# Machine Learning
# Hands-on Workshop

# Data Venture

# Data Venture

http://www.opendatalab.xyz

| | | |
|---|---|---|
| **Président** | Julien | Jerphanion |
| **Vice-Président** | Sylvain | Marchienne |
| **Trésorière** | Léna | Schofield |
| **Secrétaire** | Vincent | Dehaye |
| **Resp. Com** | Sevan | Garois |
| **Resp. Setup** | Valentin | Montupet |
| **Resp. Scientifique** | Sy Toan | Ngo |
| **Resp. Talks** | Gabriel | Hurtado |
| **Resp. Projet** | Benjamin | Rivière |
| **Resp. Partenariat** | Antoine | Jeannot |

# About me

Jonathan DEKHTIAR

2nd Year PhD. Student

Engineer @ UTC 2015
Machine Learning and Statistical Analysis

**Keeping in touch'**

Twitter: @Born2data

LinkedIn: https://www.linkedin.com/in/jonathandekhtiar/

Tech. Blog : https://www.born2data.com

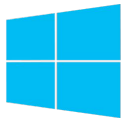RSS Feed: https://www.feedcrunch.io/@dataradar

Email: contact@jonathandekhtiar.eu and jonathan.dekhtiar@utc.fr

# How to ML

Most efficient, lightweight, hacking-ready, versatile virtualization platform.

## Jupyter Notebooks for Data Science

docker pull jupyter/datascience-notebooks *#take a coffee*

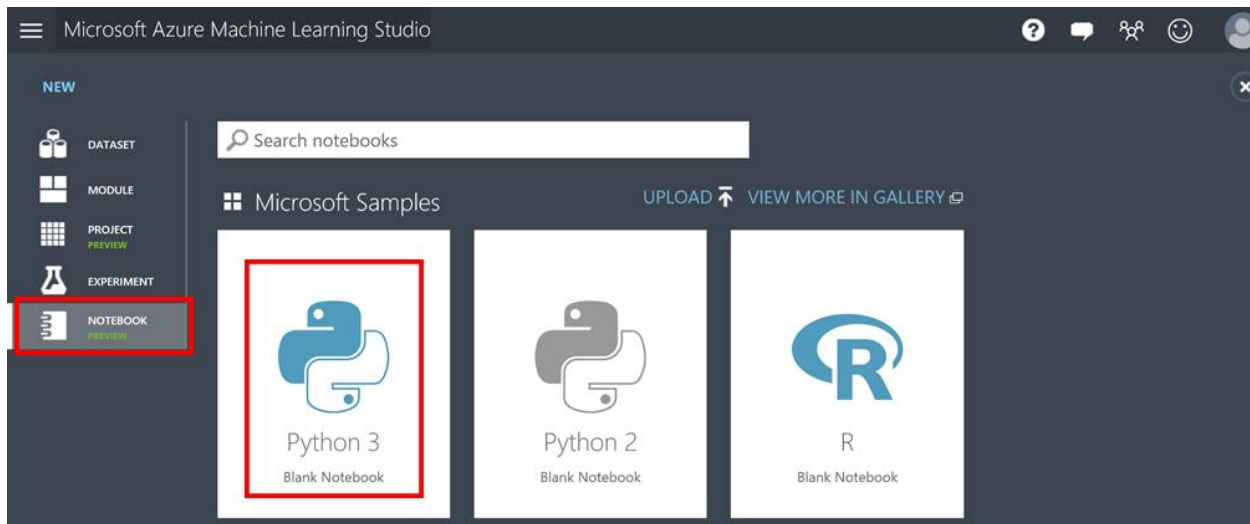docker run -d -p 8888:8888 jupyter/datascience-notebook start-notebook.sh --NotebookApp.token=''

Open Browser on: http://localhost:8888/

# How to ML

?



**Windows Azure**

https://studio.azureml.net/

Free-Tier Usage.
Quite fast and scalable.
Cloud based on - Azure SaaS

# Jupyter Notebooks ?

# Machine Learning



## Herbert Simon

**Turing Award** 1975

**Nobel Prize in Economics** 1978

"Learning is any process by which a system improves performance from experience."

"Machine Learning is concerned with computer programs that automatically improve their performance through experience. "

# Why Machine Learning ?

- Develop systems that can automatically adapt and customize themselves to individual users.
    - Personalized environment, context-awareness, etc.

- Discover new knowledge from large databases (data mining).
    - Pattern Discovery, Influential Parameters, etc.

- Ability to mimic human and replace certain monotonous tasks - which require some intelligence.
    - like analyzing imagery data 2D/3D and comparing them

- Develop systems that are too difficult/expensive to construct manually because they require specific detailed skills or knowledge tuned to a specific task (knowledge engineering bottleneck).

- Etc.

# Why Now ?

- Flood of available data (especially with the advent of the Internet)

- Increasing computational power (GPU Computing, Moore's Law, etc.)

- Growing progress in available algorithms and theory developed by researchers

- Increasing support from industries => Research fundings

- Increasing complexity of challenges induced by massive data and complex data

- etc.

# Who needs/uses ML ?

# What's the concept ?

## Learning

$$\Downarrow$$

**Improving** <u>over time</u> with <u>experience</u> at **some task**.

1. Improve over task: T
2. Measure Performance with indicators: P
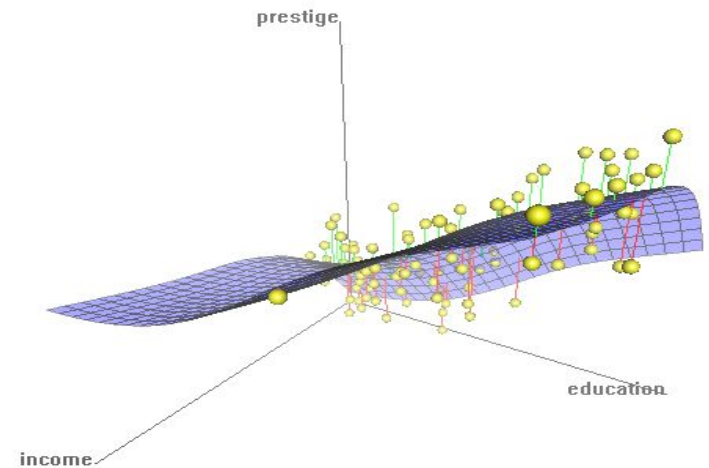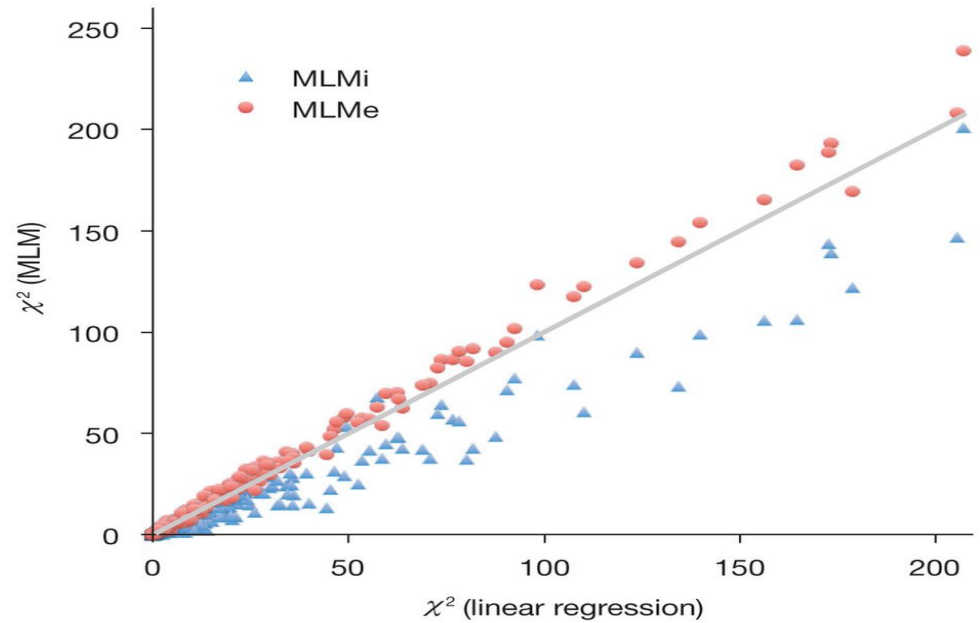3. Based on Experience: E

# Regression

**Objective:**

To determine a parameter in function of a set of features.

$$f(x_1, ..., x_n) = y$$

**A few models:**

- Linear Regression
- Logistic Regression
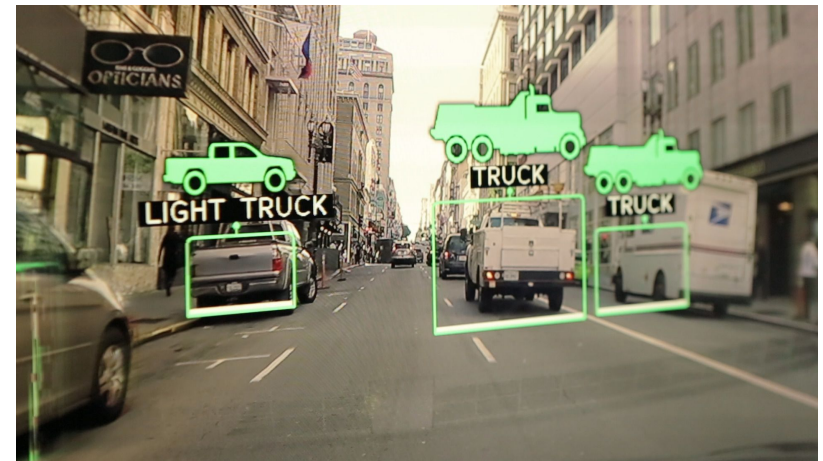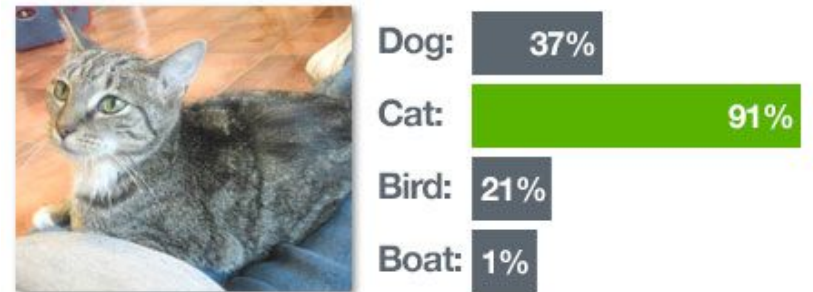- Quadratic Regression
- Bayesian Regression
- etc.

# Classification

**Objective:**

To apply a label on a given set of features.



**A few models:**

- K-Nearest Neighbors (KNN)
- Support Vector Machine (SVM)
- Random Forest
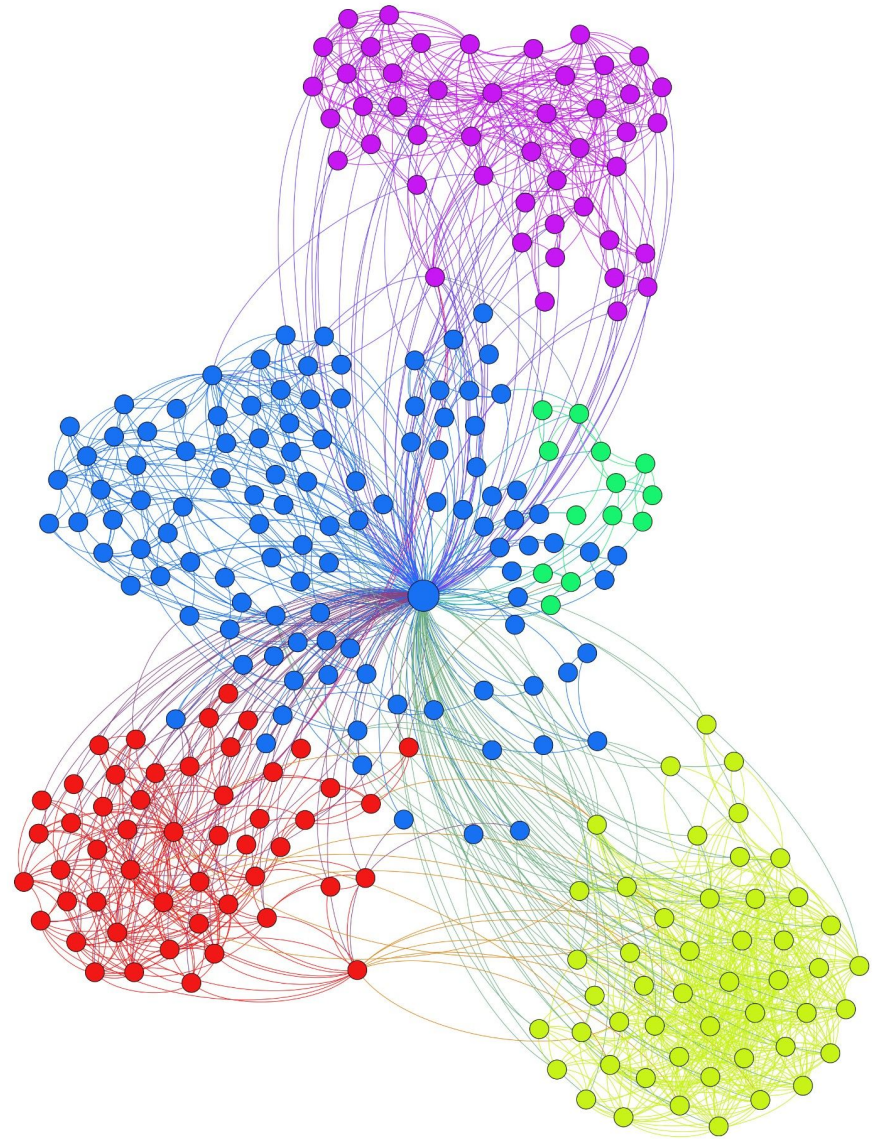- (Deep?) Neural Network (NN)
- etc.

# Clustering

**Objective:**

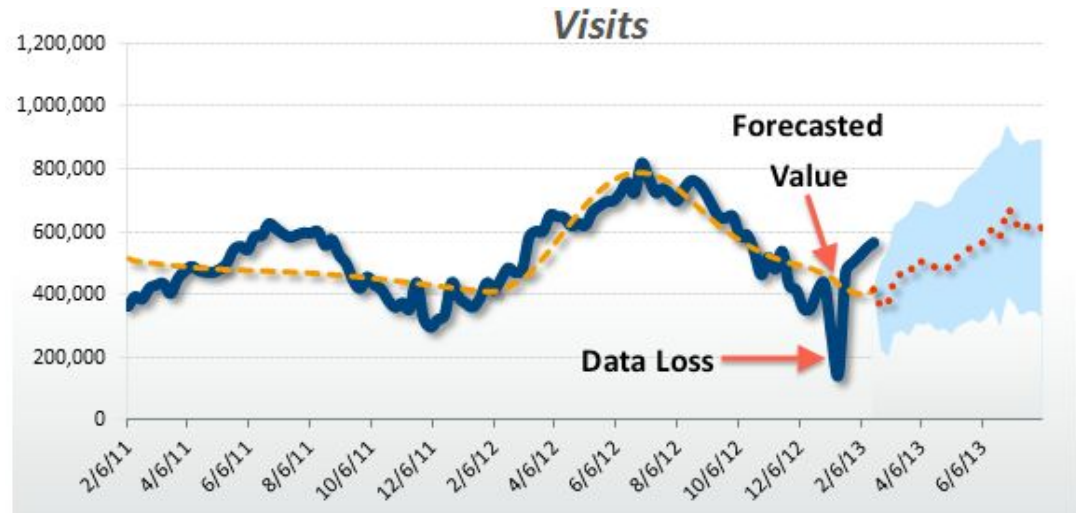Discover hidden groups or patterns inside the data

**A few models:**

- K-Means
- EM-Clustering
- Hierarchical Clustering
- Neural Networks (SOM & SOFM)
- etc.

# Time Series Prediction

**Objective:**

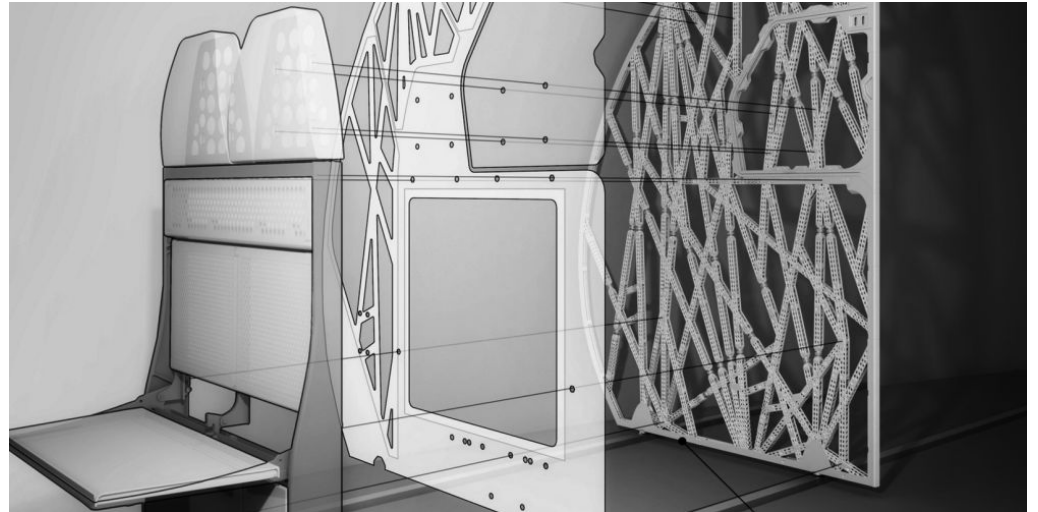Predict the next value(s) of
a time series beforehand.



**A few models:**

- Recurrent Neural Networks (RNN)
- Long Short Term Memory Neural Networks (LSTM)
- etc.

WARNING:
HARD WORK
AHEAD
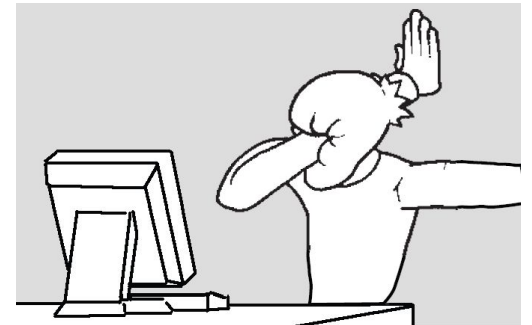
# Generative Models

**Objective:**

Generate a completely new piece of data that looks like a human-made one.
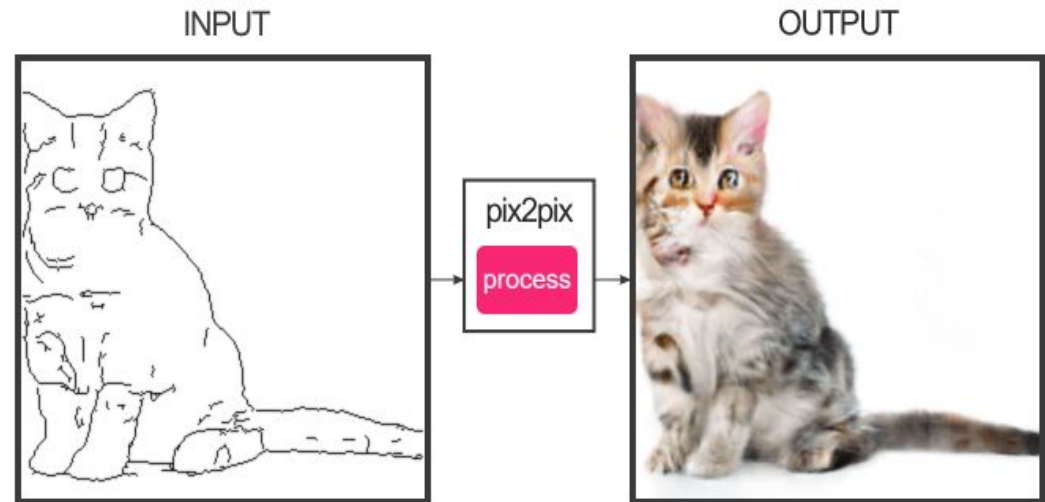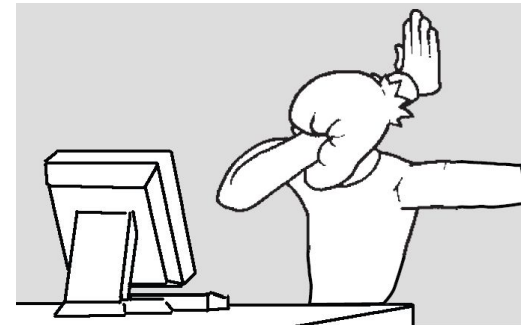
**A few models:**

- Generative Adversarial Neural Network (GAN)
- RNN
- LSTM
- etc.

# Generative Models

**Objective:**

Generate a completely new piece of data that looks like a human-made one.



INPUT

pix2pix

process

OUTPUT

**A few models:**

- Generative Adversarial Neural Network (GAN)
- RNN
- LSTM
- etc.



WARNING: HARD WORK AHEAD

# What about generalisation?

# What about generalisation?

They are "**Cats**"...
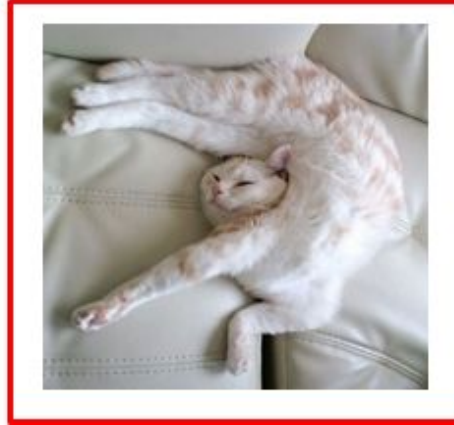
Too easy Bro !
Make it way harder !

# What about generalisation?

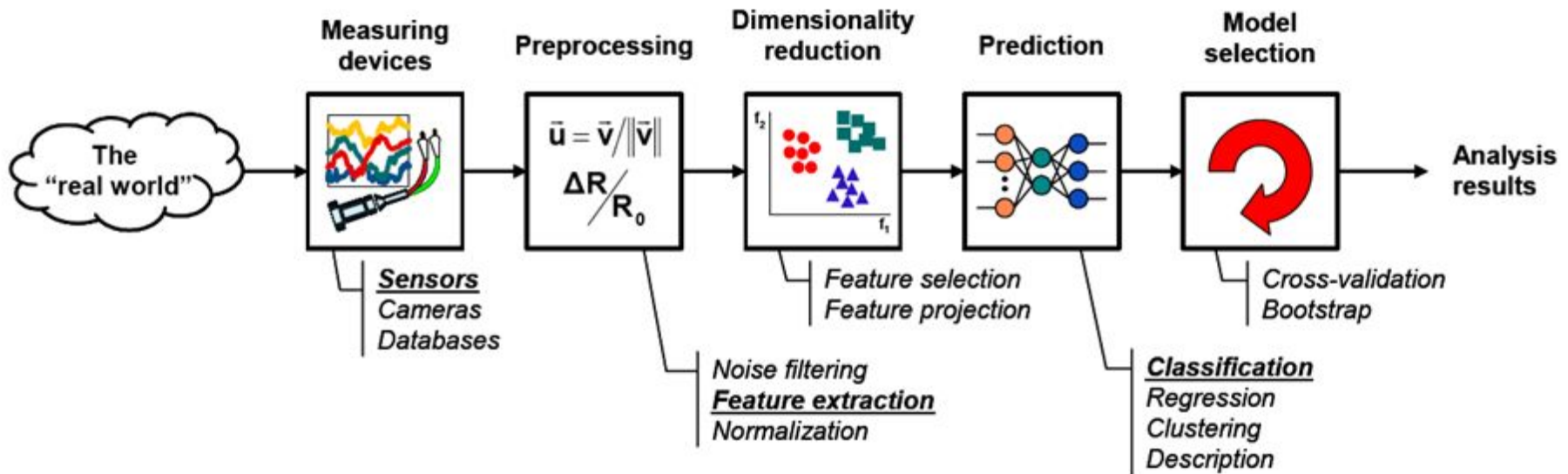# Learning workflow

# Speaking about... data !

Input Features      Label

| Number of new Recipients | Email Length (K) | Country (IP) | Customer Type | Email Type |
|---|---|---|---|---|
| 0 | 2 | Germany | Gold | Ham |
| 1 | 4 | Germany | Silver | Ham |
| 5 | 2 | Nigeria | Bronze | Spam |
| 2 | 4 | Russia | Bronze | Spam |
| 3 | 4 | Germany | Bronze | Ham |
| 0 | 1 | USA | Silver | Ham |
| 4 | 2 | USA | Silver | Spam |

Instances

Numeric      Nominal      Ordinal

# Kaggle

# AirBnB

## Where will a new guest book their first travel experience?

New users on Airbnb can book a place to stay in 34,000+ cities across 190+ countries. By accurately predicting where a new user will book their first travel experience, Airbnb can share more personalized content with their community, decrease the average time to first booking, and better forecast demand.

In this recruiting competition, **Airbnb challenges you to predict in which country a new user will make his or her first booking**.

# Hands-on Lab

https://goo.gl/O2PWM7