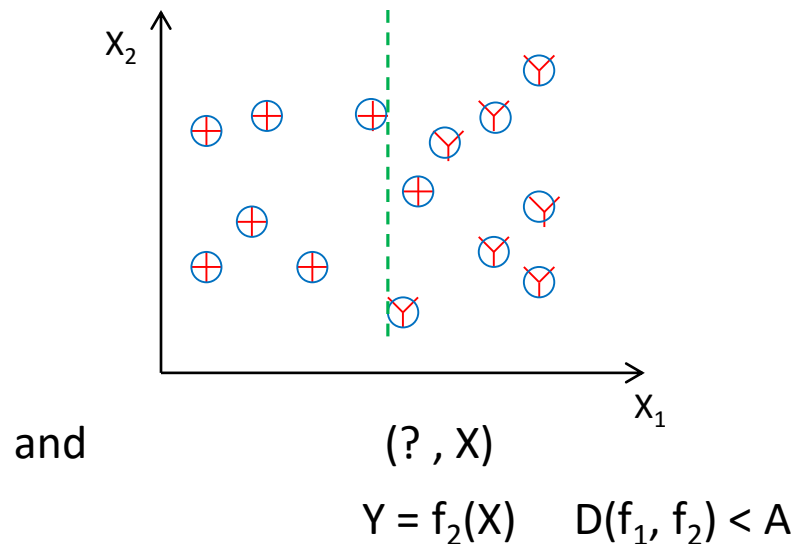
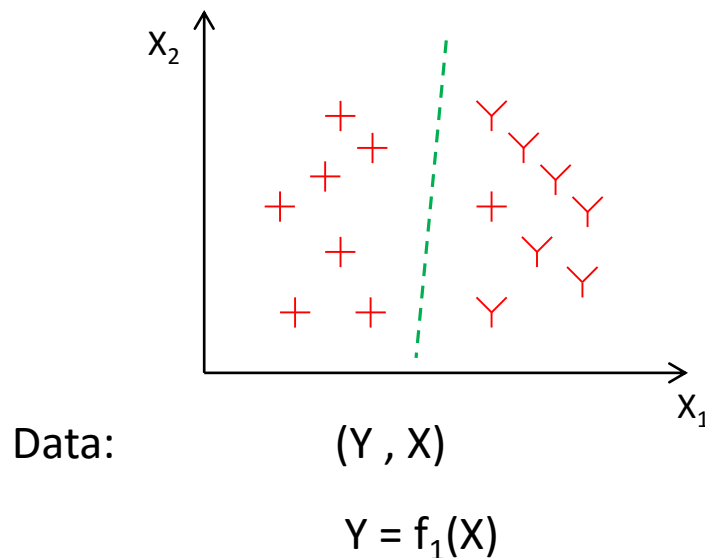


# When train and test data bases differ

**Predict  
controlled  
data**

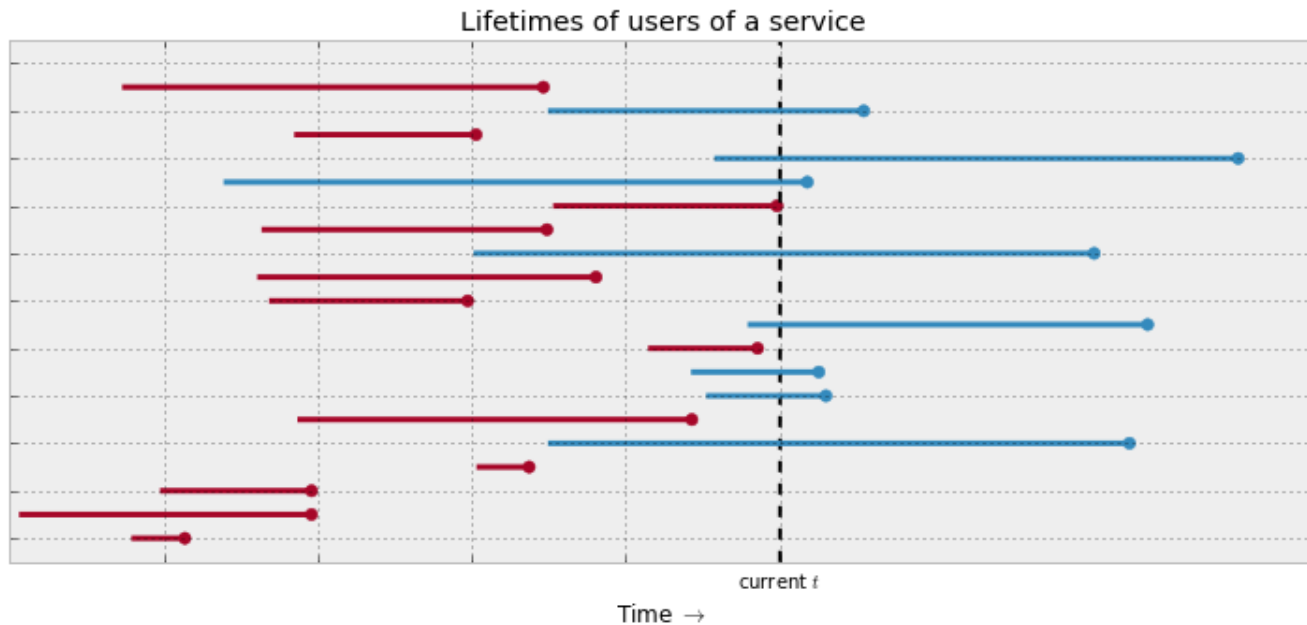
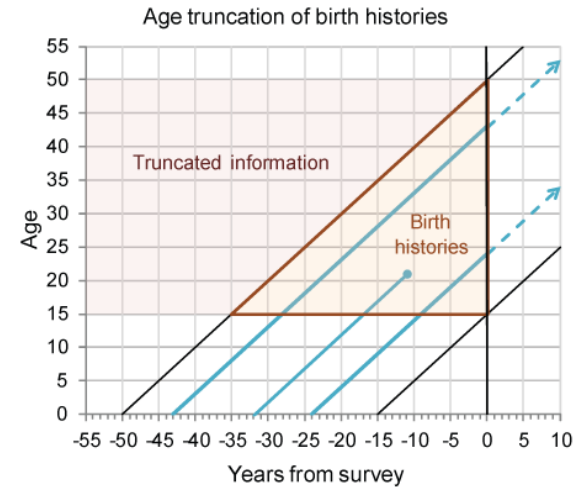
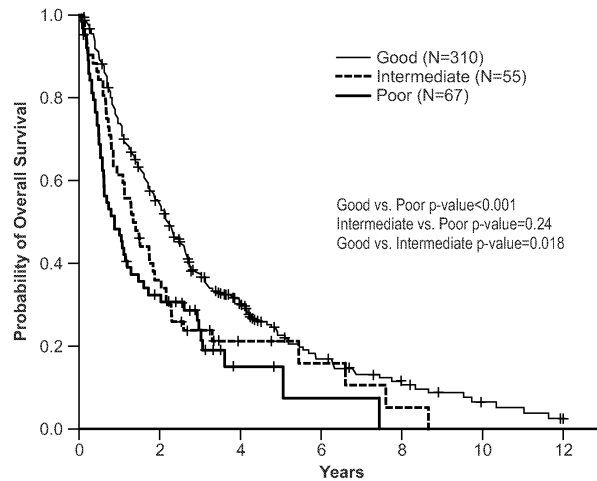


**Predict  
test data  
with a  
different  
generating  
process**



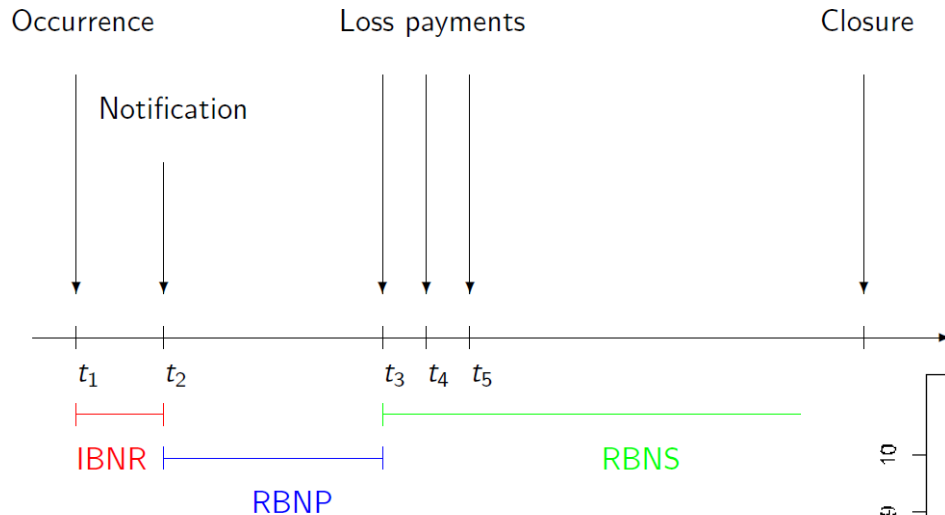
# Mining imperfect data in insurance

## Truncated / censored data

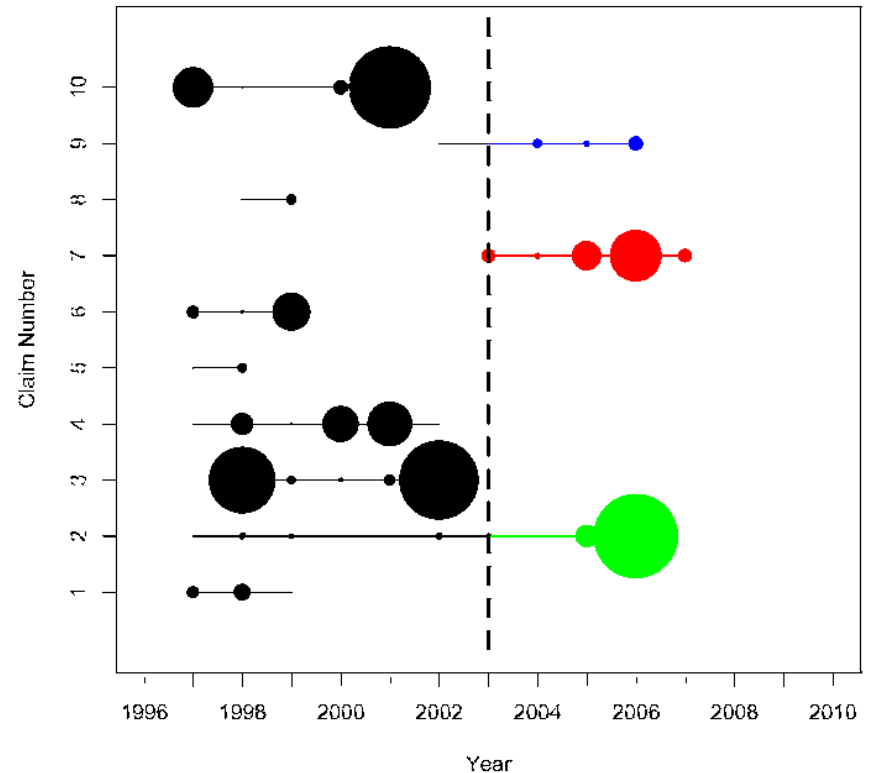


# Mining imperfect data in insurance

## Individual claim process



Incurred But Not Reported (IBNR) claims  
Reported But Not Paid (RBNP) claims  
Reported But Not Settled (RBNS) claims



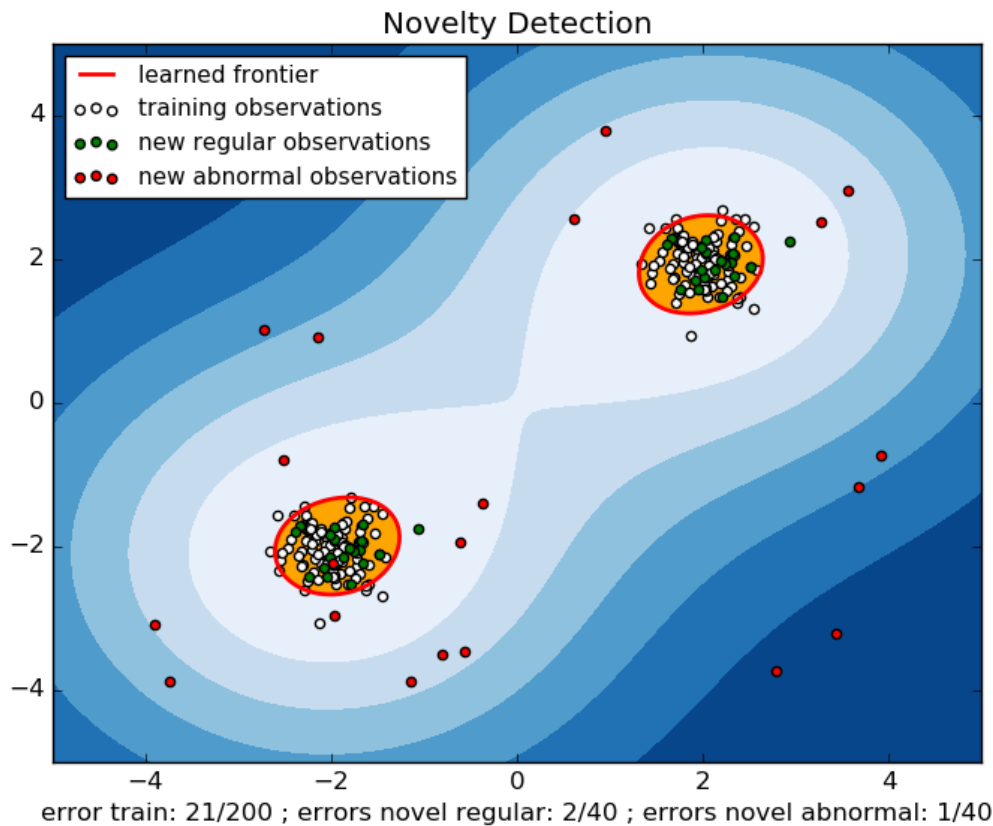
# Mining imperfect data in insurance

Insurance products with several generations of policies / customers



# Mining imperfect data in insurance

## Novelty / Fraud detection





# Machine Learning vs Statistics/Econometrics

## Subfields

**Machine Learning** is a subfield of computer science and artificial intelligence which deals with building systems that can learn from data, instead of explicitly programmed instructions.

**Statistical Modelling** is a subfield of mathematics which deals with finding relationship between variables to predict an outcome

## Data mechanism/data generating process

**Machine Learning** uses algorithmic models and treats the data mechanism as unknown.

**Statistical Modelling** assumes that the data are generated by a given stochastic data model.

## Model choice

**Machine Learning** focuses on Predictive Accuracy even in the face of lack of interpretability of models. Model Choice is based on Cross Validation of Predictive Accuracy using Partitioned Data Sets.

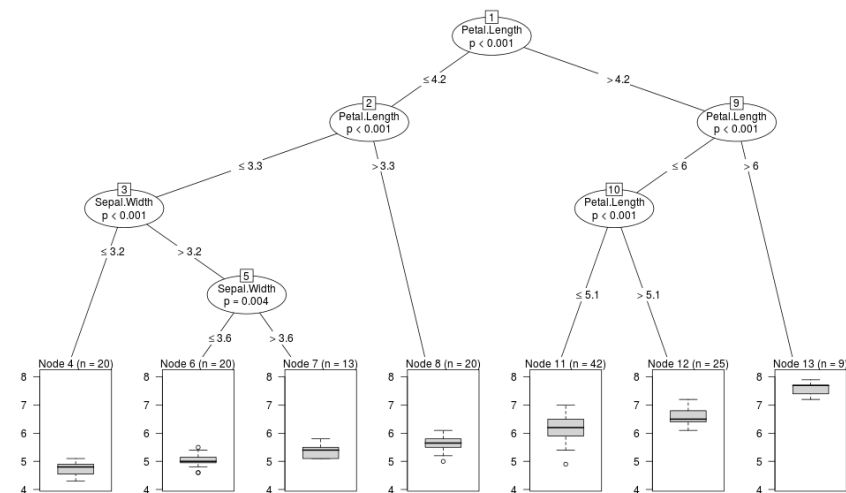
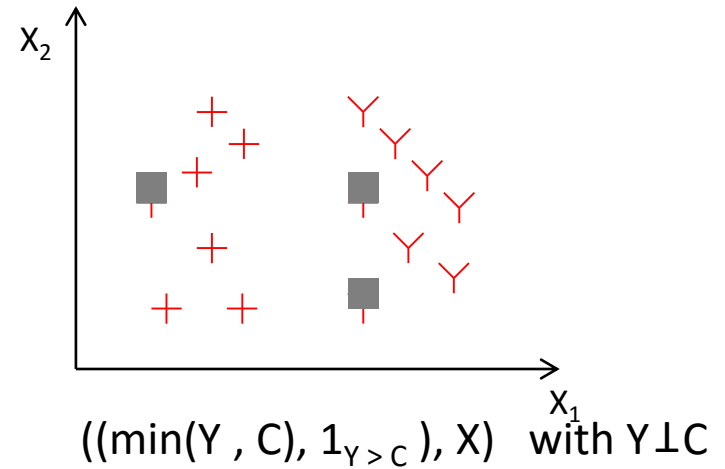
**Statistical Modelling** focuses on hypothesis testing of causes and effects and interpretability of models. Model Choice is based on parameter significance and/or confidence intervals, and In-sample Goodness-of-fit.

# Tree-based censored regression/Survival random forest

- Random forests have been extended to the survival context by Ishwaran et al. (2008), who prove consistency of **Random Survival Forests** (RSF) algorithm assuming that all variables are categorical.

- Yang et al. (2010) showed that by incorporating **kernel functions into RSF**, their algorithm KIRSF achieves better results in many situations.

- Lopez et al. (2015) used an approach that is based on the **IPCW strategy** (Inverse Probability of Censoring Weighting") and that consists in determining a **weighting scheme** that compensates the lack of complete observations in the sample.



# One-class classification

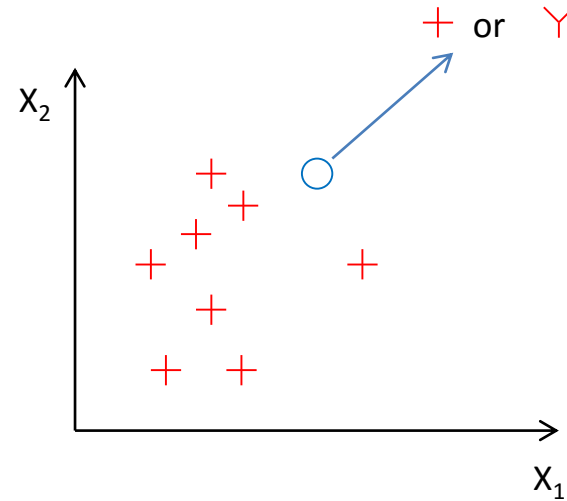
One-class classification tries to identify objects of a specific class amongst all objects, by learning from a **training set containing only the objects of that class**.

It is also known as Outlier detection, Novelty detection, Concept learning, Single class classification, or Unary classification.

*An example is the automatic diagnosis of a disease. It is relatively easy to compile positive data (all patients who are known to have a 'common' disease) but negative data may be difficult to obtain since other patients in the database cannot be assumed to be negative cases if they have never been tested, and such tests can be expensive.*

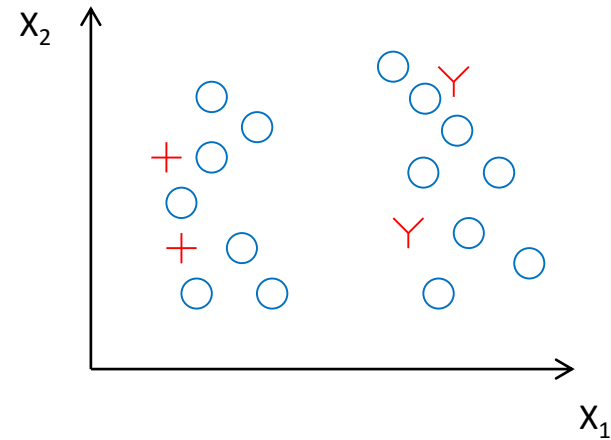
## Algorithms that can be used

- One-class Support Vector Machines (OSVMs)
- Neural networks
- Decision trees
- Nearest neighbors



# Semi-supervised learning

It is a class of supervised learning tasks and techniques that also make use of unlabeled data for training – **typically a small amount of labeled data with a large amount of unlabeled data.**



**Goal:** Using both labeled and unlabeled data to build better learners, than using each one alone.

In order to make any use of unlabeled data, it is implicitly assumed some structure to the underlying distribution of data: Smoothness assumption, Cluster assumption, Manifold assumption.

## Algorithms that can be used

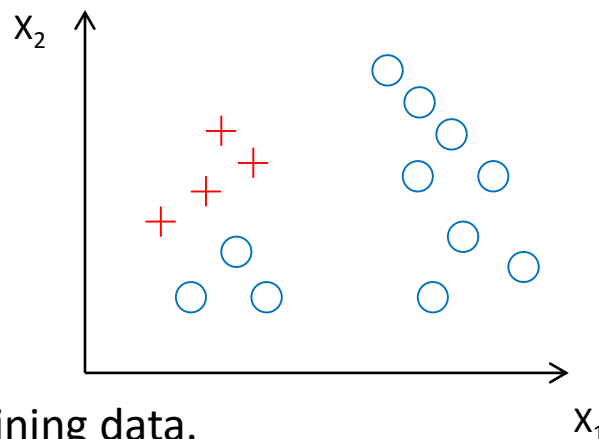
- self-training models,
- EM with generative mixture,
- co-training,
- transductive support vector machines,
- graph-based methods.

# Learning from Positive and Unlabeled data

One has a set of examples of a class  $+$ , and a set of unlabeled examples with instances of a class  $+$  and also not from  $+$  (negative examples).

**Goal:** Build a classifier to classify the unlabeled examples and/or future (test) data.

**Key feature of the problem:** no labeled negative training data.



This problem is known as PU-learning.

*An example is when a company has a database with details on its customer – positive examples, and a database with details on individuals who are not customers, but could become or not customers if they were proposed some products.*

2-step strategy for text classification

Step 1: Identifying a set of reliable negative documents from the unlabeled set.

Step 2: Building a sequence of classifiers by iteratively applying a classification algorithm and then selecting a good classifier.