# *DATA SCIENCE INTERNSHIP – MONTH 1*

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

# MOTIVATION

- Effects of machine learning algorithms are all around us
- Curiosity
- Digit recognition

# 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
- Supervised learning

**Preprocessing**
- 1303 examples, 4 features
- Missing values in Age column and PClass column
- Only one outlier in PClass column
- Dimensionality reduction

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

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

- 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

| Name | Probability of survival | Actually survived |
|---|---|---|
| Sage, Miss Constance | 0.0148 | 0 |
| Sage, Miss Dorothy | 0.0148 | 0 |
| Wilkinson, Mrs Elizabeth Anne | **0.2134** | **1** |
| Thomas, Master Assad Alexander | 0.0344 | 0 |
| Zakarian, Mr Artun | 0.0079 | 0 |
| Zakarian, Mr Maprieder | 0.0079 | 0 |
| Anderson, Mr Harry | 0.8047 | 1 |
| Andrews, Miss Kornelia Theodosia | **0.2133** | **1** |
| Brown, Mrs John Murray (Caroline Lane Lamson) | 0.8413 | 1 |
| Carter, Mr William Ernest | 0.7876 | 1 |

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

# 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

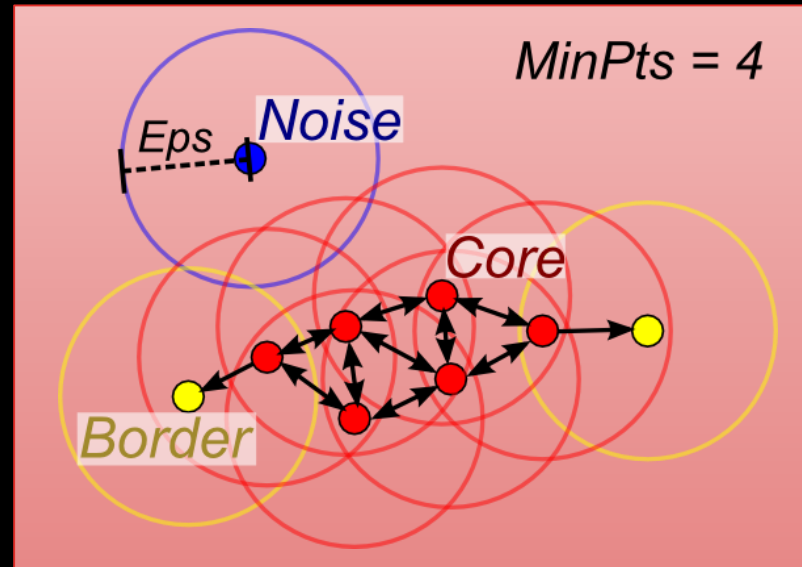- Only numeric values

- Missing values

- Outliers

# Anomaly detection

- 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
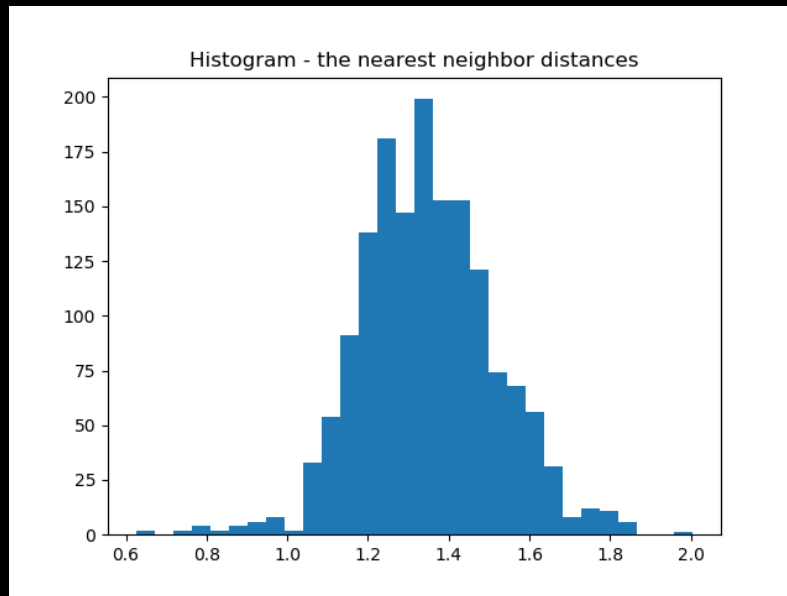  - Information theoretic
  - Spectral

# Anomaly detection

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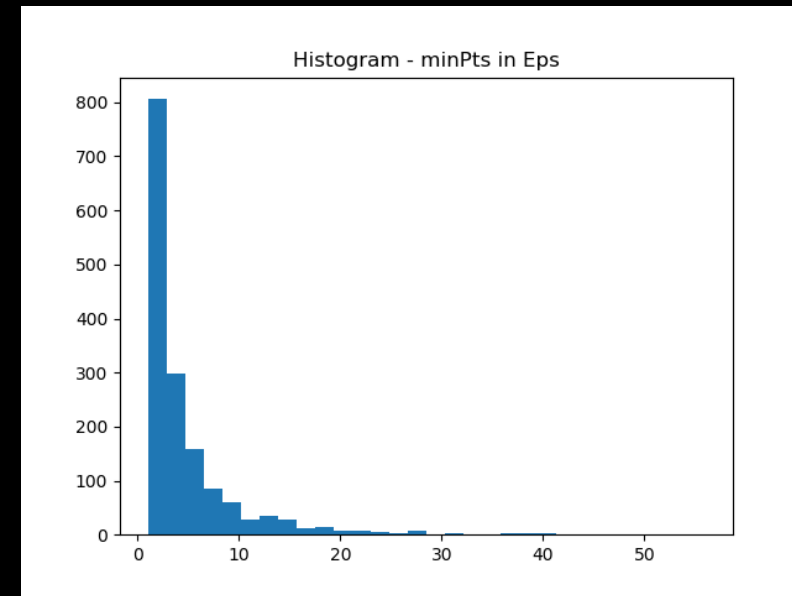  - Core point
  - Border point
  - Noise point

# 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
- It's important to detect the fault in a timely manner

|  | **Precision** | **Recall** | **F1-score** | **Support** |
|---|---|---|---|---|
| Outlier | 0.07 | 0.51 | 0.13 | 104 |
| Not outlier | 0.94 | 0.54 | 0.69 | 1463 |
| Avg/total | 0.88 | 0.54 | 0.65 | 1567 |

# CONCLUSION

- Learn while working
- Math is fun
- Advantage of creative mindset
- There is no general approach for solving machine learning problems

# THANK YOU!