

Università degli studi di Milano - Bicocca

Department of Informatics, System and Communication (DISCo)

Master's Degree in Data Science

Integrated Supervised and Unsupervised Learning for Multivariate Socio-Economic Data Analysis

Data Science Lab Project



Davide Fabio Loreti - 865309

June 26, 2025

Contents

Code Repository	1
Introduction	2
1 Exploratory Data Analysis (EDA)	3
1.1 Overview of the Dataset	3
1.2 Gender Distribution	4
1.3 Age Distribution	5
1.4 Geographic Area Distribution	6
1.5 Education Level Distribution	7
1.6 Income Bracket Distribution	8
1.7 Correlation Analysis of Numerical Variables	8
1.8 EDA Summary	10
2 Machine Learning Algorithms	11
2.1 Data Preparation	11
2.2 Construction and evaluation of supervised models	13
2.3 Best Model Evaluation	15
2.4 Clustering Analysis	16
2.5 Clustering Evaluation	18
2.6 Second Clustering Attempt with Two Clusters	19
2.7 Clustering with Algorithms for Non-Spherical or Elliptical Structures	21
2.8 Visualization and Interpretation of Results	21
2.9 Final Remarks	24
3 Final considerations	25

Code Repository

The full code used to preprocess the data, train and evaluate the models, and generate all figures is available at the following repository:

https://github.com/DavideFabioLoreti/DATA_SCIENCE_LAB/blob/main/DS%20LAB%20PROJECT.ipynb

Introduction

This report presents an analysis of data collected by the Bank of Italy concerning the financial literacy survey of the adult Italian population. The main objective is to explore the demographic and socio-economic characteristics of the participants and to apply machine learning techniques for classification and clustering tasks on relevant variables.

Initially, an exploratory data analysis (EDA) was conducted to examine the distribution of key variables such as gender, age, geographic area, education level, and income bracket. This preliminary analysis helped to understand the sample's structure and to identify patterns and potential challenges for further modeling.

For the predictive modeling phase, classification algorithms were applied to predict the gender of individuals based on other survey variables. Additionally, clustering methods were used to identify natural groupings within the data, uncovering underlying profiles among respondents.

The ability to classify gender from financial and socio-economic data holds practical relevance for financial institutions and fintech companies. Often, due to privacy constraints or incomplete records, demographic attributes such as gender may be missing or unavailable. Inferring gender indirectly through behavioral financial data can enable organizations to better tailor marketing strategies, personalize product offerings, and improve risk assessment models. This approach supports data-driven decision-making while respecting privacy, allowing businesses to address customer needs more effectively and identify potential disparities or patterns linked to gender in financial behavior.

The classification models were evaluated using metrics such as accuracy, precision, recall, and F1-score, while clustering results were interpreted to provide insights into the heterogeneity of the sample.

This combination of exploratory data analysis, classification, and clustering aims to offer a comprehensive understanding of the dataset, supporting the development of data-driven strategies and policy interventions related to financial literacy.

1 Exploratory Data Analysis (EDA)

The Exploratory Data Analysis (EDA) phase represents a crucial preliminary step in any data science project, serving to uncover the underlying structure, detect patterns, and identify potential anomalies or data quality issues within the dataset. Before applying any advanced modeling techniques, such as machine learning or deep learning algorithms, it is essential to thoroughly understand the data at hand to ensure the validity and reliability of subsequent analyses.

In this project, our focus is on a comprehensive dataset related to financial literacy among the adult Italian population. The dataset is enriched with various demographic and socio-economic variables that provide a multifaceted view of the individuals surveyed. By carefully examining these features, we aim to gain insights into the distribution and relationships between variables, which will inform the selection and tuning of predictive models later on.

1.1 Overview of the Dataset

The dataset comprises a total of 4,862 rows and 219 columns, corresponding respectively to the number of individual responses and the total number of features collected. Its shape is therefore defined as (4862, 219).

The dataset includes several key attributes describing individual respondents, including:

- **Gender**, capturing the biological sex of the participant;
- **Age**, representing the respondent's age in years;
- **Geographic Area**, indicating the regional location within Italy;
- **Education Level**, detailing the highest attained educational qualification;
- **Income Bracket**, classifying the monthly income range of the individual.

Each of these variables is subject to both visual and statistical exploration. Visualizations such as histograms, bar charts, and kernel density estimates (KDE) allow for an intuitive grasp of the data distribution, while summary statistics provide quantitative measures such as means, medians, and counts. This dual approach facilitates the detection of any irregularities, including missing values, outliers, or unexpected patterns, which could influence model performance if left unaddressed. Through this careful examination, we set the foundation for a robust and meaningful analysis.

For a complete description of all the dataset features, refer to the official documentation available at the following link: [Descrizione delle variabili - 2023 \(PDF\)](#).

1.2 Gender Distribution

Descriptive statistics:

- Unique categories: 2 (Female, Male)
- Most frequent category: **Female** (with 2446 observations against 2416)
- Missing values: 0



Figure 1: Gender distribution in the dataset.

As shown in the figure, the population is fairly balanced in terms of gender, with a slight predominance of females. This balance makes the variable suitable for future comparative analysis.

1.3 Age Distribution

Descriptive statistics:

- Mean age: 50.31 years
- Median: 50.00
- Minimum and maximum: 18 - 79
- Missing values: 219

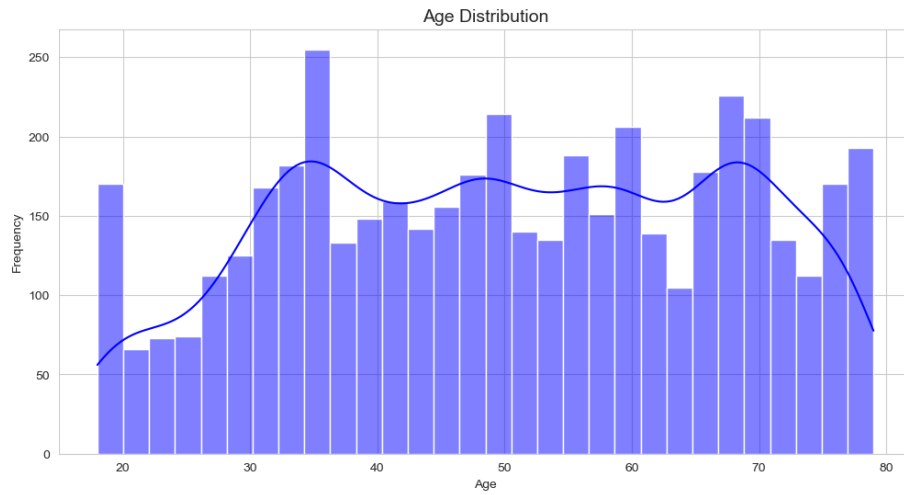


Figure 2: Age distribution with histogram and KDE curve.

The age distribution appears symmetric and centered around 50 years. The presence of 219 missing values requires proper handling (e.g., imputation or removal) before proceeding to modeling.

1.4 Geographic Area Distribution

Descriptive statistics:

- Unique categories: 5
- Most represented area: Nord - Ovest
- Missing values: 0

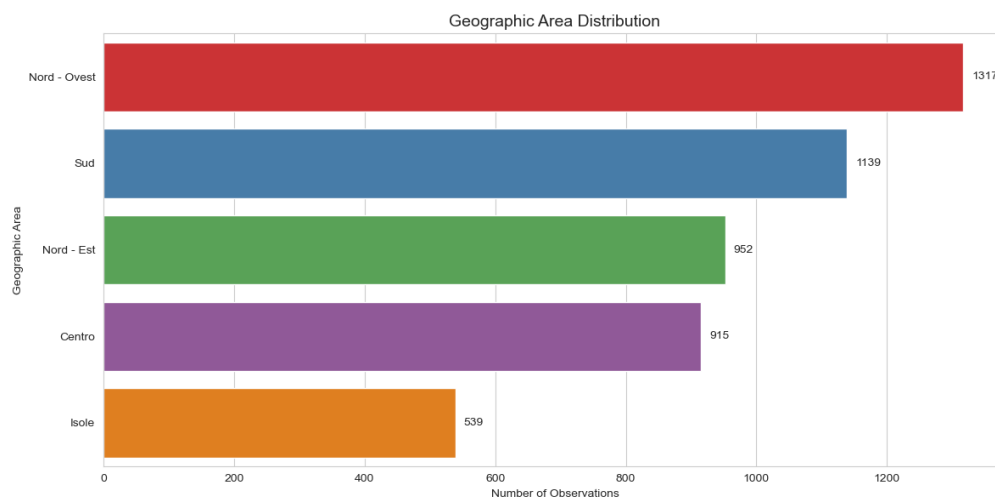


Figure 3: Distribution by geographic area.

The North-West area is the most represented in the sample. This concentration may reflect higher participation or population density in that region.

1.5 Education Level Distribution

Descriptive statistics:

- Unique categories: 10
- Most frequent category: Scuola media superiore con diploma
- Missing values: 0

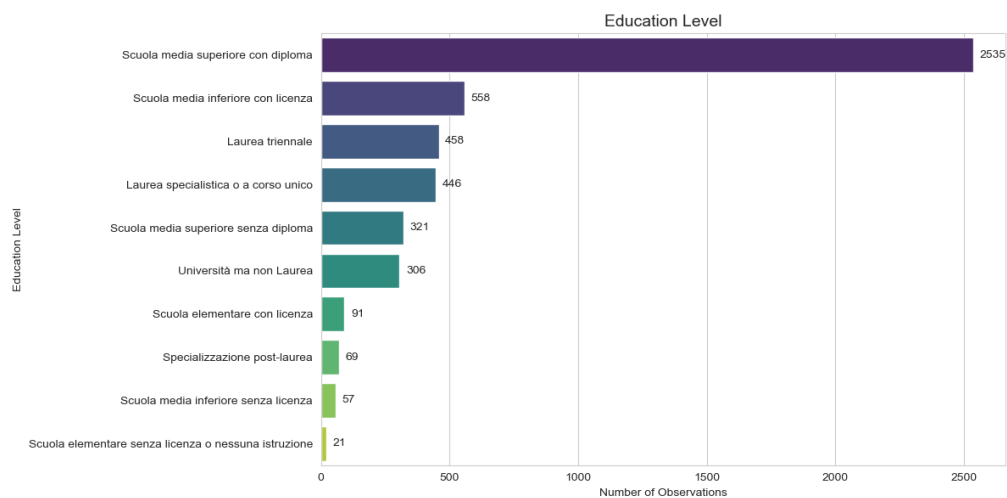


Figure 4: Education level distribution.

The most common education level is a high school diploma, followed by various other degrees. The large variety of categories makes this variable useful for segmentation and classification models.

1.6 Income Bracket Distribution

Descriptive statistics:

- Unique categories: 5
- Most frequent category: tra 1.751 Euro e 2.900 Euro al mese
- Missing values: 0

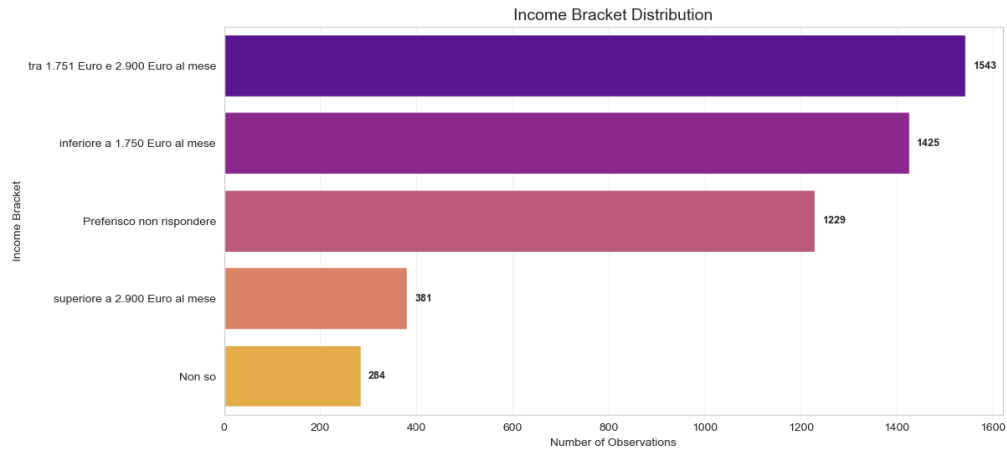


Figure 5: Monthly income bracket distribution.

The central income bracket is the most represented, suggesting a concentration in the middle class. All categories are well distributed and no missing values are present.

1.7 Correlation Analysis of Numerical Variables

This subsection presents the results of the Pearson correlation analysis conducted on a selected set of numerical variables from the dataset. The considered variables include demographic information (**Age**), Class Age **qd7_a**, six binary variables (**qd5_1** through **qd5_6**), and a sampling weight (**wght**). The objective was to identify any significant linear relationships among these variables to provide interpretative insights and guide further analysis.

The variables were selected based on their theoretical and practical relevance in the context of financial literacy and socio-demographic characteristics of the sample. In particular:

- **Age** is a key demographic factor often associated with various financial behaviors and competencies; therefore, its impact on other variables was examined.

- **class_age** represents the respondent's age bracket, categorized into seven ordinal groups ranging from 18–19 to 70–79 years old, with an additional category for non-responses. This variable is useful for analyzing demographic trends and associations with financial behavior or literacy.
- The binary variables **qd5_1** through **qd5_6** correspond to specific survey questions, likely related to particular financial behaviors or knowledge. Analyzing their interrelations helps to identify response patterns or clusters potentially influencing literacy levels.
- The **sampling weight (wght)** was included to check for possible associations with other variables and to control for the influence of sample weighting.

Figure 6 displays the heatmap of the correlation matrix described above. The color gradient emphasizes positive correlations in warm colors and negative correlations in cool colors, with numeric values annotated for clarity.

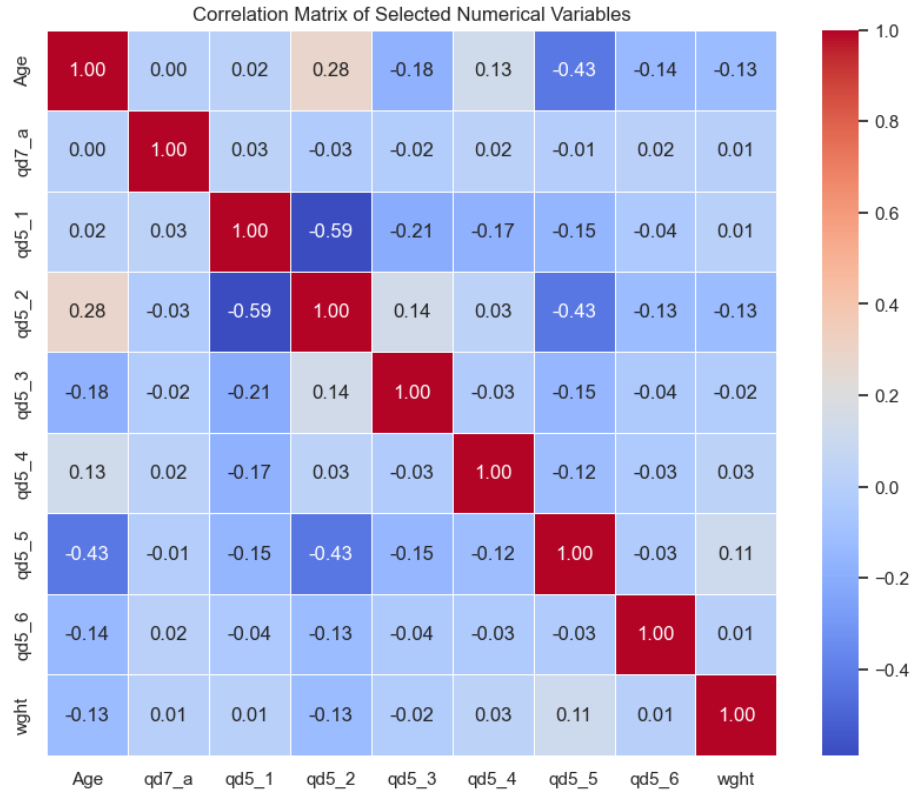


Figure 6: Heatmap of the Pearson Correlation Matrix for Selected Numerical Variables

Age shows a moderate positive correlation with the binary variable `qd5_2` ($r = 0.28$), suggesting that older respondents tend to have higher values on this indicator. Conversely, Age correlates moderately negatively with `qd5_5` ($r = -0.43$), indicating that this feature tends to decrease as age increases. Other correlations involving Age are generally weak, indicating limited linear association.

Class Age (`qd7_a`) is essentially uncorrelated with the other variables. Correlation coefficients involving `qd7_a` are close to zero.

Binary Variables (`qd5_1` to `qd5_6`) exhibit a strong negative correlation between `qd5_1` and `qd5_2` ($r = -0.59$), indicating an inverse relationship between these two indicators. Other correlations among these variables are weak to moderate and mostly negative. For example, `qd5_5` shows negative correlations with Age and `qd5_2` ($r = -0.43$), suggesting distinct behavioral or demographic patterns captured by these variables.

Sampling Weight (`wght`) shows negligible correlations with all other variables, as expected given its role as a weighting factor, ideally independent of measured characteristics.

The minimal correlations between the financial literacy score and other variables suggest that the latter do not linearly explain variation in financial literacy within the sample. This may indicate the influence of unmeasured factors or nonlinear relationships. The strong negative correlation between `qd5_1` and `qd5_2` could reflect mutually exclusive responses or contrasting categories within the survey. Correlations with Age highlight demographic trends, with indicators such as `qd5_5` more frequent among younger respondents and `qd5_2` more prevalent among older participants.

1.8 EDA Summary

The exploratory data analysis reveals a dataset in generally good condition, except for missing values in the **Age** variable. The distributions are consistent with expectations for the Italian population and show a balanced representation across key demographic dimensions. These findings support the use of the selected variables in the subsequent modeling phases using Machine Learning algorithms.

2 Machine Learning Algorithms

The previous chapter presented a thorough exploratory data analysis (EDA), which was essential to understanding the composition and key characteristics of the dataset. Through the examination of demographic and socio-economic variables, such as gender, age, geographic area, education level, and income bracket, we identified meaningful patterns and detected some data quality issues, including missing values in the age variable. These insights laid a solid foundation for building reliable and effective predictive models.

This chapter focuses on applying machine learning techniques to deepen the analysis and leverage the information contained in the dataset. Specifically, we adopted a dual approach: unsupervised clustering methods to uncover natural groupings within the data without relying on labels, and supervised classification models aimed at predicting the target variable, gender.

For classification, several well-established algorithms were employed, including Random Forest, Support Vector Machines, K-Nearest Neighbors, and XG-Boost, chosen for their robustness and ability to handle structured, heterogeneous data. Preprocessing steps such as feature scaling and feature selection were implemented to enhance model performance and generalization. Model evaluation relied on standard metrics like accuracy, precision, recall, and F1-score, alongside cross-validation techniques to ensure reliable performance estimates.

The integration of clustering and classification methods allowed us to both explore the underlying data structure and assess predictive capabilities, ultimately supporting more precise profiling and targeted interventions in the context of financial literacy.

The following sections detail the methodologies, implementation choices, and results obtained, providing a comprehensive view of the machine learning workflow applied in this study.

2.1 Data Preparation

The initial data preparation phase involved careful handling of missing values and encoding of categorical variables, essential steps to make the dataset compatible with machine learning techniques. Specifically, the **Age** column was converted to a numeric format, coercing any non-numeric entries into missing values (**NaN**). Then, numeric and categorical columns were distinguished.

Missing values in numeric columns were imputed using the median, a choice robust against potential outliers. For categorical columns, missing values were replaced with the constant string **Unknown**, preserving the information about missingness without discarding entire rows.

Regarding the encoding of categorical variables, a selective approach was applied: only columns with fewer than 50 unique categories were encoded using **LabelEncoder**. This transformation converts textual categories into integer values, enabling the variables to be used by machine learning models.

In total, 216 categorical columns were encoded, including relevant variables such as:

- **Gender** with 2 unique categories,
- **Geographic_Area** with 5 categories,
- **Work_Sector** with 10 categories,
- **Income_Bracket** with 5 categories.

After these operations, the resulting dataset consists of 4862 observations and 219 columns.

Next, the classification target was explicitly defined as the variable **Gender**. All numeric features except the target itself were considered available for training, resulting in a total of 217 numeric features. The first ten features included:

- Age, qd7_a, Geographic_Area, Province, qd5_1, qd5_2, qd5_3, qd5_4, qd5_5, qd5_6.

The presence of valid targets was verified for both classification tasks (represented by **Gender**) and regression (with the variable **Age**).

Given the sufficient number of features and the availability of the classification target, the gender prediction task was set up. Rows with missing target values were removed, yielding a final dataset of 4862 samples. Class distribution was also analyzed to ensure balanced representation.

Feature selection was not applied at this stage. Although feature selection can help reduce dimensionality and improve generalization, in our case all available features were kept for the following reasons:

- The number of numeric features (217) was not excessive for the models used, which included Random Forest, Logistic Regression, Gradient Boosting, Support Vector Machines (SVM), Naive Bayes, and K-Nearest Neighbors.
- Retaining all variables ensures that no potentially useful information is prematurely excluded.
- Some machine learning models, such as tree-based ones, can naturally handle the presence of irrelevant or less informative features.
- Keeping all features contributes to the robustness and accuracy of the models by allowing them to exploit the full information content of the data.
- The goal was to evaluate the model's performance on the full feature space before considering any dimensionality reduction techniques.

For training and evaluation, the dataset was split into a training set (80%) and a test set (20%), stratifying to preserve the distribution of the **Gender** variable. The training set contains 3889 samples, while the test set contains 973 samples, both with 217 features. To optimize performance, data were standardized by scaling to zero mean and unit variance, with the scaler fitted on the training set and then applied to the test set.

2.2 Construction and evaluation of supervised models

If the matrix `X_selected_class`, containing the features selected for classification, we proceed to train and evaluate multiple classification models. The models considered are Random Forest, Logistic Regression, Gradient Boosting, Support Vector Machine (SVM), Naive Bayes, and K-Nearest Neighbors (KNN).

For each model, training is performed either on the scaled or unscaled training data depending on the model type. In particular, tree-based models such as Random Forest and Gradient Boosting are trained on unscaled data, while models like Logistic Regression, SVM, Naive Bayes, and KNN are trained on scaled data to ensure optimal performance. After training, a 5-fold cross-validation is conducted on the training set to assess the stability and robustness of each model's performance. The final evaluation metric used to compare the models is the accuracy computed on the held-out test set.

The Random Forest algorithm is an ensemble method that builds multiple decision trees on bootstrapped subsets of the data and combines their outputs through majority voting. This method helps to reduce overfitting and improve generalization. Mathematically, the prediction of the ensemble can be written as:

$$\hat{y} = \text{majority_vote}(h_1(x), h_2(x), \dots, h_T(x)) \quad (1)$$

where $h_t(x)$ denotes the prediction of the t -th decision tree.

Logistic Regression, on the other hand, is a linear model that estimates the probability of a binary outcome using the logistic function. The model assumes a linear relationship between the input features and the log-odds of the target variable. The estimated probability is given by:

$$P(y = 1|x) = \frac{1}{1 + e^{-(\beta_0 + \beta^\top x)}} \quad (2)$$

where β is the vector of coefficients and β_0 is the intercept term.

Gradient Boosting is another ensemble method that builds a sequence of decision trees, where each tree is trained to correct the residual errors of the previous model. The general form of the model is expressed as:

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x) \quad (3)$$

where $h_m(x)$ is the m -th base learner and γ_m is the learning rate.

The Support Vector Machine (SVM) seeks the optimal hyperplane that maximizes the margin between the two classes. In the case of non-linearly separable

data, kernel functions are used to project the data into a higher-dimensional space. The optimization problem can be written as:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{subject to} \quad y_i(\mathbf{w}^\top \phi(x_i) + b) \geq 1 \quad (4)$$

where $\phi(x)$ is a kernel function mapping input features to a higher-dimensional space.

Naive Bayes is a probabilistic classifier based on Bayes' theorem, which assumes that features are conditionally independent given the class label. The prediction rule is given by:

$$P(y|x_1, \dots, x_n) \propto P(y) \prod_{i=1}^n P(x_i|y) \quad (5)$$

Despite its simplicity, this model can be effective in high-dimensional spaces and when the independence assumption approximately holds.

Finally, K-Nearest Neighbors (KNN) is a non-parametric method that assigns a class label to a sample based on the majority class among its k closest training examples. The predicted class is given by:

$$\hat{y} = \text{mode}(y_i \mid x_i \in \mathcal{N}_k(x)) \quad (6)$$

where $\mathcal{N}_k(x)$ denotes the set of k nearest neighbors of x in the training set.

The table below summarizes the accuracy and cross-validation results for each model:

Random Forest	: Accuracy = 0.7544, CV = 0.7457 \pm 0.0370
Logistic Regression	: Accuracy = 0.6639, CV = 0.6457 \pm 0.0389
Gradient Boosting	: Accuracy = 0.9476, CV = 0.9182 \pm 0.0199
SVM	: Accuracy = 0.6763, CV = 0.6698 \pm 0.0426
Naive Bayes	: Accuracy = 0.5920, CV = 0.5793 \pm 0.0547
K-Nearest Neighbors	: Accuracy = 0.6166, CV = 0.6153 \pm 0.0330

Among all models, the Gradient Boosting classifier achieves the best performance, with a test set accuracy of 94.76% and a cross-validation mean of 91.82%. This result suggests that Gradient Boosting is particularly well-suited to capture complex, non-linear relationships in the data. Although Random Forest also performs well, it does not reach the same level of precision. Simpler linear models such as Logistic Regression and Naive Bayes exhibit lower accuracy, likely due to their limited capacity to model interactions and non-linearities. SVM and KNN show moderate performance but are more sensitive to feature scaling and data sparsity. Overall, the empirical results support the selection of Gradient Boosting as the most effective supervised learning method for the classification task in this study.

2.3 Best Model Evaluation

Using the Gradient Boosting classifier, we generate predictions on the test set and compute a detailed classification report. The results show excellent performance across all metrics with precision, recall, and F1-score all approximately 0.95 for both classes.

Moreover, feature importance analysis reveals that variables such as `wght`, `Age`, and `Work_Sector` play a dominant role in the model's decisions.

		Predicted	
		0	1
Actual	0	464	26
	1	25	458

Figure 7: Confusion Matrix of the Gradient Boosting model on the test set.

The confusion matrix shows that out of the total test samples, 464 true negatives and 458 true positives were correctly classified. The model misclassified 26 samples as positive when they were actually negative (false positives), and 25 samples as negative when they were actually positive (false negatives). This indicates a balanced and precise classification performance.

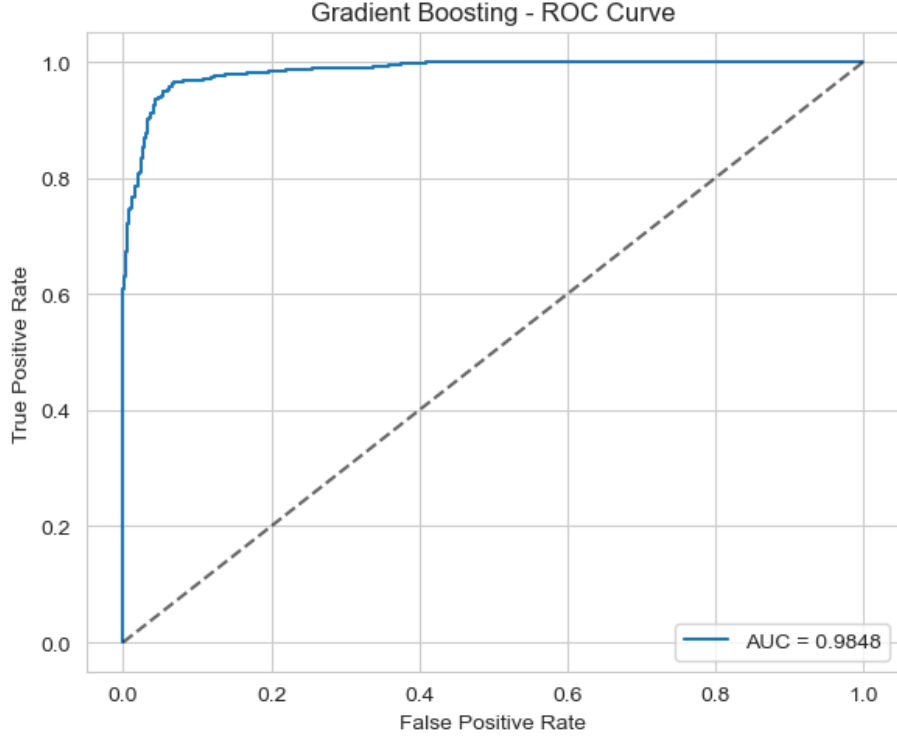


Figure 8: ROC Curve of the Gradient Boosting model on the test set. The Area Under the Curve (AUC) is 0.9848, indicating excellent discriminatory ability between the two classes.

The ROC curve further confirms the model's strong performance with an AUC of 0.9848, suggesting near-perfect separation of the classes. This high AUC complements the accuracy and precision metrics, reinforcing the effectiveness of Gradient Boosting in this classification task.

In summary, the Gradient Boosting classifier not only achieves the highest accuracy but also demonstrates robust predictive power and interpretability through feature importance and evaluation metrics, making it the preferred choice for the gender prediction task.

2.4 Clustering Analysis

The clustering analysis began by selecting all the numerical variables available in the dataset, resulting in a total of 218 features. To reduce complexity and noise, only the first ten most significant numerical features were used: **Gender**, Age, qd7_a, Geographic_Area, Province, qd5_1, qd5_2, qd5_3, qd5_4, and

qd5_5. After removing samples with missing values in these features, the final dataset for clustering comprised 4862 observations.

K-means clustering was applied with the number of clusters set to three (`n_clusters=3`). Prior to clustering, the features were scaled using the `StandardScaler` to ensure all variables contributed equally, avoiding bias due to different scales.

To facilitate visualization and interpretation, dimensionality reduction was performed using Principal Component Analysis (PCA), retaining the first two principal components. These two components explain approximately 19.0% of the total variance, enabling a two-dimensional representation that aids in understanding the cluster structure.

The K-means algorithm produced the following cluster distribution: Cluster 0 contains 1794 samples (36.9%), Cluster 1 is the largest with 2541 samples (52.3%), and Cluster 2 is the smallest with 527 samples (10.8%).

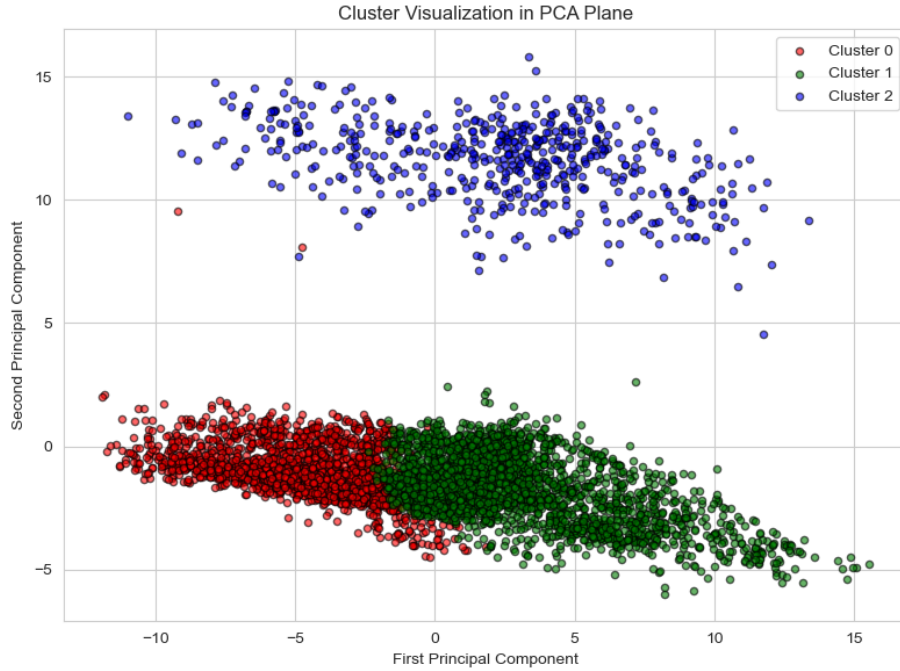


Figure 9: Clusters visualized in the two-dimensional PCA space defined by the first two principal components.

The scatter plot of the clusters in the PCA space (Figure 9) reveals a clear separation among the groups. Notably, Cluster 2 forms a distinct cloud positioned higher along the second principal component axis, with values approximately between 10 and 15, highlighting its unique profile relative to the rest of the data.

Clusters 0 and 1 are denser and located closer together in the lower region of the plot, with values between 0 and -5 on the second principal component axis. These two clusters are mainly separated along the first principal component, with Cluster 0 predominantly on the left and Cluster 1 on the right, showing a slight downward slope from left to right. This spatial arrangement suggests that while these two clusters share some similarities, they differ significantly on other characteristics, justifying their separation.

Overall, the clustering analysis successfully segmented the dataset into three distinct profiles, uncovering patterns of variability in numerical features that warrant further investigation to better understand group differences.

2.5 Clustering Evaluation

The quality of the clustering obtained with the K-means algorithm set to 3 clusters was assessed using the silhouette score, which measures how well each sample fits within its assigned cluster compared to other clusters.

The average silhouette score is 0.0788, a rather low value indicating poor internal cohesion within clusters and limited separation between groups. Silhouette scores close to zero suggest that many samples lie near the boundary between clusters, while negative or very low scores imply that some samples may be misclassified.

The silhouette plot (see Figure 10) further reveals that some components are likely misclassified. Specifically, cluster 2 shows a portion of silhouette values shifted towards the left (negative values), suggesting these samples might be better assigned to cluster 1. Similarly, cluster 1 has some values skewed towards the left, indicating a potential reassignment to cluster 0.

This overlap and unclear cluster boundary confirms the difficulty of the algorithm in clearly separating distinct groups, highlighting the need to reconsider the number of clusters, feature selection, or clustering approach.

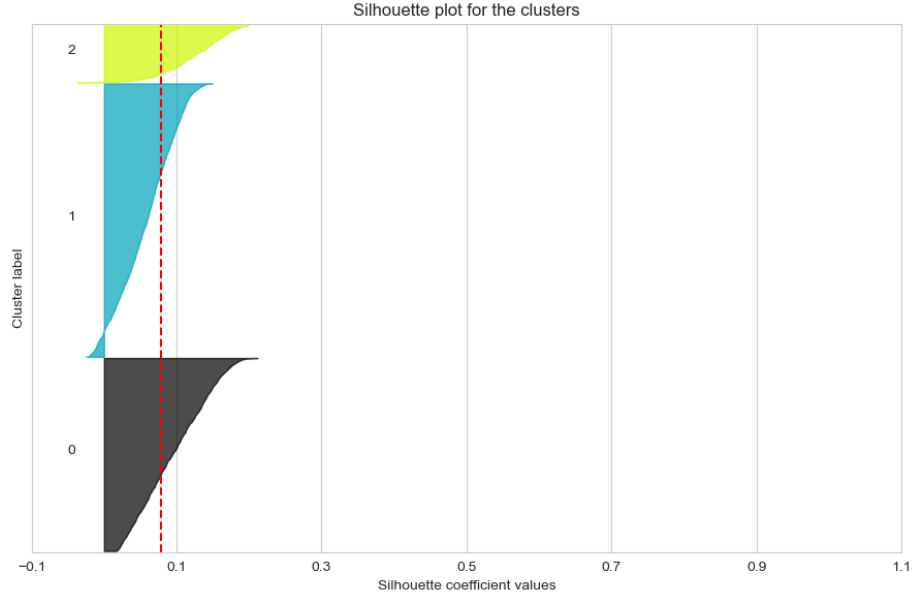


Figure 10: Silhouette plot of the clustering results for 3 clusters.

2.6 Second Clustering Attempt with Two Clusters

To further investigate the cluster structure and improve clustering quality, the K-means algorithm was rerun with the number of clusters set to two (`n_clusters=2`). The same scaled features and PCA transformation were used for consistency.

The two-cluster solution produced a slightly higher average silhouette score of 0.19, indicating an improved, yet still moderate, cluster separation compared to the three-cluster case.

Figure 11 shows the scatter plot of the two clusters in the PCA space, highlighting a broader grouping compared to the previous three-cluster solution.

The silhouette plot for the two-cluster solution (Figure 12) confirms the slightly better cluster compactness and separation, although the score remains relatively low, suggesting overlapping clusters or insufficient feature discrimination.

In conclusion, while the two-cluster approach marginally improves the clustering evaluation metric, the overall low silhouette scores suggest that the dataset's structure may be inherently difficult to segment into well-separated groups using K-means with the selected features.

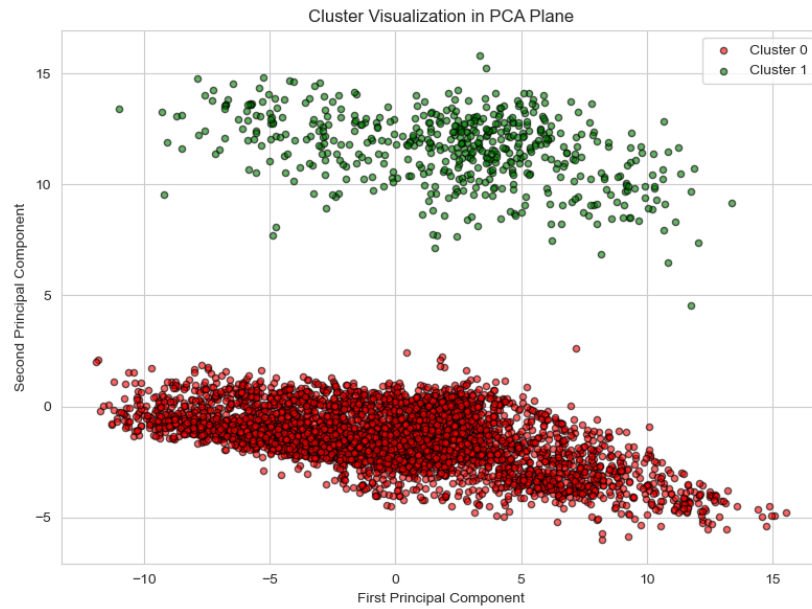


Figure 11: Clusters visualized in the two-dimensional PCA space for the two-cluster solution.

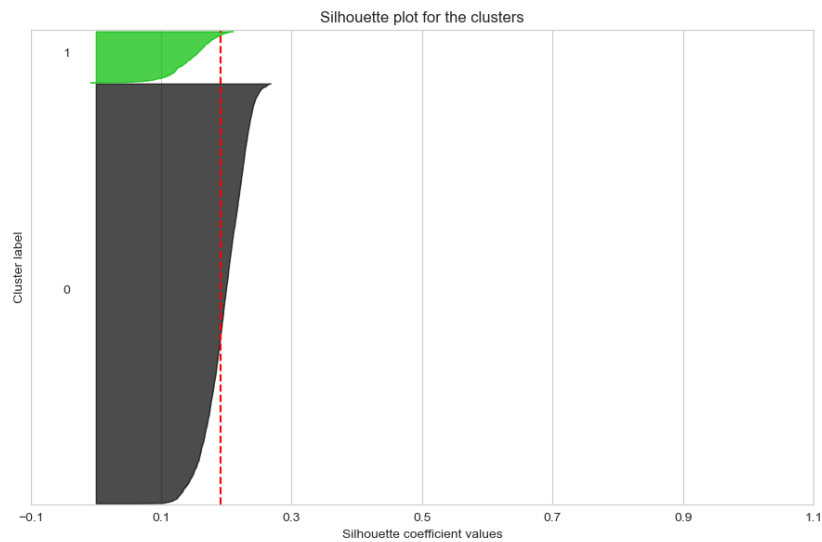


Figure 12: Silhouette plot of the clustering results for 2 clusters.

2.7 Clustering with Algorithms for Non-Spherical or Elliptical Structures

Given the limitations observed with K-Means, an alternative strategy was adopted based on clustering algorithms capable of identifying non-spherical, non-linear, or density-based structures. Specifically, the following methods were explored:

- **DBSCAN (Density-Based Spatial Clustering of Applications with Noise)**: this algorithm is based on point density and can identify arbitrarily shaped clusters, also distinguishing noise as outliers. It proved effective in detecting dense and well-separated groups, with a silhouette score of 0.262.
- **Gaussian Mixture Model (GMM)**: a probabilistic model assuming a Gaussian distribution for each cluster. Unlike K-Means, GMM allows for elliptical clusters and returns a probability of cluster membership. In this case, it yielded a silhouette score of 0.330.
- **Spectral Clustering**: an algorithm that leverages graph spectral theory to detect even non-convex clusters. It provided the best performance in terms of silhouette score, with a value of 0.452, indicating a good degree of cluster separation.

2.8 Visualization and Interpretation of Results

To better understand the nature of the resulting clusters and to compare the performance of the different algorithms, several informative visualizations were produced, as described below.

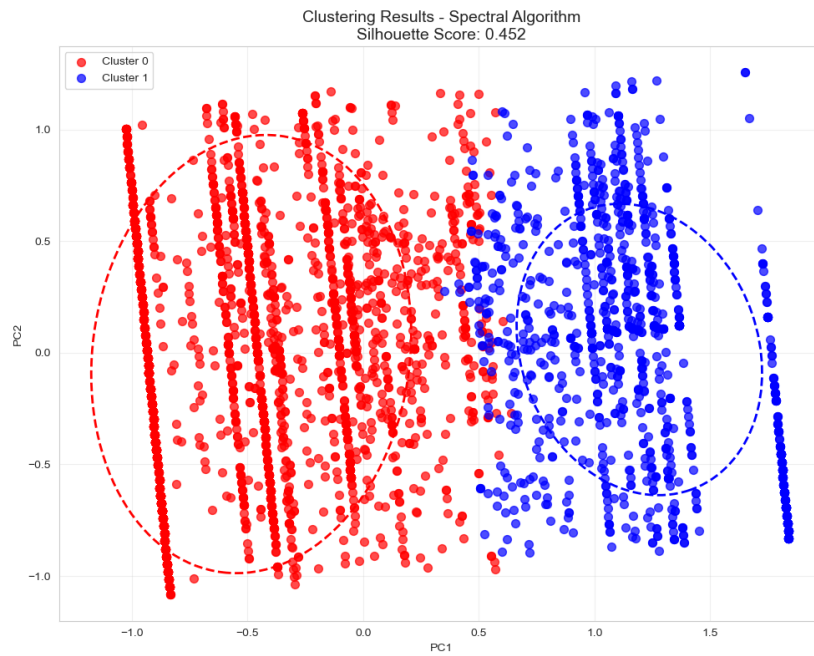


Figure 13: Distribution of data in clusters 0 and 1 according to Spectral Clustering, with a silhouette score of 0.452.

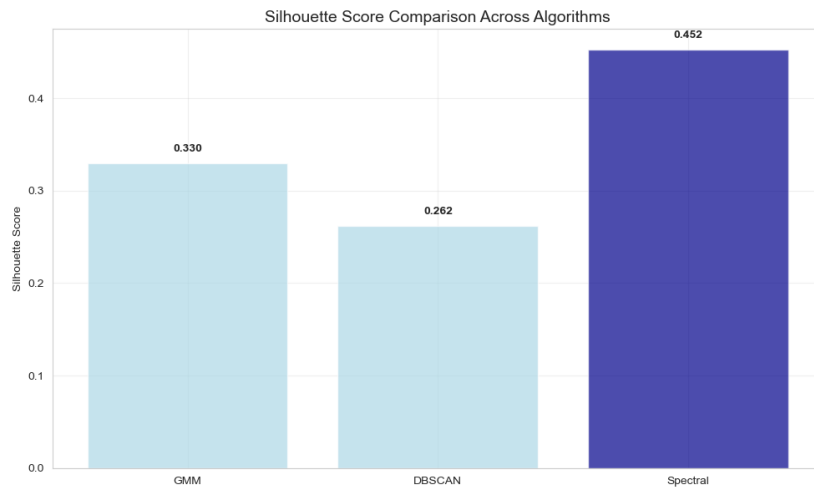


Figure 14: Comparison of silhouette scores for the three algorithm

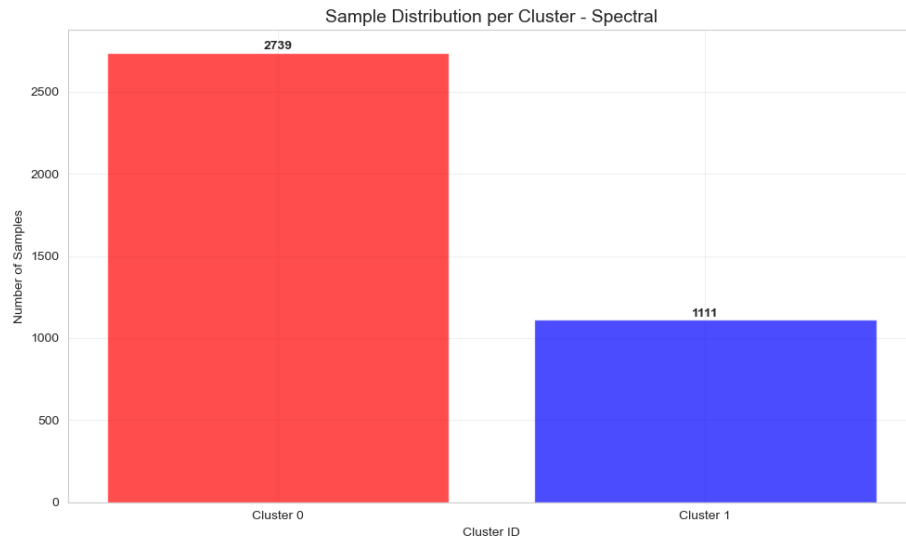


Figure 15: Sample distribution by cluster: cluster 0 contains 2739 observations, and cluster 1 contains 1111 observations.

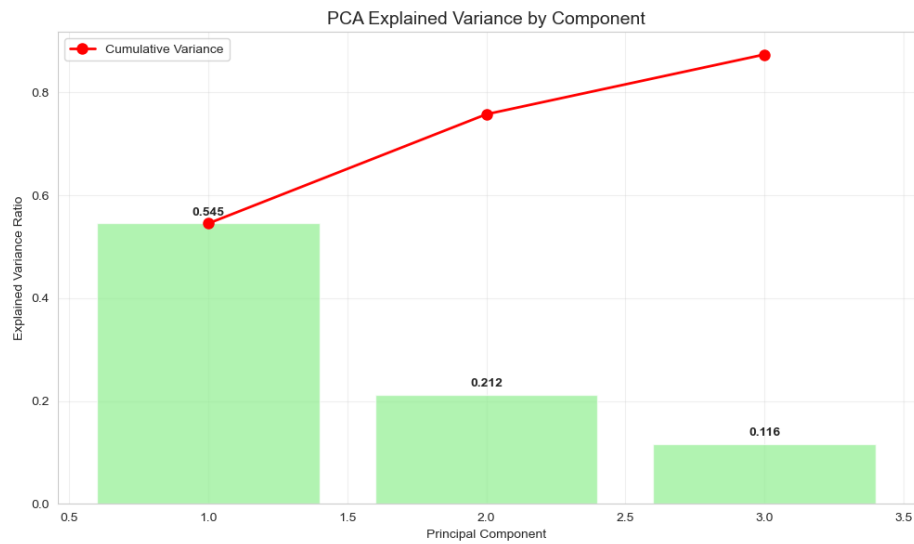


Figure 16: Explained variance per principal component. The first component explains 54.5% of the variance, the second 21.2%, and the third 11.6%. The total explained variance by the first three components is 87.3%.

The two clusters identified by the algorithm differ primarily in terms of the participants' age and the way they responded to the questionnaire items. One

group tends to include older individuals who gave higher scores on questions reflecting financial behaviors or attitudes, such as risk propensity or the ability to manage complex economic situations. The other group, by contrast, includes generally younger subjects who showed lower scores on the same questions, suggesting a potentially less experienced profile or one still in the process of developing financial literacy. This separation is not based on a single variable, but emerges from a combination of age, financial knowledge, and individual attitudes toward specific financial situations.

The analysis, supported by the high quality of the separation and the use of highly explanatory principal components, clearly shows that the algorithm distinguished two groups characterized by different levels of financial maturity and personal experience. The first cluster appears to represent a more structured, mature, and financially aware profile, while the second describes a younger segment of the population, potentially less prepared to effectively handle complex economic decisions. This type of segmentation may be extremely useful in designing targeted educational policies or specific interventions aimed at bridging knowledge gaps across different population segments.

In conclusion, the two clusters emerging from the analysis are not only statistically distinct but also interpretatively coherent with demographic and behavioral factors. The findings suggest that age and experience play a central role in shaping the responses, confirming that the segmentation performed by the algorithm has a strong and relevant foundation from an applicative perspective.

2.9 Final Remarks

This analysis highlighted how a multi-model approach, combined with visualization techniques and quantitative evaluation methods such as the silhouette score, is essential for achieving effective clustering on complex data. The failure of K-Means, documented both by performance metrics and visual inspection, led to an informed decision to apply more sophisticated alternatives. Among these, Spectral Clustering yielded the most reliable results in the context of financial literacy segmentation.

3 Final considerations

This project provided a comprehensive analysis of a financial literacy dataset through both supervised and unsupervised machine learning techniques. The dual approach allowed for meaningful insights into demographic and behavioral patterns within the Italian adult population.

In the supervised learning phase, several classification models were evaluated with the goal of predicting the respondent’s gender based on financial and socio-economic variables. Among the models tested, Gradient Boosting achieved the best results, with an outstanding accuracy of 94.76% and an AUC of 0.9848. This suggests a strong capacity to capture complex, non-linear relationships in the data. Feature importance analysis revealed that age, employment sector, and sampling weights significantly contributed to the predictive power of the model.

The unsupervised learning component aimed to explore the natural structure of the data without predefined labels. Initial attempts using K-Means clustering led to unsatisfactory results, as indicated by low silhouette scores and unclear cluster separation. Subsequent efforts with algorithms better suited to non-spherical distributions, including DBSCAN, Gaussian Mixture Models, and Spectral Clustering, yielded more promising outcomes. Spectral Clustering, in particular, stood out with a silhouette score of 0.452 and clearly differentiated two respondent profiles based on age and financial behavior, suggesting a meaningful segmentation between financially mature and less experienced individuals.

Overall, the project successfully demonstrated the value of combining different machine learning paradigms to explore complex socio-economic data. While the models provided interpretable and actionable insights, future enhancements could include dimensionality reduction through feature selection, the application of ensemble clustering methods, or the integration of time-aware or sequential data if available. These refinements would further strengthen the reliability and applicability of the findings in policy-making and financial education strategies.