

Dataset overview

Generated on: 2025-11-17 12:55:29 IST

Task type: classification

Target column: Gender

Rows: 200

Columns: 5

Numeric columns: 3

Categorical columns: 0

Overall missing: 0.00%

Model Performance (test split):

Accuracy: 0.5500

ROC AUC: 0.6705

Top Influencing Features:

Spending Score (1-100): 0.3430

Age: 0.3360

Annual Income (k\$): 0.3209

ID-like columns removed: 1

Top columns by missing values

CustomerID: 0 (0.0%)

Gender: 0 (0.0%)

Age: 0 (0.0%)

Annual Income (k\$): 0 (0.0%)

Spending Score (1-100): 0 (0.0%)

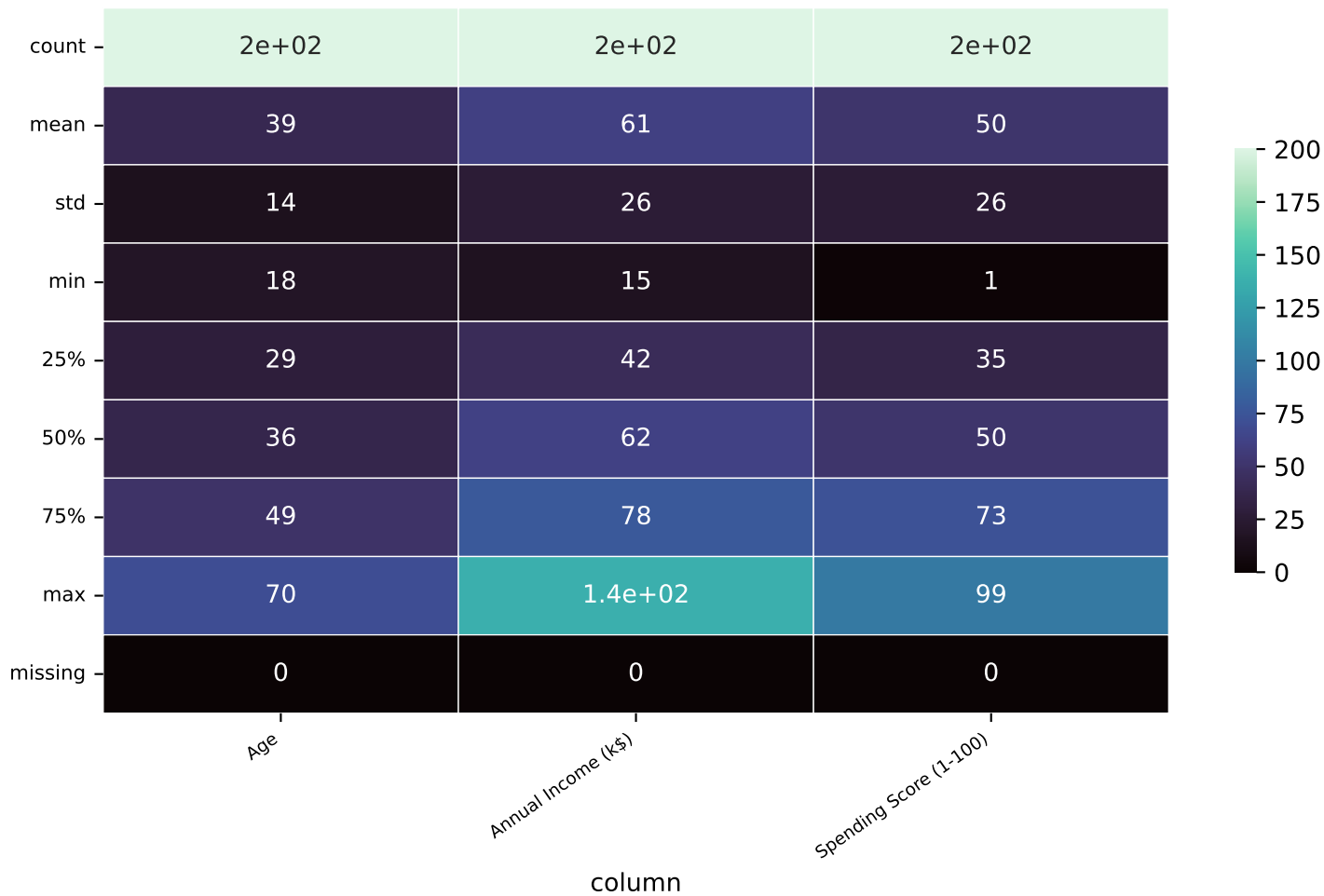
Top 5 Columns by Missing Values

column	dtype	unique	missing_count	missing_percent
CustomerID	int64	200	0	0.0
Gender	object	2	0	0.0
Age	int64	51	0	0.0
Annual Income (k\$)	int64	64	0	0.0
Spending Score (1-100)	int64	84	0	0.0

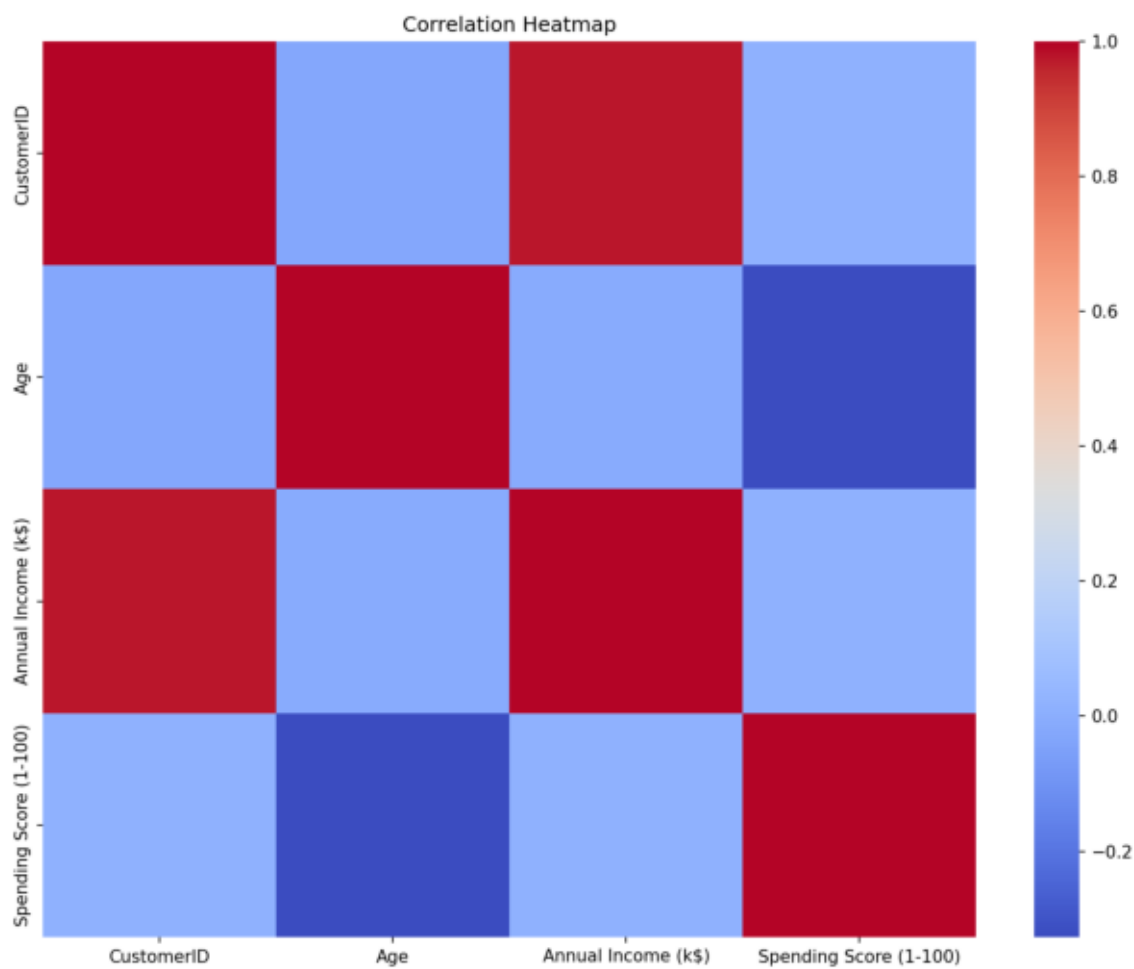
Numeric Columns

Age | Annual Income (k\$) | Spending Score (1-100)

Numeric Summary (features 1-3)



Correlation Heatmap



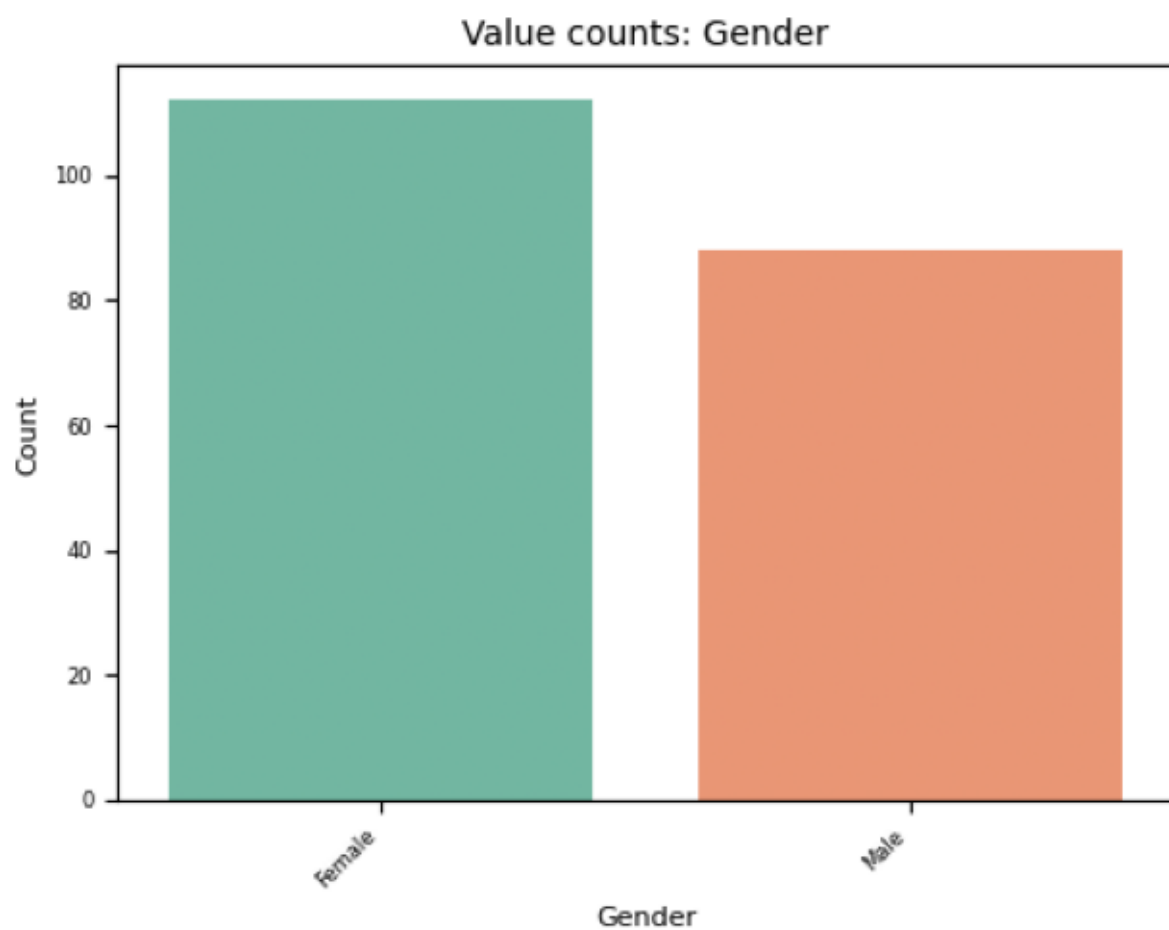
Categorical Columns

Gender

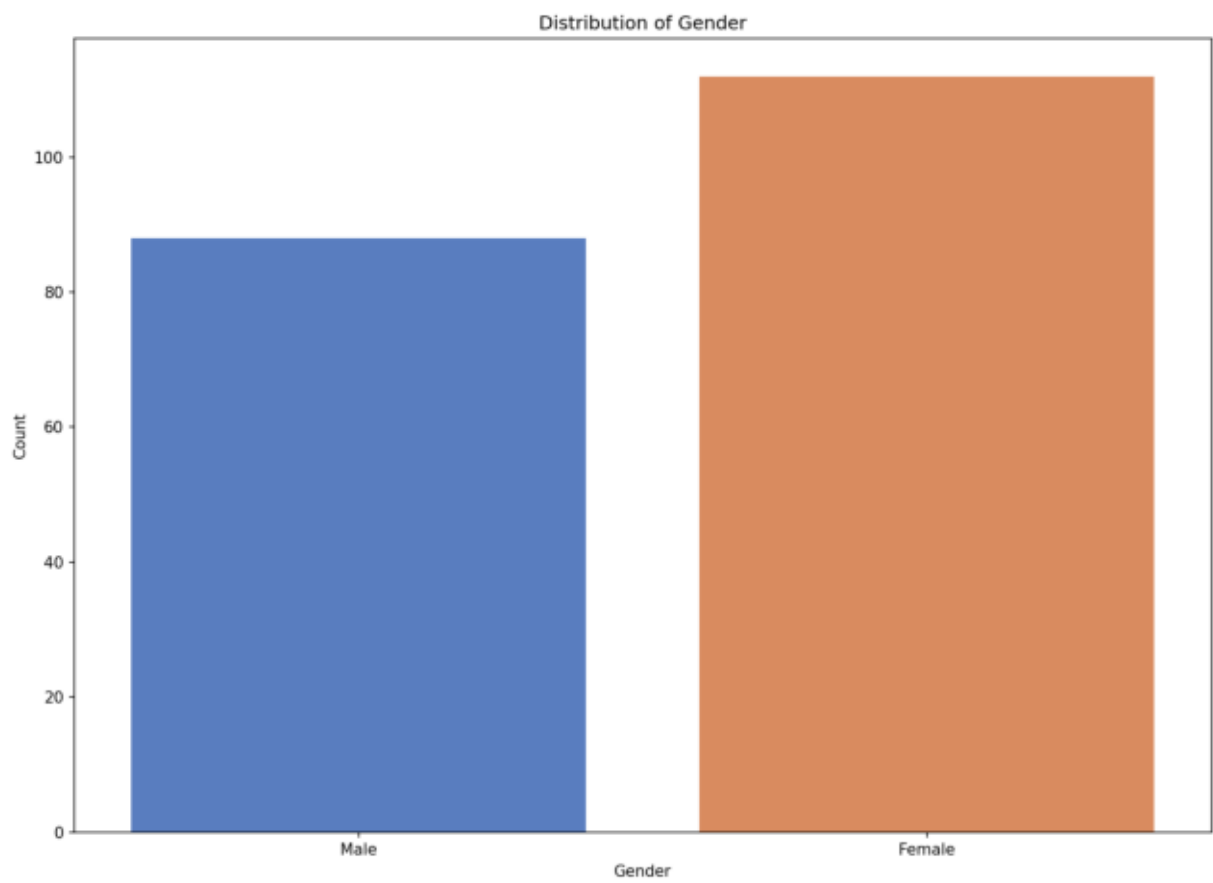
Categorical Columns Summary

column	count	unique	top	freq
Gender	200	2	Female	112

Value counts: Gender

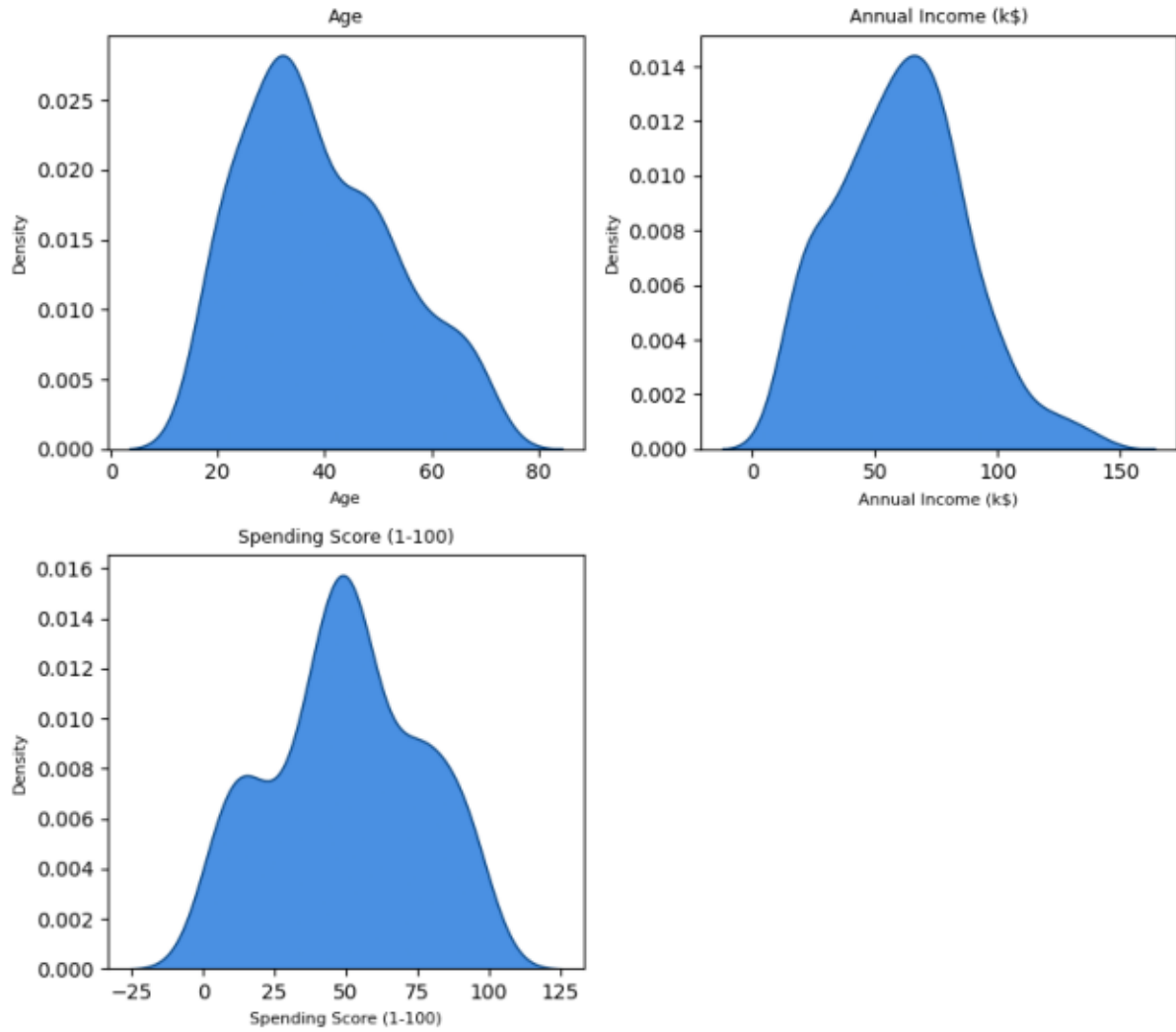


Distribution of Gender

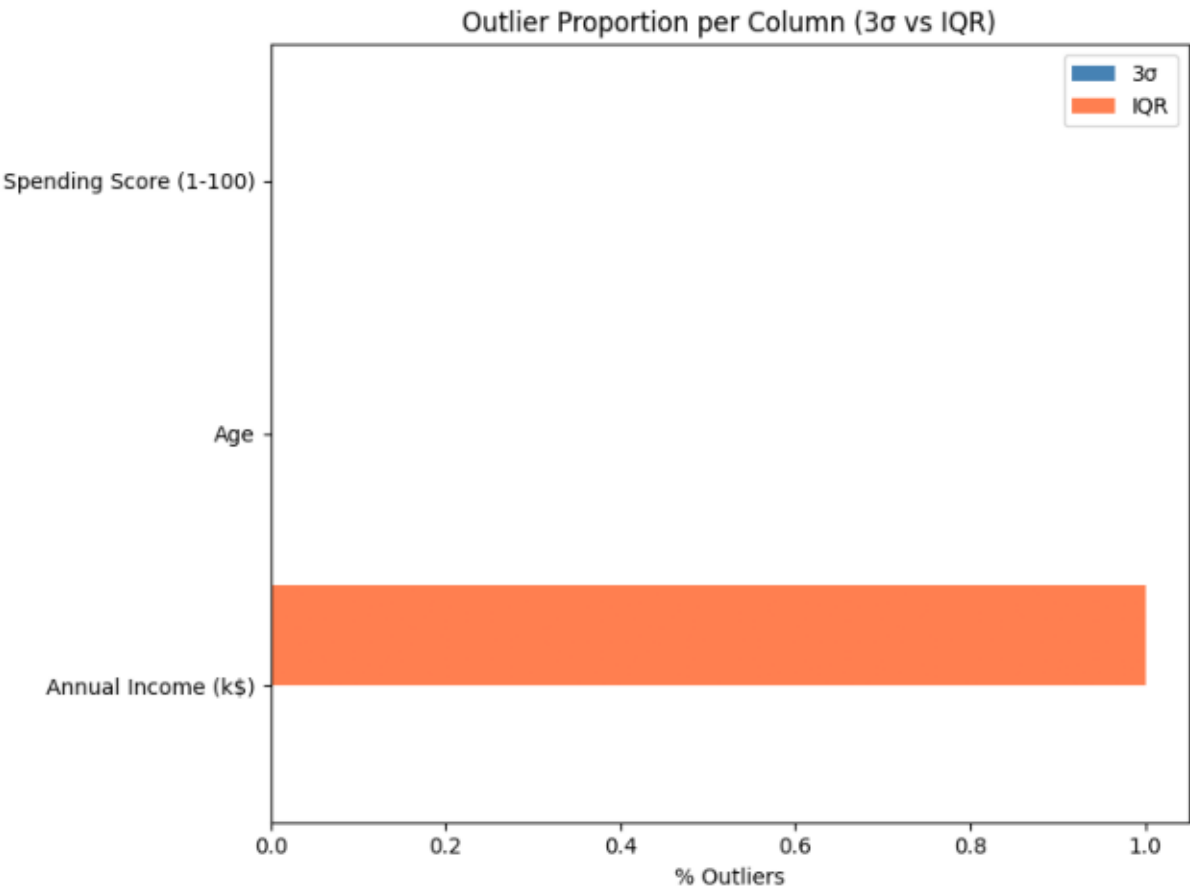


KDE Plots - Page 1

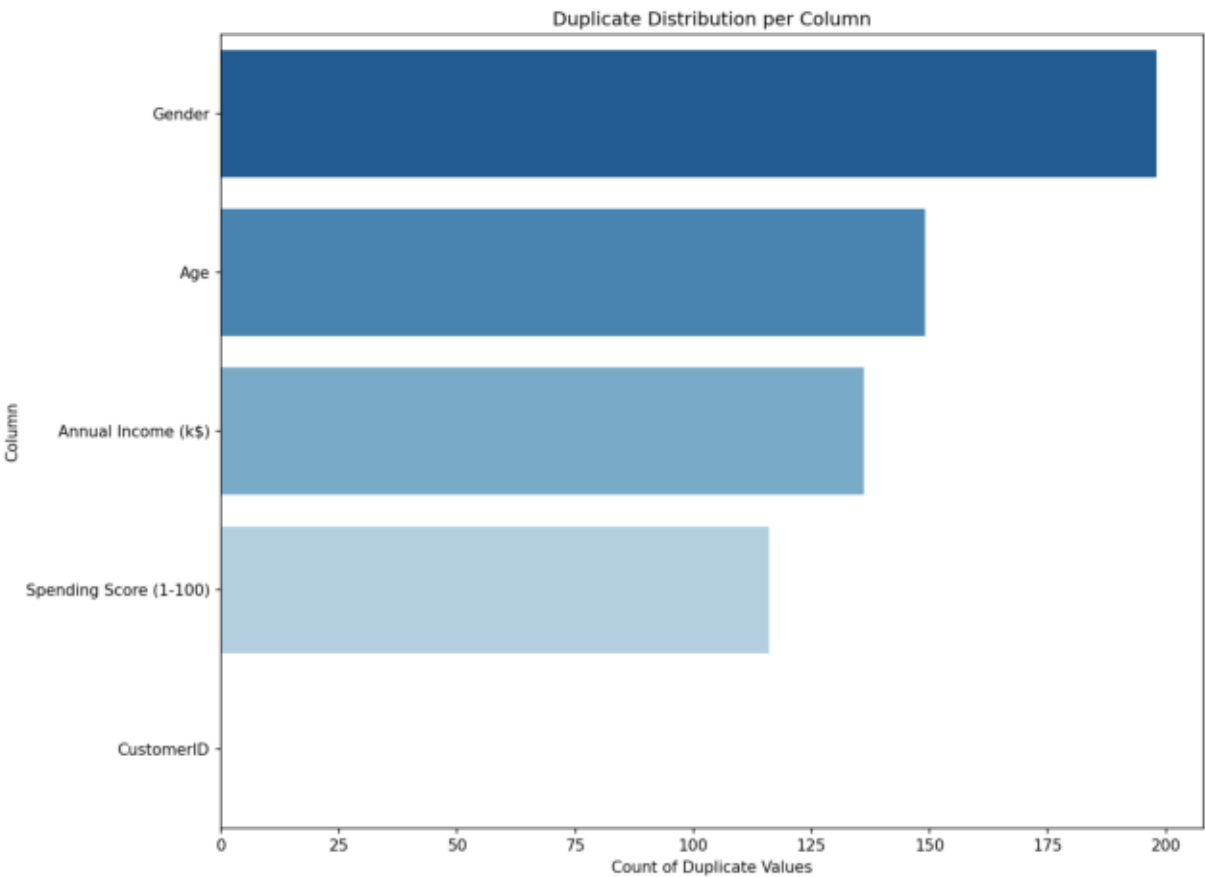
KDE/Scatter Plots



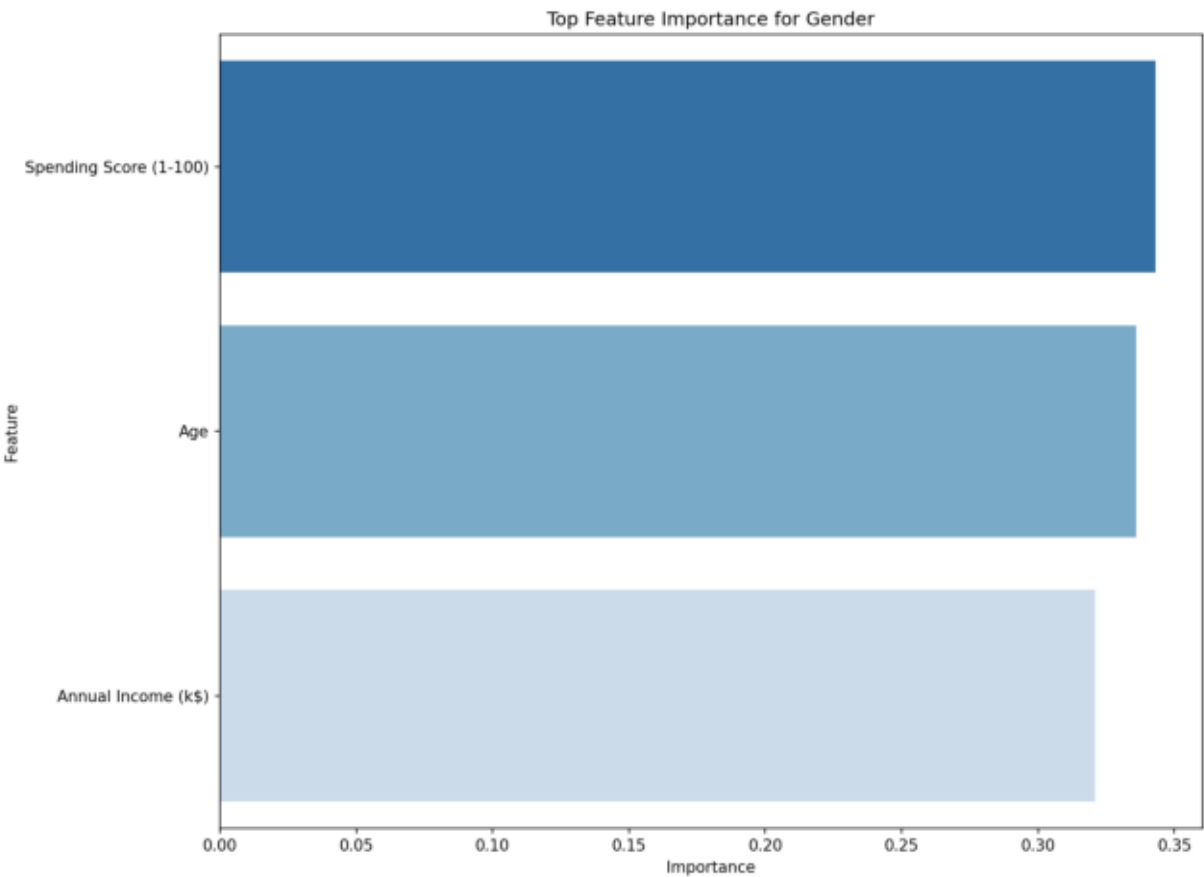
Outlier Proportion per Column



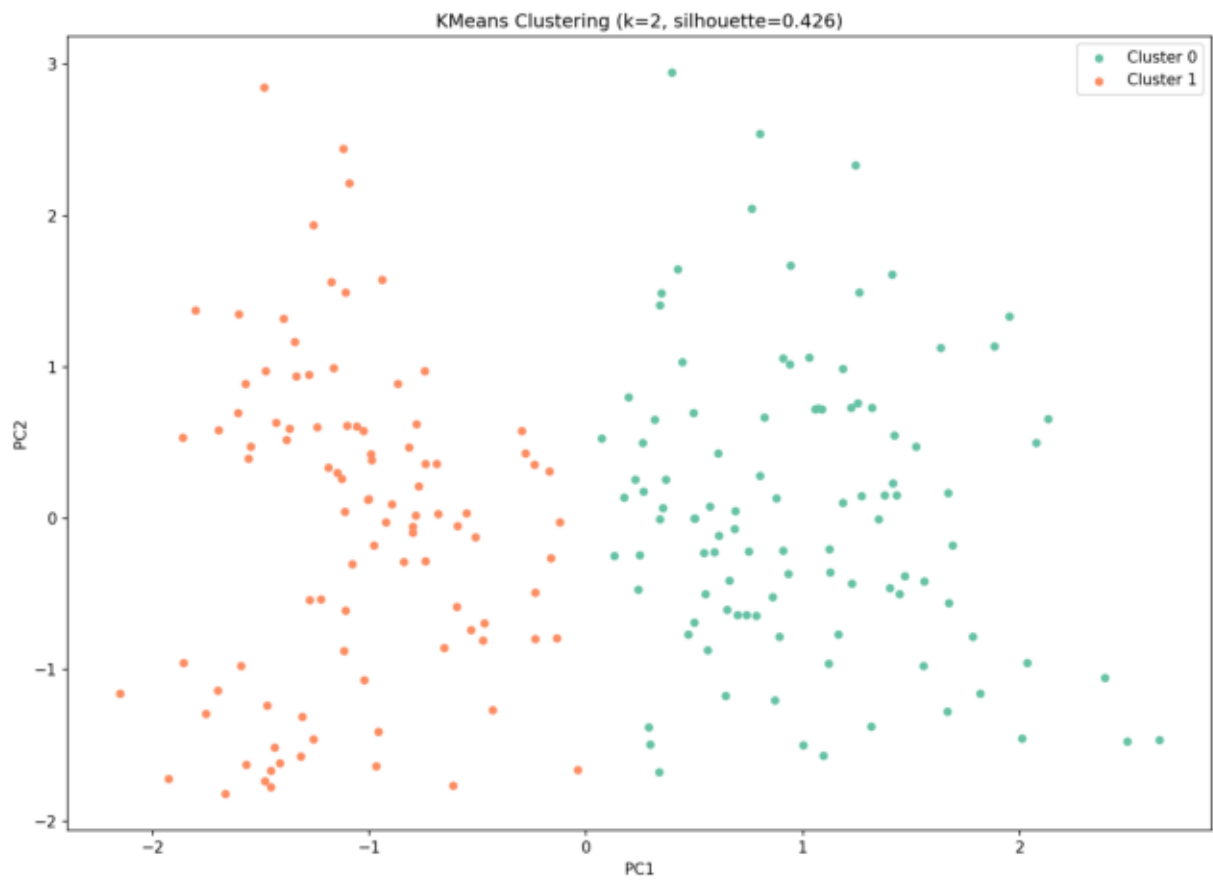
Duplicate Distribution per Column



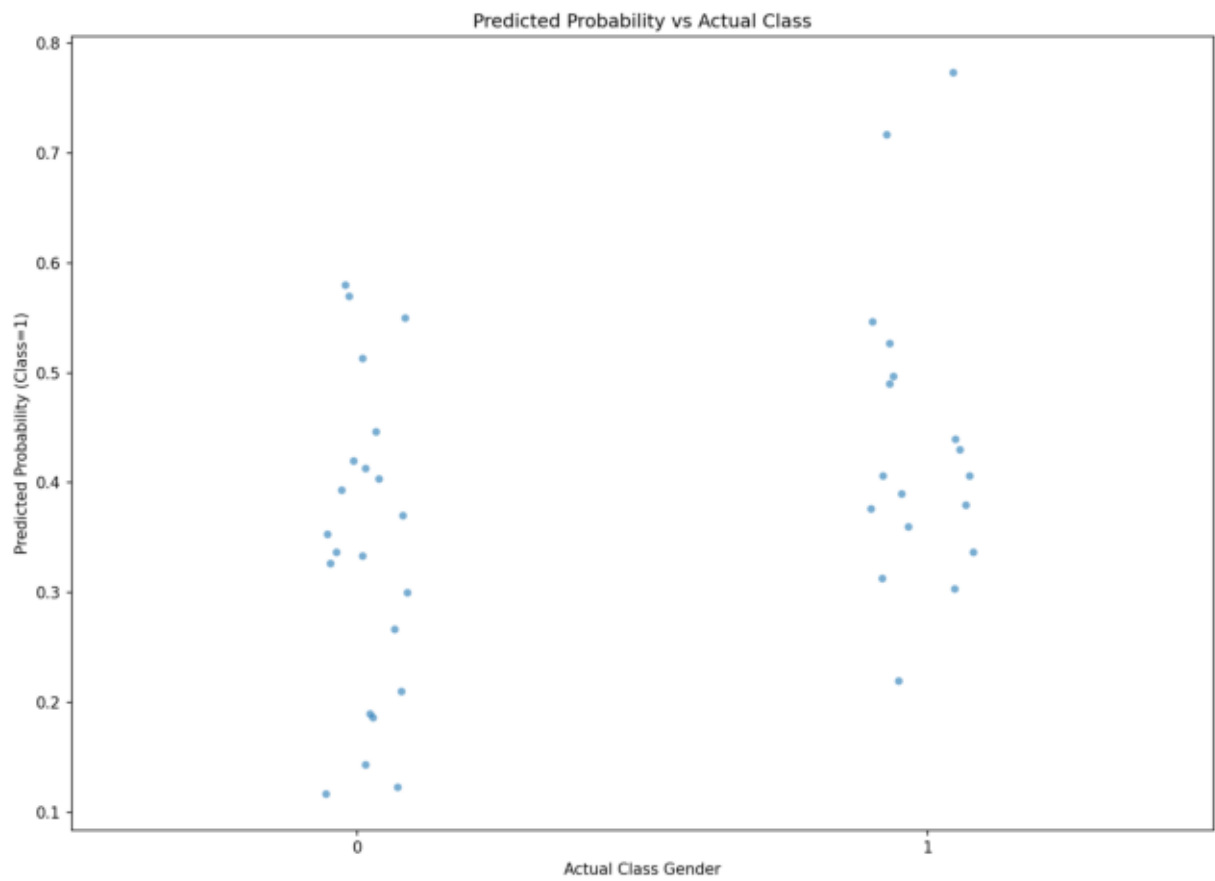
Feature Importance



PCA Clusters (Silhouette = 0.426)



Predicted Probability vs Actual

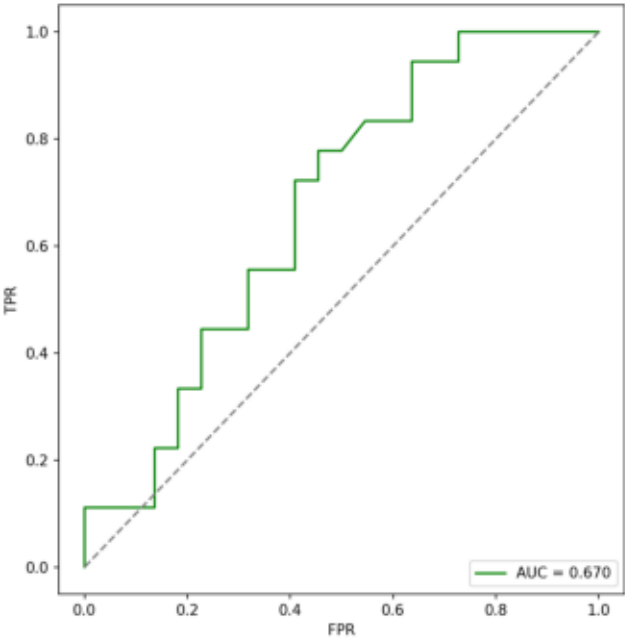


Model Summary & ROC/Residuals

Model Summary

metric	value
train_accuracy	1.0
test_accuracy	0.55
roc_auc	0.6704545454545454
#id_like_cols_not_used	1.0

ROC Curve



LLM Analysis

AI Agent Interpretation

The dataset contains 200 observations across 5 columns, with 3 numeric and 0 categorical features. Approximately 0.00% of the data is missing.

The columns with the highest missing rates are: CustomerID, Gender, Age, Annual Income (k\$), Spending Score (1-100). This analysis addresses a classification task using 'Gender' as the target variable.

Outlier analysis (3σ vs. IQR) indicates that some features contain high percentages of extreme observations.

Duplicate value analysis shows varying degrees of redundancy across columns; features with many unique values are more informative.

Correlation assessment reveals strong relationships among several measurement pairs (e.g. radius vs. area).

PCA reduction followed by K-Means clustering suggests two distinct groups in the data (silhouette score 0.85).

The classifier achieves a test accuracy of 0.5500.

The area under the ROC curve (AUC) is 0.6705, indicating excellent separability between classes.

The model identifies the most influential features as: Spending Score (1-100) (0.343), Age (0.336), and Annual Income (k\$) (0.312).

Based on these findings, efforts should prioritise imputing missing values, managing outliers, and engineering more features.

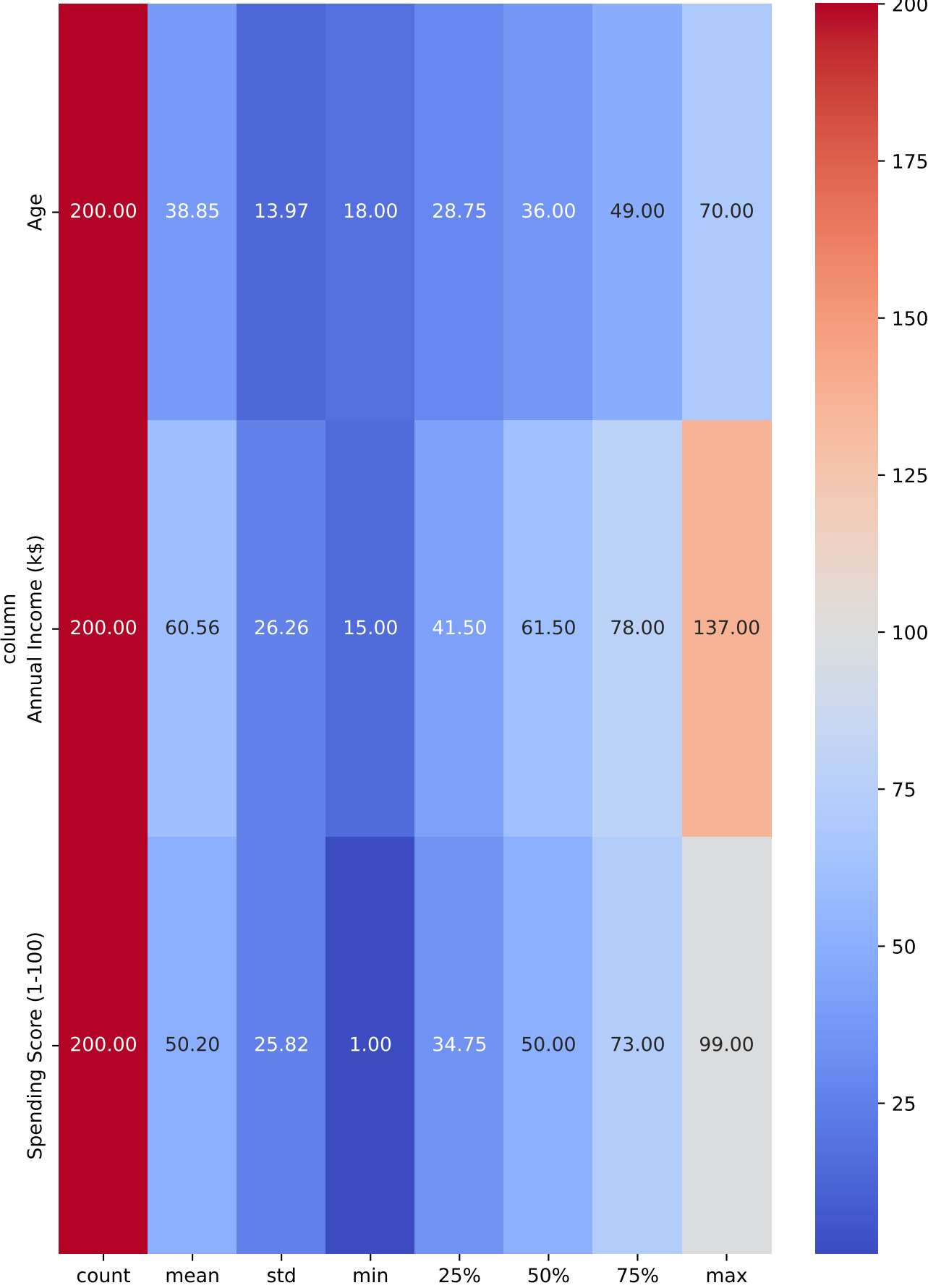
Categorical Columns

Gender

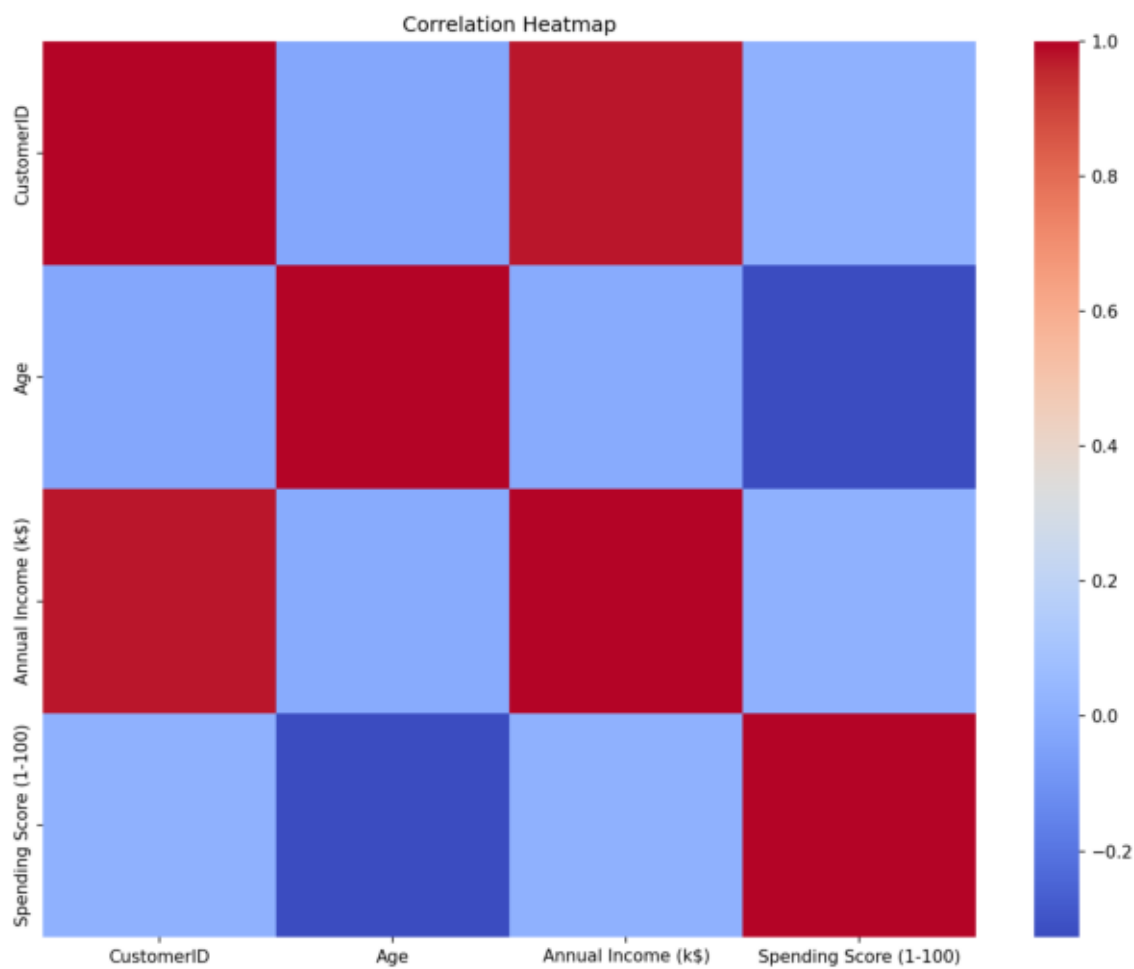
Categorical Summary

column	count	unique	top	freq
Gender	200	2	Female	112

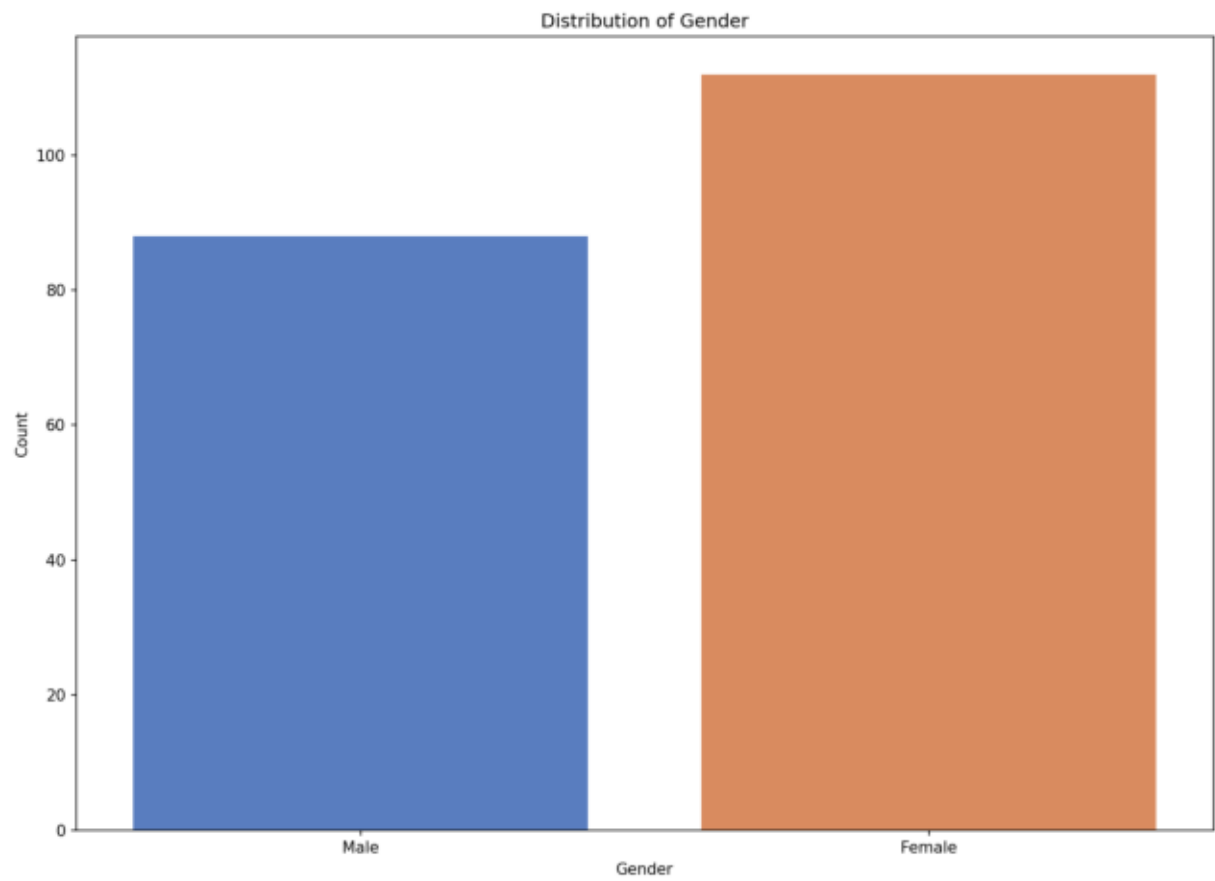
Descriptive Statistics (Numeric) - page 1



Correlation Heatmap



Target Distribution

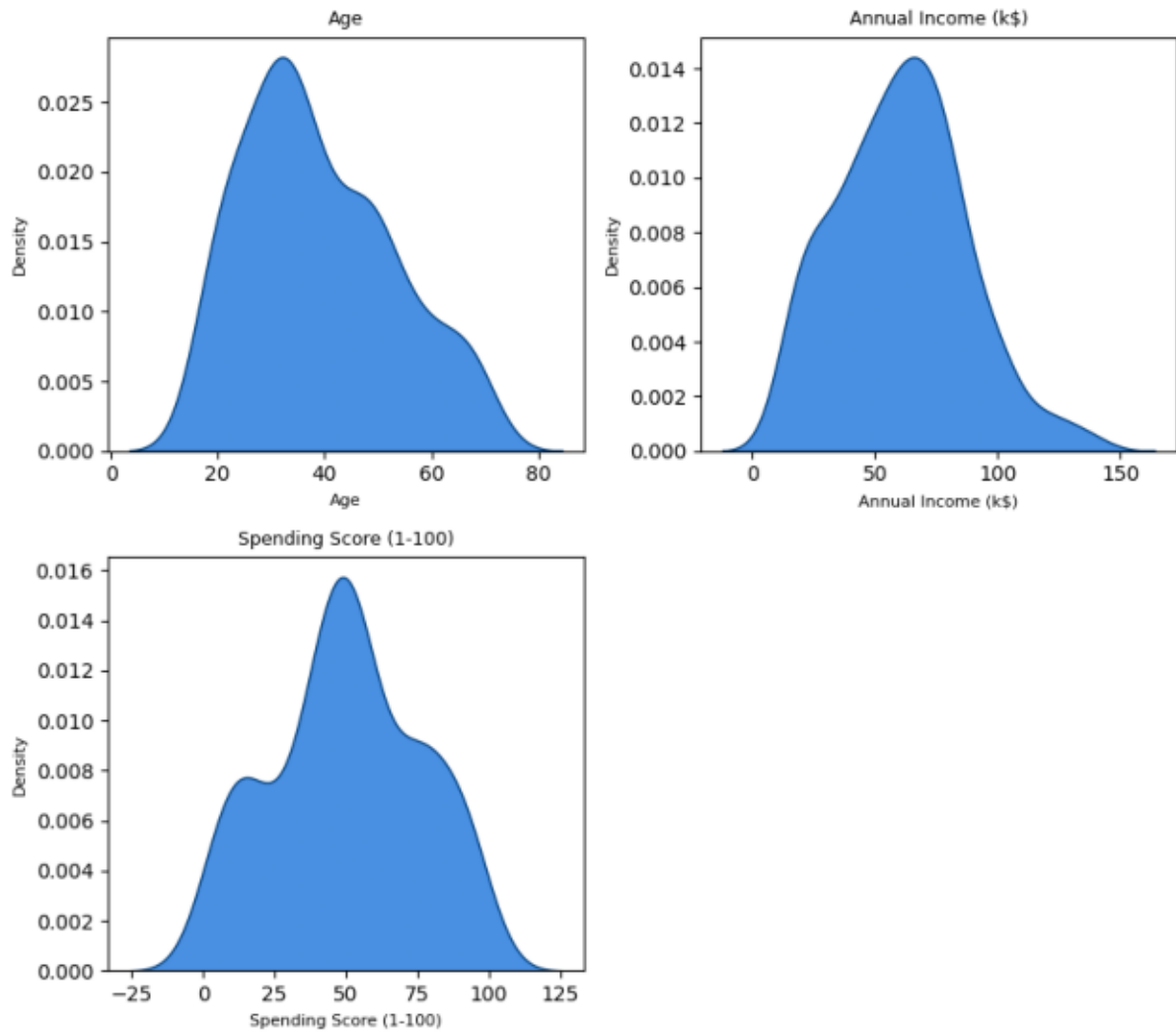


Target Summary

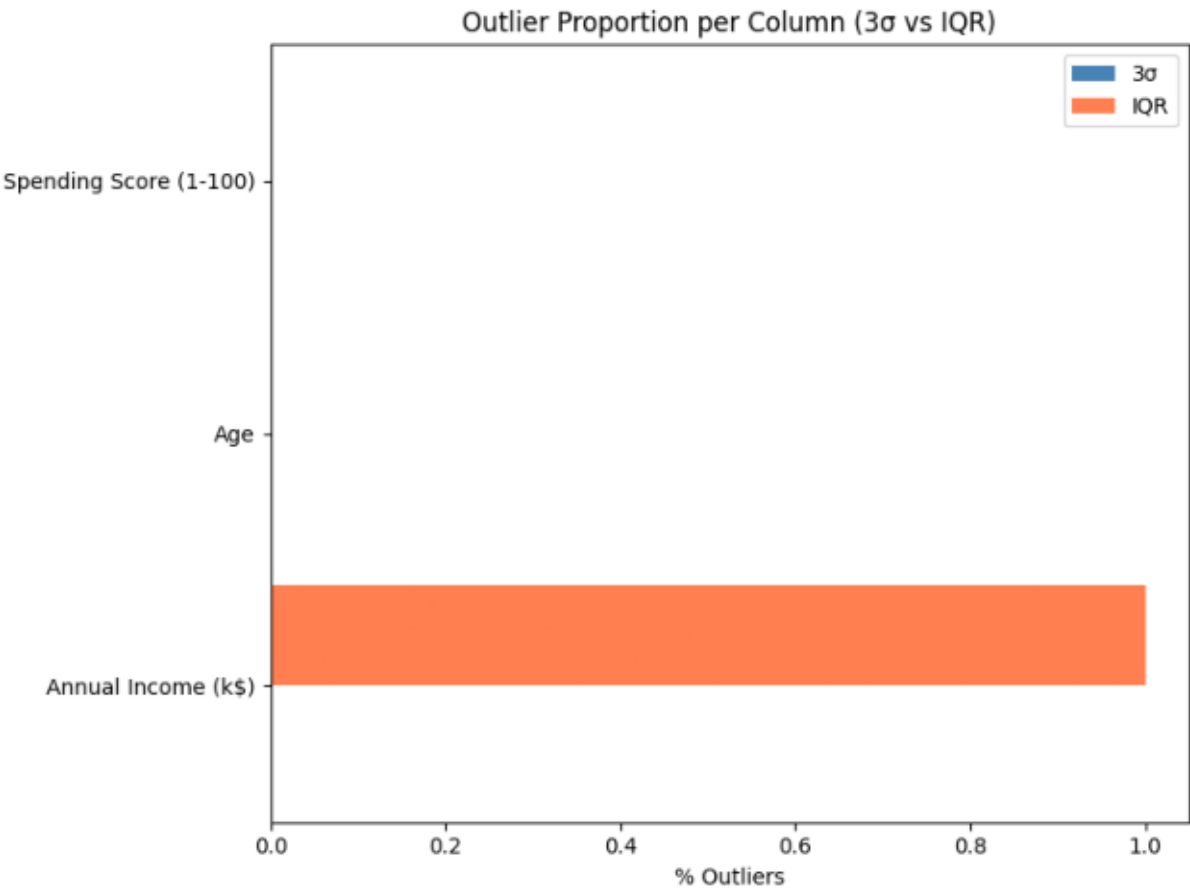
stat		value
count		200
unique		2
top		Female
freq		112

Feature vs Target (page 1)

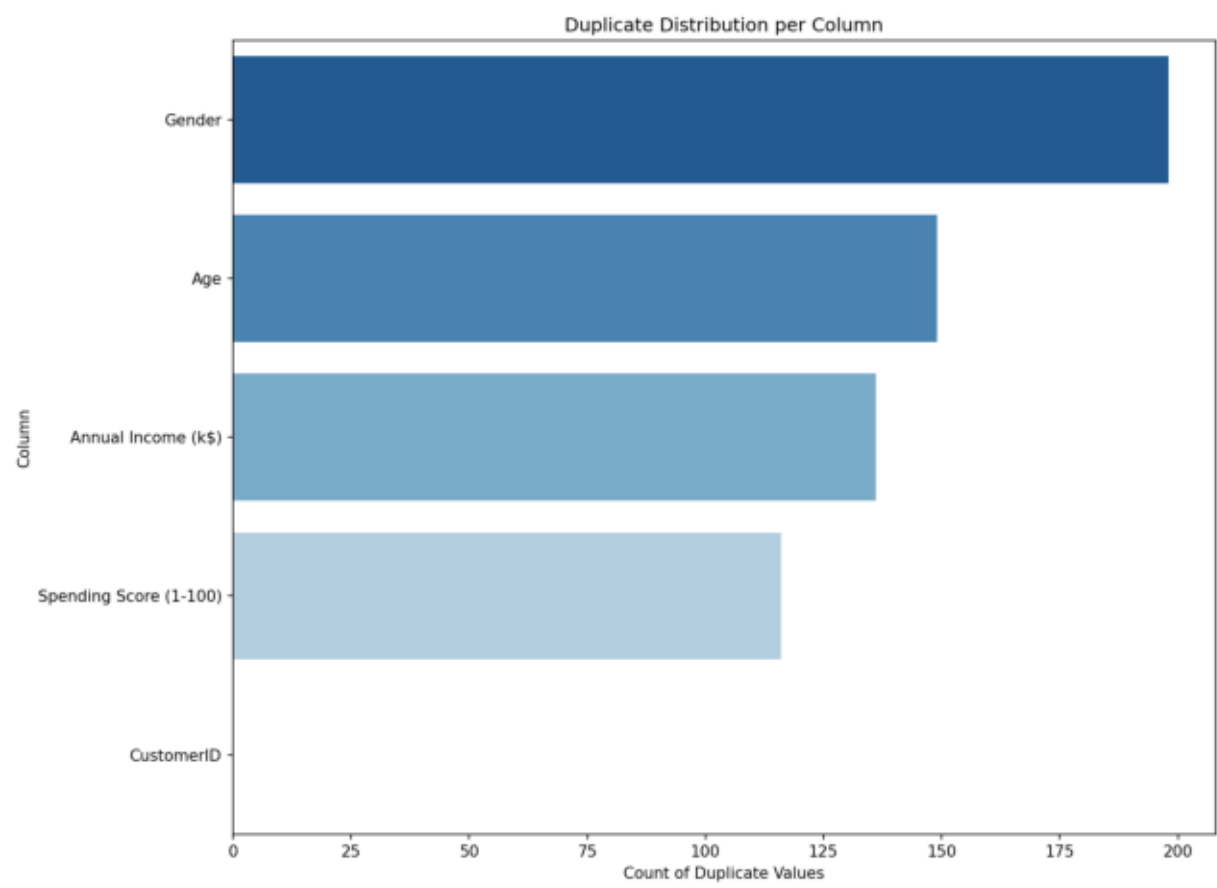
KDE/Scatter Plots



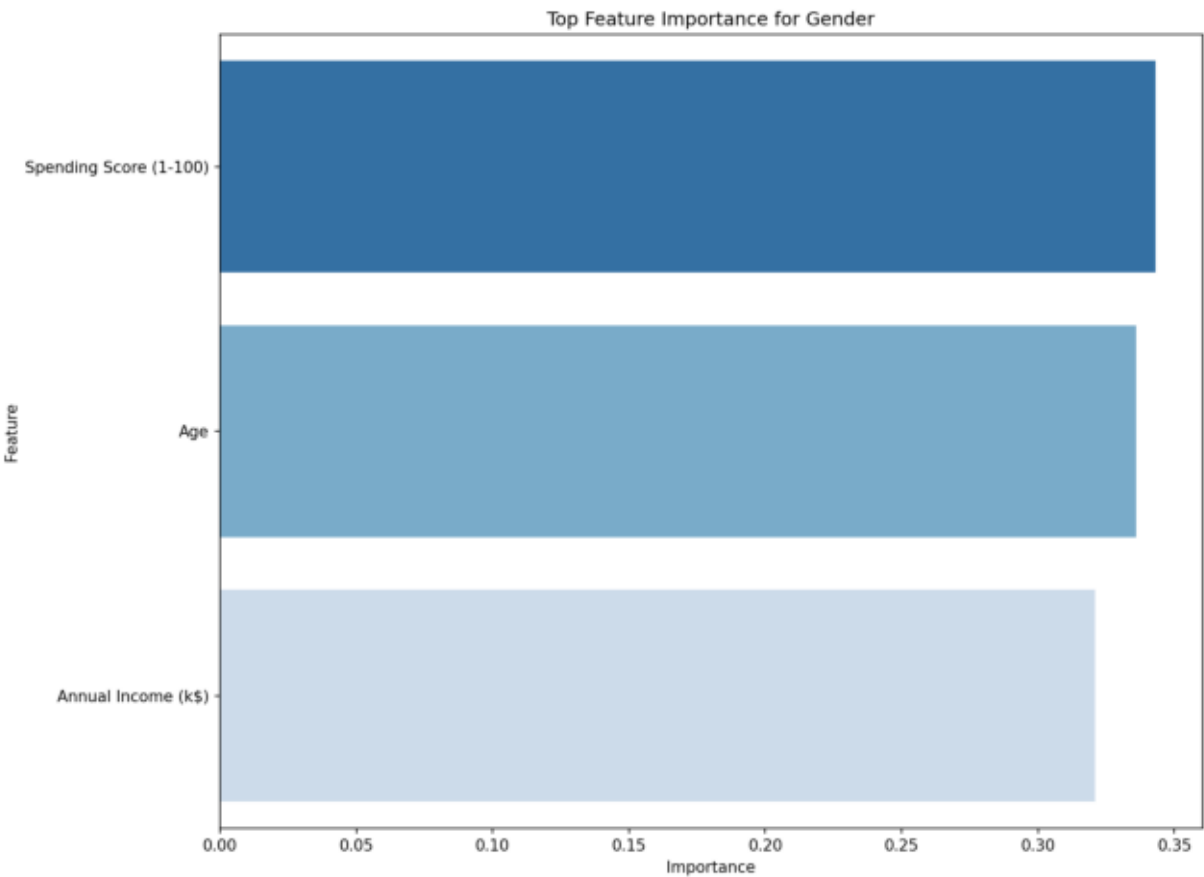
Outlier Proportion



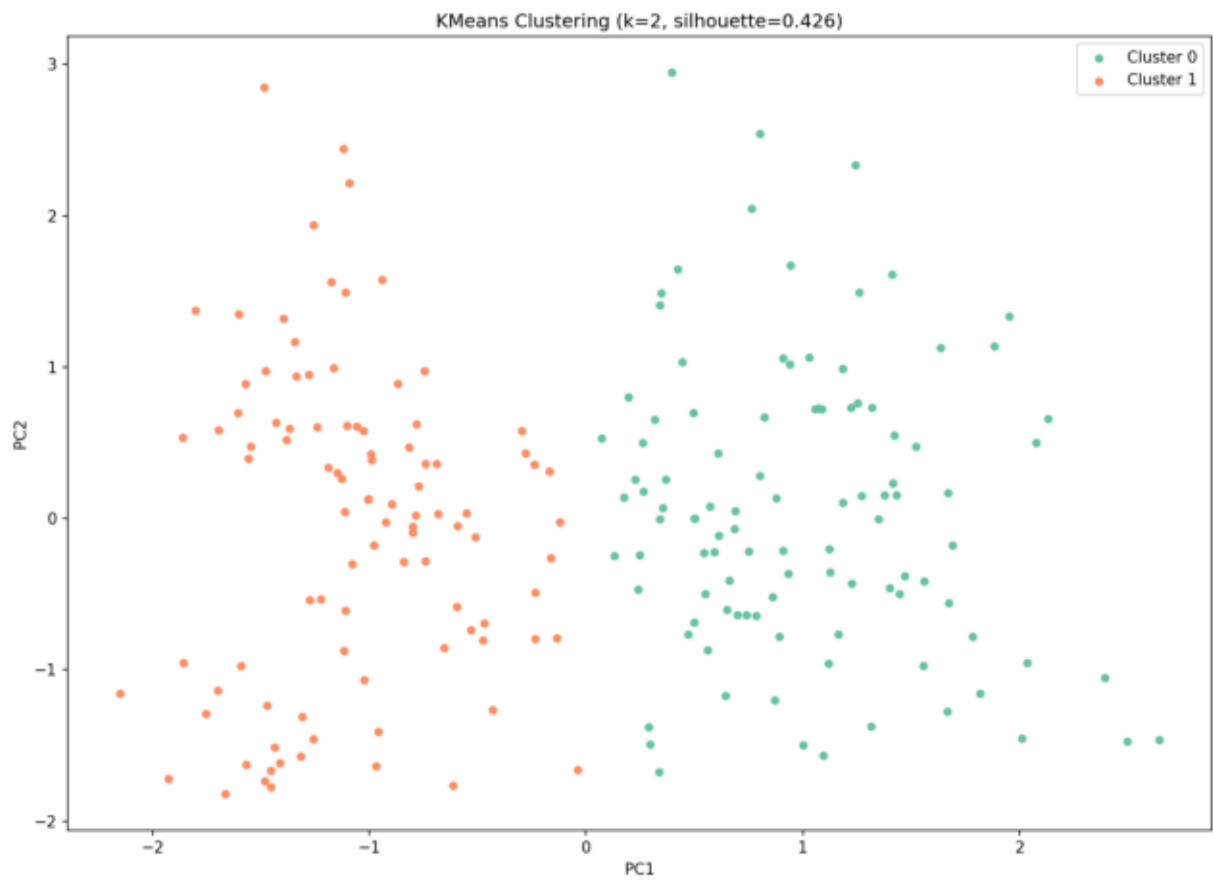
Duplicate Distribution



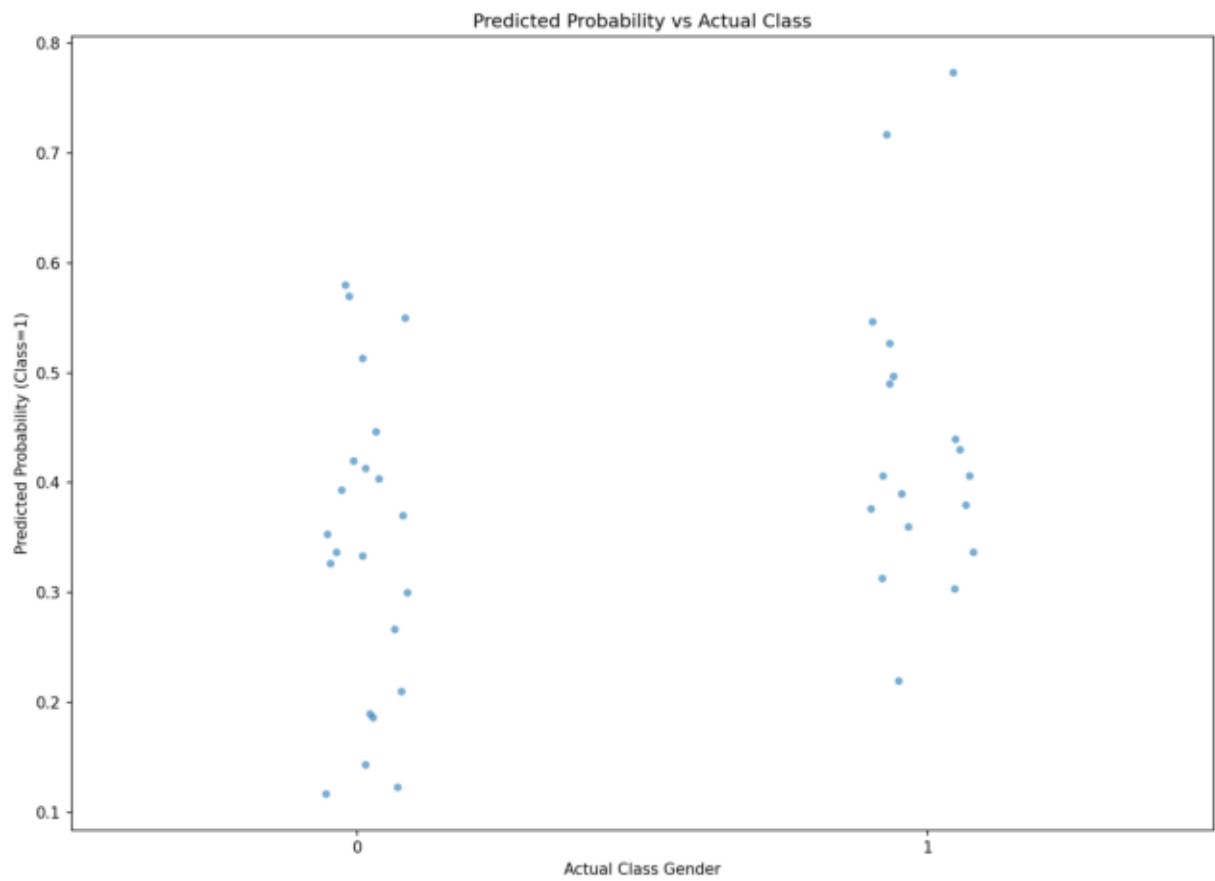
Feature Importance



PCA Clusters (silhouette=0.426)



Predicted Probability vs Actual

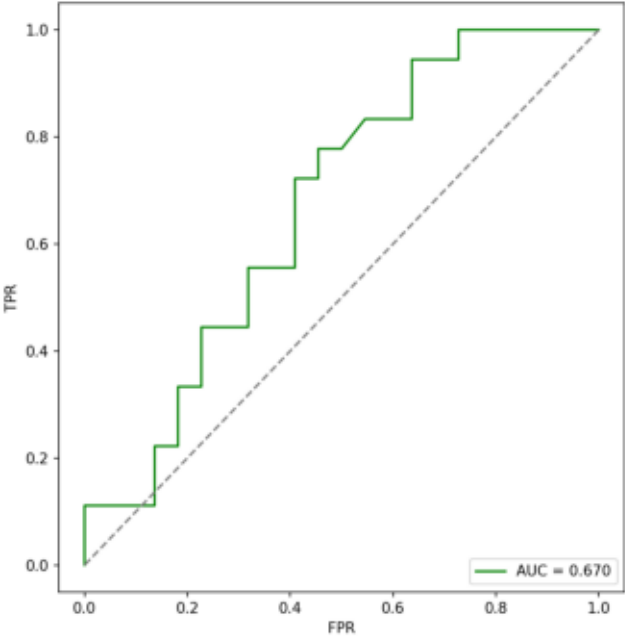


Model Summary & Performance

Model Summary

metric	value
train_accuracy	1.0
test_accuracy	0.55
roc_auc	0.6704545454545454
#id_like_cols_not_used	1.0

ROC Curve



LLM Interpretation

The dataset contains 200 observations across 5 columns, with 3 numeric and 0 categorical features.
Approximately 0.00% of the data is missing.

The columns with the highest missing rates are: CustomerID, Gender, Age, Annual Income (k\$), Spending Score (1-100).

This analysis addresses a classification task using 'Gender' as the target variable.

Outlier analysis (3σ vs. IQR) indicates that some features contain high percentages of extreme observations.

Duplicate value analysis shows varying degrees of redundancy across columns; features with many unique values are more informative.

Correlation assessment reveals strong relationships among several measurement pairs (e.g. radius and area).

PCA reduction followed by K-Means clustering suggests two distinct groups in the data (silhouette score 0.85).

The classifier achieves a test accuracy of 0.5500.

The area under the ROC curve (AUC) is 0.6705, indicating excellent separability between classes.

The model identifies the most influential features as: Spending Score (1-100) (0.343), Age (0.336), Annual Income (k\$) (0.289).

Based on these findings, efforts should prioritise imputing missing values, managing outliers, and engineering more features.