

Phase 4 Report

Megan Joseph

April 27, 2025

1 Summary

The goal of this phase is to study feature importance. The features of my chosen dataset are Age, Gender, AnnualIncome, NumberOfPurchases, ProductCategory, TimeSpentOnWebsite, and LoyaltyProgram. First, I created seven models trained on a single feature and plotting the validation accuracies of each.

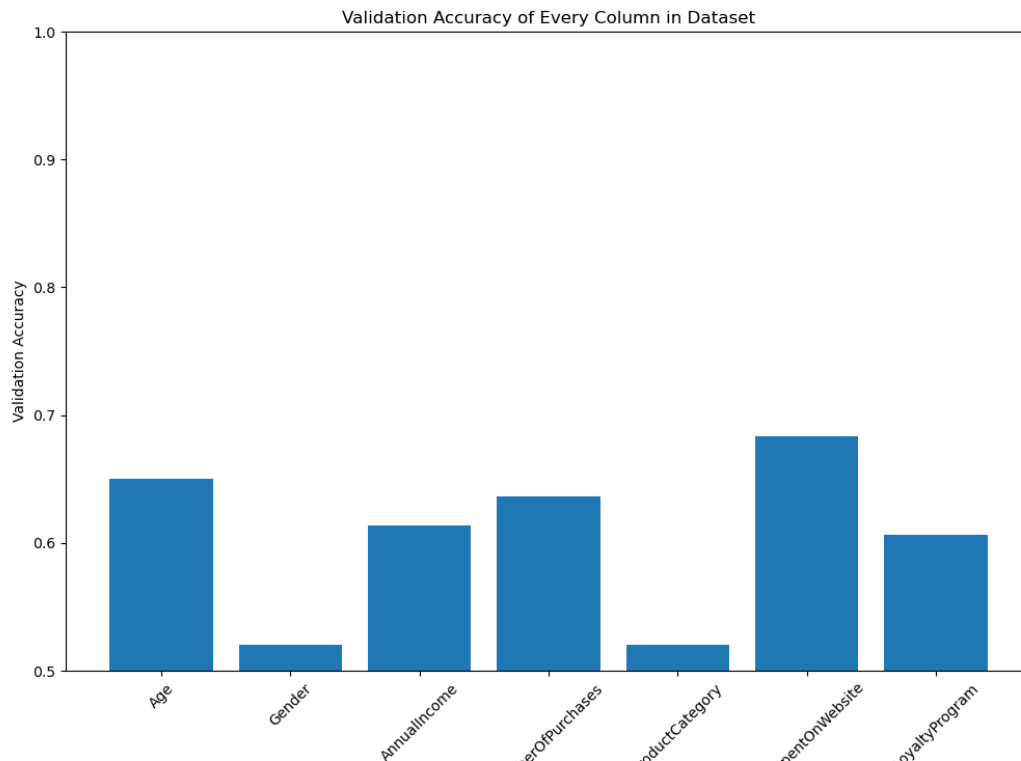


Figure 1: Validation accuracies for models trained on a single feature in the dataset

Next, I started removing the least important features (the features where the respective model had the lowest validation accuracy) from the full model and found the validation accuracy. I did this until the model only had the feature that gave the highest accuracy. In this case, I removed Gender, then ProductCategory, and so on until I was left with only TimeSpentOnWebsite. Then, I plotted the validation accuracies for each. The x-axis labels were too long, so I put them in a legend to make the graph easier to read.

After this, I compared the accuracies of the feature-reduced model with the best validation accuracy to the best model with full features. The feature-reduced model has an accuracy of 0.803 while the full-feature model has an accuracy of 0.817.

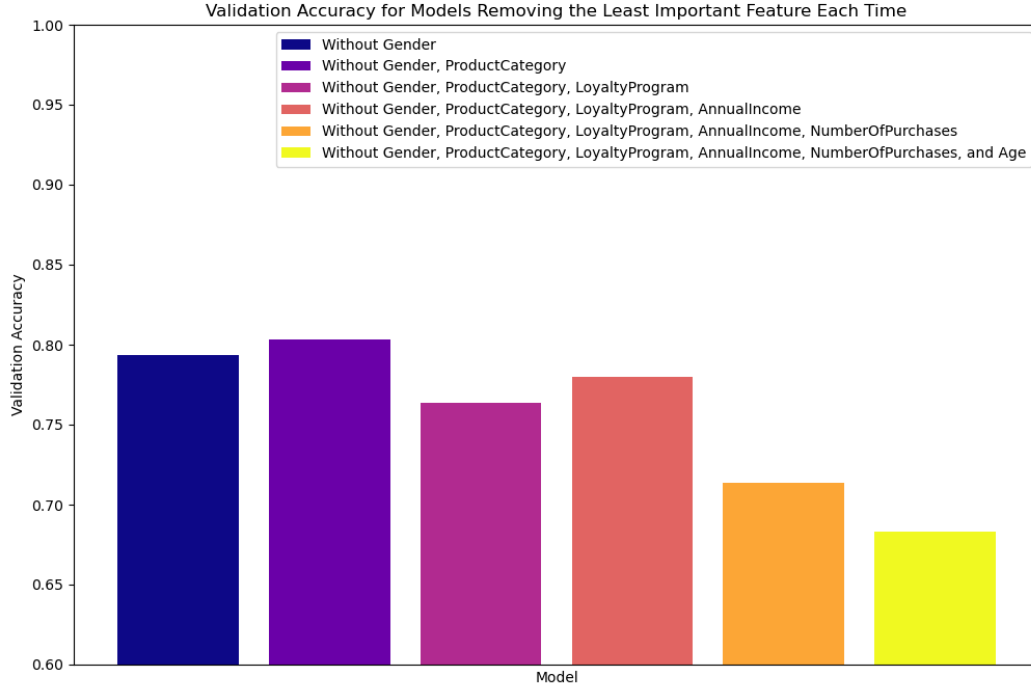


Figure 2: Validation accuracies for models trained on all features and removing the least important feature each time

To investigate feature importance further, I found the Shapley values. I chose to do this instead of LIME because I didn't have too many features, I didn't need a prediction model at the end, and just wanted to investigate the strengths of each feature. I just did a single iteration because I my Jupyter Notebook kept crashing if I did too many, but if I had the compute power, I would average over the entire dataset. The Shapley value I calculated is around 0.79. The column that is most important is AnnualIncome and the least important column is Age. Below is the code used to find the Shapley value:

```

1 def shapley_value(best_model, X_train):
2     instance = X_train.iloc[[1], :]
3     train_pred = best_model.predict(X_train).flatten()
4     base_pred = np.mean(train_pred)
5
6
7     cols = range(len(X_train.columns))
8     subsets = [[]]
9     for col in cols:
10         subsets += [sub + [col] for sub in subsets]
11
12     predictions = {}
13     for sub in subsets:
14         mod_instance = instance.copy()
15         for i in cols:
16             if i not in sub:
17                 mod_instance.iloc[0, i] = 0
18
19         pred = best_model.predict(mod_instance).flatten()
20         predictions[tuple(sub)] = pred
21
22     avg_marg_contr = {}

```

```

23     for col in cols:
24         total_contr = 0
25         num = 0
26         for sub in subsets:
27             if col not in sub:
28                 pred_wo_feature = predictions[tuple(sub)]
29                 sub_w_feature = sub + [col]
30                 pred_w_feature = predictions[tuple(sorted(sub_w_feature))]
31                 marg_contr = pred_w_feature - pred_wo_feature
32                 total_contr += marg_contr
33                 num += 1
34
35         avg_contr = total_contr / num
36         avg_marg_contr[col] = avg_contr
37
38     final_pred = base_pred
39     for col in cols:
40         final_pred += avg_marg_contr[col]
41
42     return final_pred, avg_marg_contr

```

Listing 1: Function to calculate the Shapley value for the model and determine the order of importance based on marginal contribution of each feature.