atlantbh

DATA SCIENCE PROJECT

# *DEMO PRESENTATION*

**Mentor:** Belma Ibrahimović
**Intern:** Mehmed Kadrić

# MOTIVATION

DATA SCIENCE

MACHINE LEARNING

# TABLE OF CONTENTS

- TASK 1: Python for Data Science
- TASK 2: Classification in Data Science
- TASK 3: Anomaly Detection
- Conclusion

# PYTHON FOR DATA SCIENCE

- Important for DS
- Easy to use language
- Most common libraries used in data science:
  - Numpy
  - Pandas
  - Scikit-learn
  - Matplotlib

„Started from the bottom, now I'm here"

# CLASSIFICATION IN DATA SCIENCE

- Use-case scenario:
  - Titanic sank after colliding with an iceberg, killing 1502 out of 2224 passengers and crew (32% survival rate)
  - Latest discovery have shown that there were additional 10 passengers
- Goal:
  - Create the system which would predict probability of their survival
- Dataset was given

# TITANIC

- Supervised learning
- Preprocessing:
  - 1303 examples, 4 features

| | Name | PClass | Age | Sex | Survived |
|---|---|---|---|---|---|
| 1 | Allen, Miss Elisabeth Walton | 1st | 29 | female | 1 |

- Missing values in Age column and PClass column
- Only one outlier in PClass column
- Dimensionality reduction

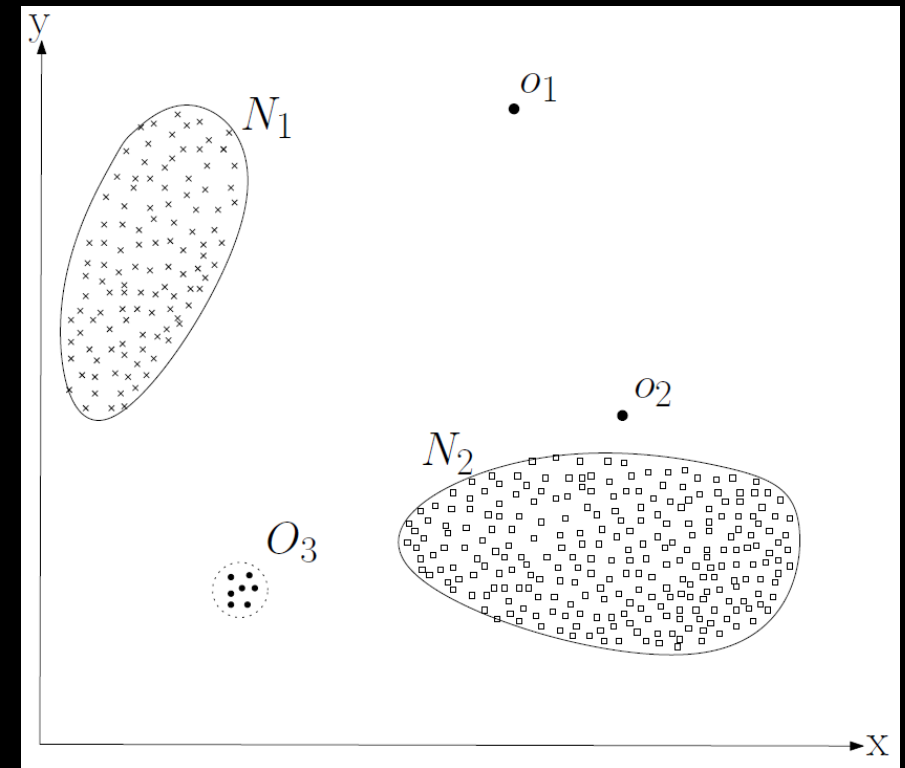| | Title | PClass | LifeStage | Survived |
|---|---|---|---|---|
| 1 | Miss | 1st | 2 | 1 |

# TITANIC

- There are many classification algorithms such as: Logistic Regression, Naive Bayes classifier, SVM, KNN, Decision Trees, etc.

- Support Vector Machine Classifier gave the best results

| Actual output | Model output | Probability 0 | Probability 1 |
|---------------|--------------|---------------|---------------|
| 0 | 0 | 0.98511669 | 0.01488331 |
| 0 | 0 | 0.98511669 | 0.01488331 |
| 1 | 0 | 0.78656724 | 0.21343276 |
| 0 | 0 | 0.96557864 | 0.03442136 |
| 0 | 0 | 0.99201528 | 0.00798472 |
| 0 | 0 | 0.99201528 | 0.00798472 |
| 1 | 1 | 0.19527331 | 0.80472669 |
| 1 | 0 | 0.78665345 | 0.21334655 |
| 1 | 1 | 0.15862112 | 0.84137888 |
| 1 | 1 | 0.21237841 | 0.78762159 |

*Accuracy*
*Precision*
*Recall*
*F1-score*

# Anomaly detection

- Anomalies are patterns in data that do not conform to a well defined notion of normal behavior

- Anomaly detection techniques:
  - Classification based
  - Nearest Neighbor based
  - Clustering based
  - Statistical
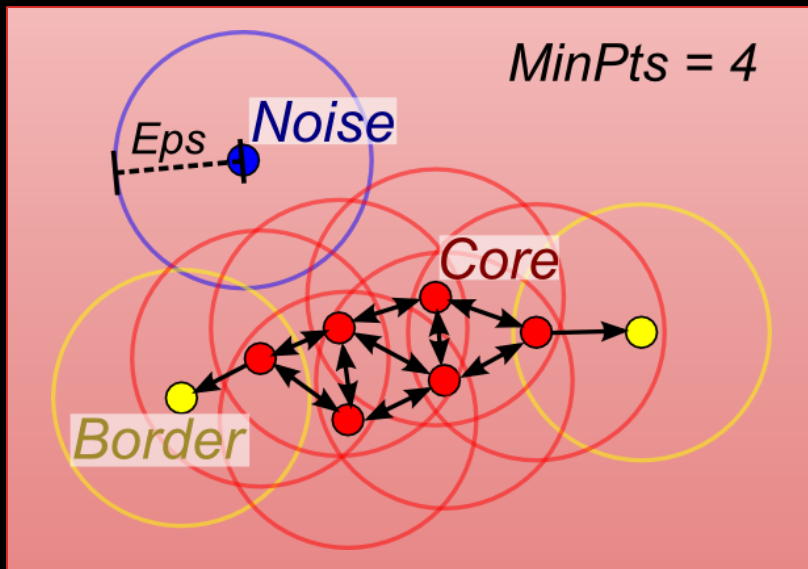  - Information theoretic
  - Spectral

# Anomaly detection

- Develop system which can be used for smart production monitoring and fault detection in manufacturing processes

- Given dataset: 1567 instances, 591 features

- Unbalanced dataset: 104 fails

- Only numeric values

- Missing values

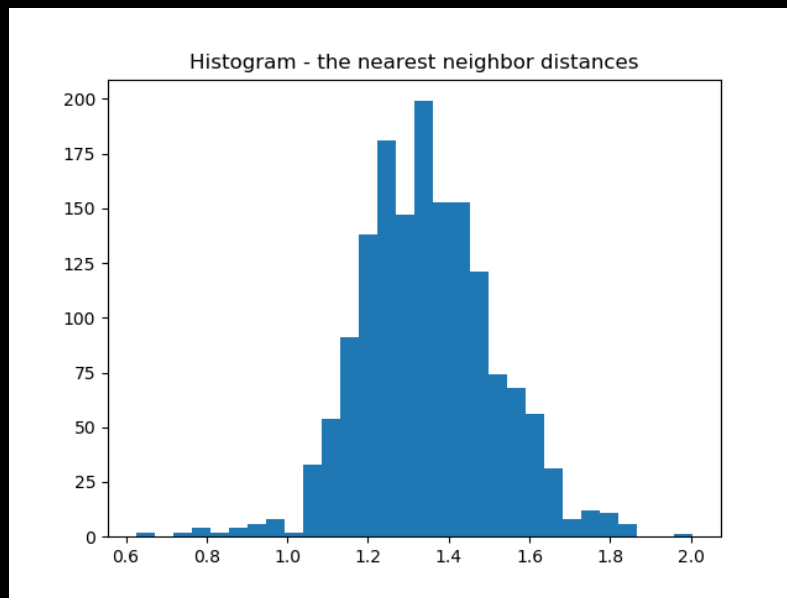- Outliers

# Anomaly detection

- DBSCAN (Density-based spatial clustering of applications with noise) algorithm
  - Core point
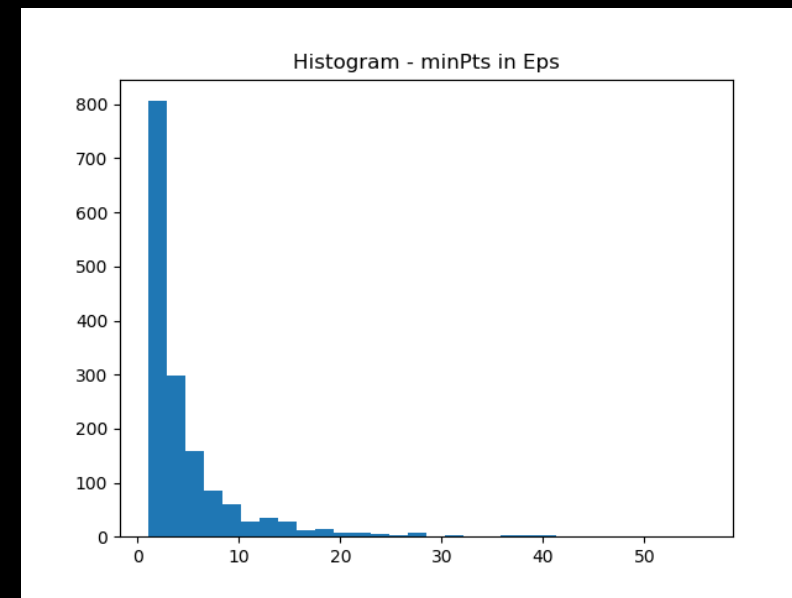  - Border point
  - Noise point



- Advantages:
  - Does not require number of clusters
  - Can find arbitrarily shaped clusters
  - Requires two parameters

- Disadvantages:
  - border points that are reachable from more than one cluster can be part of either cluster, depending on the order the data are processed
  - Can't cluster data sets well with large differences in densities

# Anomaly detection

- Parameter tuning:
  - In DBSCAN learning algorithm, thare are two parameters: eps and minPts



Eps = 1.45



minPts = 6

# Anomaly detection

- Unbalanced dataset – accuracy is not a good measure

|  | Precision | Recall | F1-score | Support |
|---|---|---|---|---|
| -1 | 0.07 | 0.51 | 0.13 | 104 |
| 0 | 0.94 | 0.54 | 0.69 | 1463 |
| Avg/total | 0.88 | 0.54 | 0.65 | 1567 |

# Anomaly detection

- Possible applications:
  - Fraud detection
  - Fault detection
  - System health monitoring

# THANK YOU!