**Deadline: May 18, Noon! Discussions: May 20. Each group must attend based on schedule!**

## Final Project description

In this project, you will apply advanced unsupervised learning techniques and anomaly detection methods on assigned real-world datasets. You must analyze, cluster, visualize, validate, and extract insights from complex, high-dimensional data.

You are expected to deeply explore the data, perform extensive sensitivity analysis, and properly document your methodology and findings.

## Assigned Datasets

You **must** use the following datasets (both required):

1. **Credit Card Fraud Detection Dataset** (Kaggle)

   o Transactions labeled as fraud or not fraud.

   o Highly imbalanced, anonymized data with 30 features.

2. **Mall Customer Segmentation Dataset**

   o Data on customer age, income, spending score, etc.

   o A simpler dataset to practice clustering more visibly.

## Techniques and Algorithms You Must Apply

**1. Dimensionality Reduction**

- Apply **PCA** on both datasets:

  o Analyze explained variance.

  o Choose the minimum number of components preserving 90–95% variance.

  o Visualize reduced space (2D scatter plots).

- Apply **t-SNE** for visualization:

  o Visualize clusters formed after PCA.

  o Create 2D t-SNE plots **before** and **after** clustering

## 2. Clustering Analysis

- **K-Means Clustering**:

    o Use the **Elbow method** and **Silhouette score** to determine best k.

    o Plot **Inertia vs. *k*** and **Silhouette score vs. *k***.

    o Create **Silhouette diagrams** for at least 3 values of k.

- **K-Means++ Initialization**:

    o Compare standard K-Means and K-Means++.

    o Show how initialization impacts convergence and results.

- **Mini-Batch K-Means**:

    o Apply on both datasets.

    o Compare speed, memory usage, and accuracy with standard K-Means.

- **DBSCAN Clustering**:

    o Tune *eps* and *min_samples* carefully.

    o Analyze how noise points are classified.

    o Plot clusters and compare to K-Means clusters.

## 3. Anomaly Detection

- Apply to **Credit Card Fraud Dataset** ONLY:

    o **Isolation Forest** for anomaly detection.

    o **One-Class SVM** for anomaly detection.

- Evaluate using:

    o **Precision**, **Recall**, **F1-Score**.

    o **Confusion Matrix** (use actual fraud labels).

- Analyze:

    o How unsupervised models detect fraud vs. true labels.

- o   Which method performs better, and why.

## Sensitivity Analysis Requirements

For **each clustering technique** (KMeans, MiniBatch KMeans, DBSCAN):

- Vary important hyperparameters:

  - o   For KMeans: k = 2 to 20

  - o   For DBSCAN: grid search on eps and min_samples

- Plot and interpret:

  - o   **Inertia graphs**

  - o   **Elbow curves**

  - o   **Silhouette score plots**

  - o   **Silhouette diagrams**

**Additionally**:

- Compare clustering performance **with and without** dimensionality reduction (PCA).

- Compare clustering in full space vs. reduced 2D space.

## Special Tasks (Mandatory)

- Create a **comparison table** summarizing:

  - o   Inertia values

  - o   Silhouette scores

  - o   Execution times (for KMeans vs MiniBatch)

  - o   Number of clusters found (for DBSCAN)

- Document:

  - o   How scaling (StandardScaler) affects clustering.

  - o   How random seed affects results (especially for t-SNE and KMeans).

- Feature engineering:

  - o   If applicable, create at least **two new features** based on existing ones.

o   Justify their usefulness.

## Deliverables

1. **Technical Report** (embedded in Jupyter itself):

   o   Introduction (motivation, datasets description)

   o   Data preprocessing and feature engineering

   o   Dimensionality reduction (PCA, t-SNE)

   o   Clustering experiments (KMeans, MiniBatch KMeans, DBSCAN)

   o   Anomaly detection (Isolation Forest, One-Class SVM)

   o   Sensitivity analysis results

   o   Tables, charts, graphs, interpretation

   o   Final insights and discussion

   o   Challenges and how you solved them

   o   Findings!

   o   Work load distribution between team members! If we find something different from what you say, all team members might get the lowest grade in the team!

2. Fully executable Jupyter Notebook:

   o   Well-commented

   o   Organized sections

   o   Code that runs from start to finish

## General Important Notes:

- We KNOW that you can get some suggested codes from Kaggle AND other sources, HOWEVER, be very careful that *the requirements of the project are VERY detailed*! Every step counts!

- You will have to do sufficient preprocessing, as needed, tables, figures, plots, and then discussions that support your argument!

- We will check in between the lines, so be ready to be tested in whatever you write and around these concepts to check your level of understanding!

- **Failure to attend the presentation will result in a ZERO no matter what your solution is!**
- **There will be NO MAKEUPS!**