# Customer Segmentation Classification

Violet Huang: Data Preprocess, Feature Encoding, Data Exploration, Data Analysis
April Wang: Feature Encoding, Logistic Regression Model, Random Forest Model
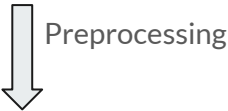Zenthia Song: Decision Tree Classifier Model, K Nearest Neighbor Model, Model Comparison

# Overview

- TASK
  - We aim to train a multi-class classification model to divide new customers into existing segments (A, B, C, D) for a automobile company to predict the potential groups of new customers.
- DATASET
  - https://www.kaggle.com/datasets/kaushiksuresh147/customer-segmentation?datasetId=841888&sortBy=voteCount
- METHOD
  - Logistic Regression, Random Forest, Decision Tree, and K Nearest Neighbours

# Summary Data

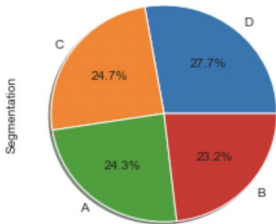8068  training data set

2627 testing data set

⬇ Preprocessing

7308 training data set

2355 testing data set

| | ID | Gender | Ever_Married | Age | Graduated | Profession | Work_Experience | Spending_Score | Family_Size | Var_1 | Segmentation |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 462809 | Male | No | 22 | No | Healthcare | 1.0 | Low | 4.0 | Cat_4 | D |
| 1 | 462643 | Female | Yes | 38 | Yes | Engineer | NaN | Average | 3.0 | Cat_4 | A |
| 2 | 466315 | Female | Yes | 67 | Yes | Engineer | 1.0 | Low | 1.0 | Cat_6 | B |
| 3 | 461735 | Male | Yes | 67 | Yes | Lawyer | 0.0 | High | 2.0 | Cat_6 | B |
| 4 | 462669 | Female | Yes | 40 | Yes | Entertainment | NaN | High | 6.0 | Cat_6 | A |

| | Gender | Ever_Married | Graduated | Spending_Score | Segmentation | Artist | Doctor | Engineer | Entertainment | Executive | ... |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | 0 | 0 | 0 | 0 | 3 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... |
| | 1 | 1 | 1 | 1 | 0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | ... |
| | 1 | 1 | 1 | 0 | 1 | 1.0 | 0.0 | 1.0 | 0.0 | 0.0 | ... |
| | 0 | 1 | 1 | 2 | 1 | NaN | 0.0 | 0.0 | 0.0 | 0.0 | ... |
| | 1 | 1 | 1 | 2 | 0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | ... |

| Cat_1 | Cat_2 | Cat_3 | Cat_4 | Cat_5 | Cat_6 | Cat_7 | Age_normalized |
|---|---|---|---|---|---|---|---|
| 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.058824 |
| 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.294118 |
| 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.720588 |
| 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.720588 |
| 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.323529 |

| Work_Experience_normalized | Family_Size_normalized |
|---|---|
| 0.1 | 0.500000 |
| 0.3 | 0.333333 |
| 0.1 | 0.000000 |
| 0.0 | 0.166667 |
| 0.3 | 0.833333 |

```
D    2127
C    1896
A    1867
B    1779
Name: Segmentation, dtype: int64

<AxesSubplot:ylabel='Segmentation'>
```
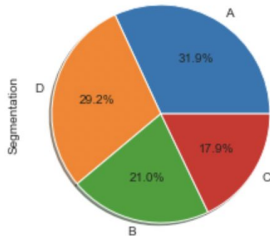
```
A    794
D    726
B    523
C    445
Name: Segmentation, dtype: int64

<AxesSubplot:ylabel='Segmentation'>
```
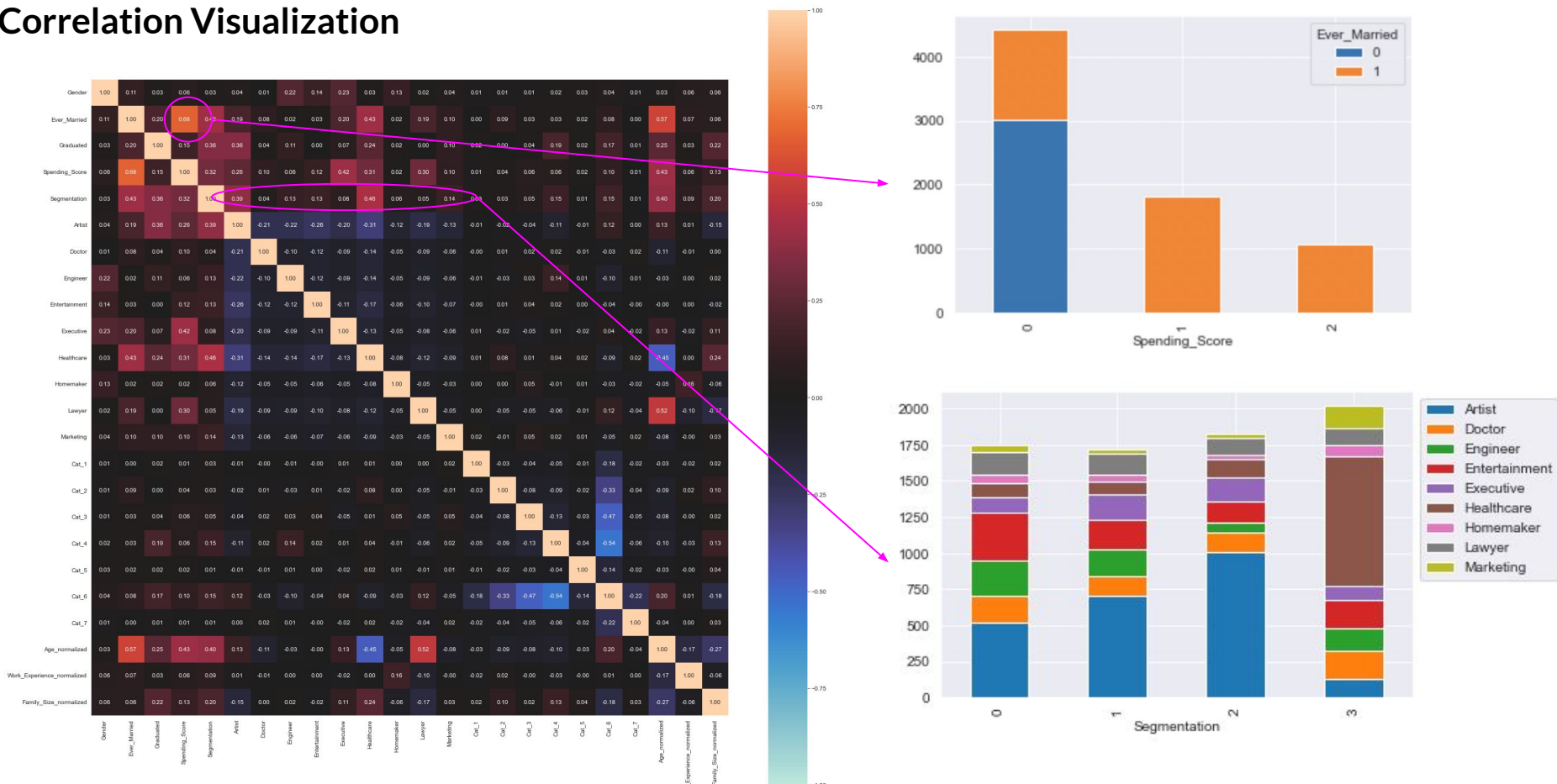
# Correlation Visualization

# Model: Logistic Regression

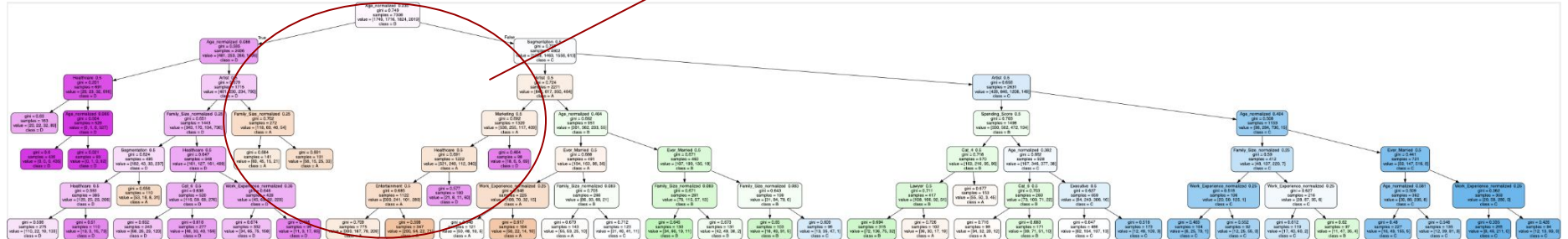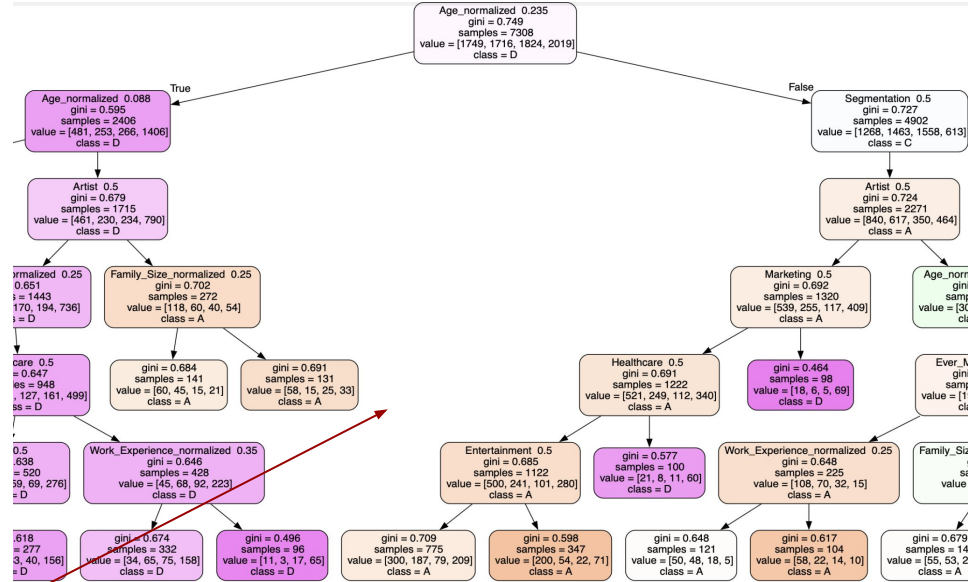| Type of Model | Logistic Regression |
|---|---|
| Parameter | solver='newton-cg'<br>random_state=42 |
| Accuracy | 32.4% |

# Model: Random Forest

| Type of Model | Random Forest |
|---|---|
| Parameter | n_estimators=200<br>criterion="entropy"<br>max_depth=5<br>random_state=42 |
| Accuracy | 33.9% |

# Model: Decision Tree


🏆 WINNER

| Type of Model | Decision Tree |
|---|---|
| Parameter | criterion='gini'<br>random_state=42<br>min_samples_leaf=90<br>max_step=6 |
| Accuracy | 34.5% |

# Model: K Neighbors

| Type of Model | K Neighbors |
|---|---|
| Parameter | radius=1.44 |
| Accuracy | 32.5% |

# Comparing Model Performance

| | Model | Accuracy |
|---|---|---|
| 2 | Decision Tree | 34.5 |
| 1 | Random Forest | 33.9 |
| 3 | K Neighbors | 32.5 |
| 0 | Logistic Regression | 32.4 |

- The model with the highest accuracy is Decision Tree because there are a large number of categories in the feature set.
- Based on these four accuracy scores, we found that 33% is an average score, which is not ideally high.
- Our RF's accuracy score is lower than our DT's, which potentially mean DT outperforms RF. This is possible because aggregated/ensemble models are not universally better than their "single" counterparts, they are better if and only if the single models suffer from instability (a Machine Learning System that produces large differences in generalization patterns when small changes are made to its initial conditions).

# Error Analysis

```
Gender                      1.000000
Ever_Married                1.000000
Graduated                   1.000000
Spending_Score              0.000000
Artist                      0.000000
Doctor                      0.000000
Engineer                    1.000000
Entertainment               0.000000
Executive                   0.000000
Healthcare                  0.000000
Homemaker                   0.000000
Lawyer                      0.000000
Marketing                   0.000000
Cat_1                       0.000000
Cat_2                       0.000000
Cat_3                       0.000000
Cat_4                       0.000000
Cat_5                       0.000000
Cat_6                       1.000000
Cat_7                       0.000000
Age_normalized              0.268657
Work_Experience_normalized  0.000000
Family_Size_normalized      0.000000
Name: 0, dtype: float64
```

Predicted Segmentation: B

Actual Segmentation: A

# Error Analysis

- On the whole, the accuracy of the four models are not high.
- Potential reasons:
  - Lack of data: the amount of data may not capture the complexity of the problem
    - Solution: acquiring more data of the characteristics of old customers
  - Non-representative data: the newly acquired customers may behave inherently different from the old customers
    - Solution: using labeled new customer data to predict unlabeled new customer data
  - Irrelevant Features: We may selected irrelevant features. For Decision Tree model, most of our features have importances smaller than 0.1
    - Solution: Filter out unimportant features. Add other dimension to the data set if possible
  - Algorithm Selection Issue: The choice of four algorithms may also impact the accuracy of the model. Specifically, some more complex algorithms may be better suited this multi-class classification problem.
    - Solution: Deep learning, Naive Bayes