# Extra Credit Option 1

*Emma Wang (06648801), wangj422@wfu.edu, 12/12/2022*
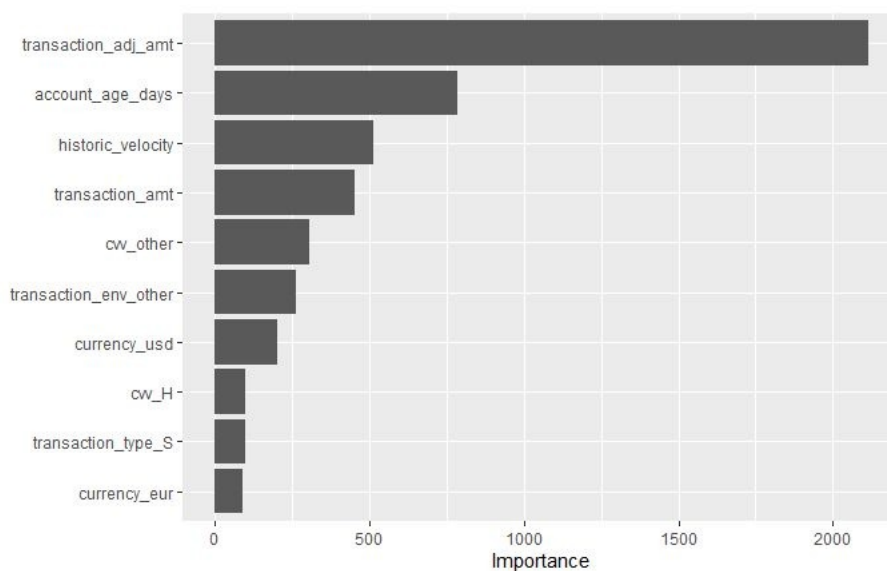
## Model's Performance (From 2nd Challenge in class BAN6053)

| Model | accuracy | Log loss | Roc_auc | precision | recall |
|-------|----------|----------|---------|-----------|--------|
| Logistic Regression - Train | 0.973 | 0.0911 | 0.941 | 0.861 | 0.606 |
| Logistic Regression – Test | 0.973 | 0.0928 | 0.941 | 0.863 | 0.609 |
| Random Forest (10 trees) – Train | 0.988 | 0.041 | 0.999 | 0.987 | 0.795 |
| Random Forest (10 trees) – Test | 0.974 | 0.203 | 0.936 | 0.906 | 0.591 |
| Random Forest (100 trees) - Train | 0.989 | 0.039 | 1.000 | 0.994 | 0.809 |
| Random Forest (100 trees) - Test | 0.976 | 0.101 | 0.946 | 0.930 | 0.611 |

Based on the metrics above, we will choose Random Forest (100 trees) as our best model. And we will explore this model's interpretability further.
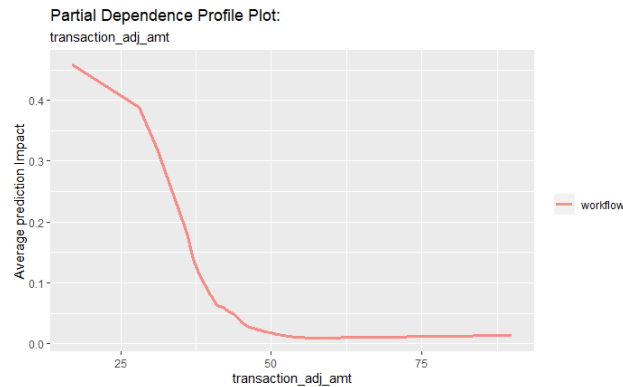
## Global Interpretability

- *Use VIP function to find top 5 most important variables.*
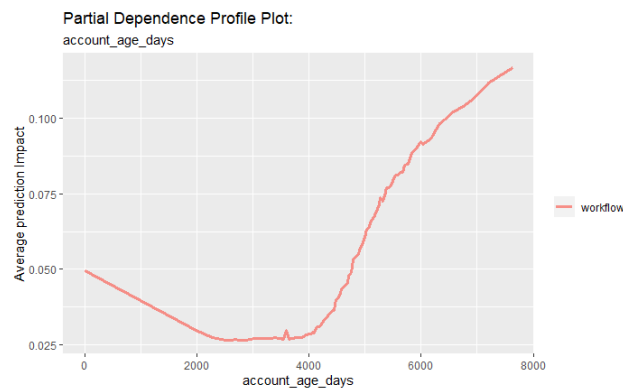


Based on the chart above, the top 5 most important variables are: "transaction_adj_amt", "account_age_days", "historic_velocity", "transaction_amt", and "cvv".

- *Generate Partial Dependence Plot for top 5 variables.*

**Partial Dependence Profile Plot:**
transaction_adj_amt

When "transaction_adj_amt" is below 30, the probability of fraud will be much higher, while if the value is above 50, the probability of fraud will be very low.

**Partial Dependence Profile Plot:**
account_age_days

The relationship between "account_age_days" and the probability of fraud is a "U-shape". The probability of fraud will be higher if the value is too low or too high. However, when the age of the account is above 5000, the probability of fraud will gradually reach the highest.

**Partial Dependence Profile Plot:**
historic_velocity

**Partial Dependence Profile Plot:**
transaction_amt

A higher historical velocity (above 7000), or transaction amount (above 3500) will bring a higher probability of fraud.

Partial Dependence Profile Plot:
cvv

The probability of fraud will be higher if the CVV is coded as "E", "M", "N", "S", "T", or "V".

- *How SHAPLEY VALUEs are calculated.*

Essentially, the Shapley value is the average marginal contribution of a feature considering all possible combinations. From the Shapley value, we can learn how a specific variable influences the prediction result, that is to say, we can learn the direction of the impact from this variable, and, how important this variable is, compared to other variables. We will easily learn the combination of variables' contribution.

## Local Interpretability

- *Find best estimates and worst estimates.*
  - Top 10 true positives: filter event_label = "fraud", arrange .pred_fraud descending, and leave the top 10 records as the most correct predictions.
  - Top 10 false positives: filter event_label = "legit", arrange .pred_fraud descending, leave top 10 records as the most wrong predictions which take negatives as positives.
  - Top 10 false negatives: filter event_label = "fraud", arrange .pred_fraud ascending, we have a tie, 17 records have event_label = "fraud" while the .pred_fraud = 0. We choose to leave 10 records among them, as the most wrong predictions which take positives as negatives.

- *Top 10 TP – Shapley*

From Shapley values, all of the top 10 TP predictions have a "transaction_adj_amt" below 50, and almost all of them have this variable in the top 3 that contributed the most variables. Besides, the majority of them have an "account_age_days" higher than 5500.

- *Top 10 FP - Shapley*

From Shapley value, most of these negatives are assigned as positives because they have a very low transaction adjusted amount. Some of them also have a high account age days.

- *Top 10 FN – Shapley*

According to Shapley value, a lot of these positives are detected as negatives incorrectly because they have a relatively normal account age and transaction-adjusted amount. Even though there are some other variables play parts in it, but these two contributed the most.

- *Top 10 FP – Breakdown*

**Break Down profile**

workflow

| Feature | Contribution |
|---|---|
| intercept | 0.051 |
| transaction_adj_amt = 31 | +0.315 |
| transaction_type = 20 | +0.196 |
| signature_image = 10 | +0.086 |
| historic_velocity = 3817 | +0.004 |
| transaction_amt = 2968 | +0.055 |
| cvv = 4 | +0.009 |
| transaction_env = 4 | -0.013 |
| currency = 1 | +0.005 |
| billing_state = 23 | +0 |
| account_age_days = 4444 | +0.058 |
| + all other factors | +0 |
| prediction | 0.765 |

**Break Down profile**

workflow

| Feature | Contribution |
|---|---|
| intercept | 0.051 |
| transaction_adj_amt = 37 | +0.142 |
| transaction_type = 20 | +0.193 |
| account_age_days = 5395 | +0.235 |
| billing_state = 19 | +0.013 |
| historic_velocity = 4003 | -0.023 |
| transaction_env = 7 | -0.003 |
| cvv = 3 | +0.011 |
| currency = 1 | +0.028 |
| transaction_amt = 2689 | +0.07 |
| signature_image = 22 | +0.032 |
| + all other factors | +0 |
| prediction | 0.751 |

**Break Down profile**

workflow

| Feature | Contribution |
|---|---|
| intercept | 0.051 |
| transaction_adj_amt = 37 | +0.142 |
| signature_image = 14 | +0.214 |
| transaction_type = 8 | +0.1 |
| account_age_days = 5149 | +0.159 |
| cvv = 24 | -0.032 |
| transaction_amt = 2868 | +0.061 |
| transaction_env = 25 | -0.018 |
| currency = 1 | +0.025 |
| billing_state = 37 | +0.014 |
| historic_velocity = 5655 | +0.012 |
| + all other factors | +0 |
| prediction | 0.729 |

**Break Down profile**

workflow

| Feature | Contribution |
|---|---|
| intercept | 0.051 |
| transaction_adj_amt = 36 | +0.184 |
| account_age_days = 5607 | +0.202 |
| transaction_type = 22 | +0.146 |
| cvv = 24 | -0.03 |
| signature_image = 20 | +0.038 |
| transaction_amt = 2994 | +0.069 |
| transaction_env = 25 | +0.014 |
| currency = 1 | +0.054 |
| billing_state = 14 | +0.002 |
| historic_velocity = 5546 | -0.01 |
| + all other factors | +0 |
| prediction | 0.721 |

**Break Down profile**

workflow

| Feature | Contribution |
|---|---|
| intercept | 0.051 |
| transaction_adj_amt = 34 | +0.231 |
| transaction_type = 23 | +0.168 |
| account_age_days = 5293 | +0.211 |
| cvv = 6 | +0 |
| transaction_env = 4 | -0.023 |
| transaction_amt = 2765 | +0.063 |
| currency = 1 | +0.036 |
| historic_velocity = 4894 | -0.015 |
| signature_image = 8 | +0.029 |
| billing_state = 6 | -0.047 |
| + all other factors | +0 |
| prediction | 0.705 |

**Break Down profile**

workflow

| Feature | Contribution |
|---|---|
| intercept | 0.051 |
| transaction_adj_amt = 28 | +0.37 |
| cvv = 22 | +0.064 |
| transaction_env = 8 | +0.043 |
| account_age_days = 4961 | +0.104 |
| transaction_type = 6 | +0.066 |
| billing_state = 33 | -0.002 |
| historic_velocity = 4667 | +0.038 |
| transaction_amt = 2538 | +0.044 |
| signature_image = 7 | -0.01 |
| currency = 4 | -0.064 |
| + all other factors | +0 |
| prediction | 0.703 |

- *Top 10 FP – Breakdown*

**Break Down profile**

workflow

| Feature | Contribution |
|---|---|
| intercept | 0.051 |
| transaction_type = 22 | +0.025 |
| account_age_days = 4064 | -0.027 |
| transaction_amt = 1736 | -0.016 |
| transaction_adj_amt = 48 | -0.009 |
| transaction_env = 24 | -0.004 |
| cvv = 4 | -0.004 |
| historic_velocity = 4092 | -0.006 |
| currency = 1 | -0.002 |
| signature_image = 25 | -0.006 |
| billing_state = 40 | -0.002 |
| + all other factors | +0 |
| prediction | 0 |

**Break Down profile**

workflow

| Feature | Contribution |
|---|---|
| intercept | 0.051 |
| account_age_days = 5565 | +0.049 |
| transaction_adj_amt = 67 | -0.06 |
| cvv = 25 | -0.009 |
| transaction_amt = 2932 | +0.013 |
| transaction_env = 3 | -0.005 |
| transaction_type = 4 | -0.014 |
| billing_state = 33 | -0.011 |
| signature_image = 22 | -0.002 |
| historic_velocity = 5687 | -0.004 |
| currency = 4 | -0.009 |
| + all other factors | +0 |
| prediction | 0 |

**Break Down profile**
workflow

| | |
|---|---|
| intercept | 0.051 |
| account_age_days = 5912 | +0.059 |
| transaction_adj_amt = 54 | -0.05 |
| signature_image = 23 | +0.014 |
| cvv = 3 | -0.011 |
| transaction_env = 9 | -0.004 |
| transaction_type = 3 | -0.023 |
| transaction_amt = 2354 | +0.015 |
| historic_velocity = 5114 | +0.003 |
| billing_state = 14 | -0.009 |
| currency = 4 | -0.046 |
| + all other factors | +0 |
| prediction | 0 |

**Break Down profile**
workflow

| | |
|---|---|
| intercept | 0.051 |
| billing_state = 48 | +0.043 |
| account_age_days = 4156 | -0.043 |
| transaction_amt = 1788 | -0.02 |
| transaction_type = 18 | +0.005 |
| transaction_adj_amt = 48 | -0.012 |
| cvv = 3 | -0.005 |
| transaction_env = 9 | -0.002 |
| currency = 1 | +0 |
| signature_image = 25 | -0.014 |
| historic_velocity = 4548 | -0.004 |
| + all other factors | +0 |
| prediction | 0 |

**Break Down profile**
workflow

| | |
|---|---|
| intercept | 0.051 |
| transaction_adj_amt = 59 | -0.028 |
| transaction_type = 6 | +0.005 |
| transaction_amt = 1493 | -0.007 |
| transaction_env = 3 | -0.006 |
| cvv = 3 | -0.003 |
| currency = 1 | +0 |
| account_age_days = 4568 | -0.001 |
| historic_velocity = 4423 | -0.003 |
| billing_state = 44 | -0.001 |
| signature_image = 27 | -0.007 |
| + all other factors | +0 |
| prediction | 0 |

**Break Down profile**
workflow

| | |
|---|---|
| intercept | 0.051 |
| transaction_type = 23 | +0.098 |
| transaction_adj_amt = 60 | -0.053 |
| billing_state = 10 | +0.006 |
| transaction_amt = 1955 | -0.043 |
| cvv = 4 | -0.002 |
| transaction_env = 4 | -0.008 |
| currency = 1 | +0.004 |
| account_age_days = 4398 | -0.018 |
| historic_velocity = 4406 | -0.011 |
| signature_image = 4 | -0.024 |
| + all other factors | +0 |
| prediction | 0 |

**Break Down profile**
workflow

| | |
|---|---|
| intercept | 0.051 |
| account_age_days = 5150 | +0.031 |
| transaction_adj_amt = 68 | -0.05 |
| signature_image = 20 | +0.017 |
| transaction_amt = 3086 | +0.004 |
| cvv = 4 | -0.009 |
| transaction_env = 4 | -0.003 |
| transaction_type = 4 | -0.022 |
| currency = 1 | -0.009 |
| historic_velocity = 5732 | +0.004 |
| billing_state = 14 | -0.014 |
| + all other factors | +0 |
| prediction | 0 |

**Break Down profile**
workflow

| | |
|---|---|
| intercept | 0.051 |
| transaction_adj_amt = 62 | -0.028 |
| transaction_type = 22 | +0.013 |
| account_age_days = 3178 | -0.011 |
| historic_velocity = 6248 | +0.008 |
| transaction_amt = 1484 | -0.01 |
| cvv = 3 | -0.003 |
| currency = 1 | -0.001 |
| billing_state = 37 | -0.002 |
| transaction_env = 27 | -0.006 |
| signature_image = 27 | -0.01 |
| + all other factors | +0 |
| prediction | 0 |

**Break Down profile**
workflow

| | |
|---|---|
| intercept | 0.051 |
| cvv = 22 | +0.03 |
| transaction_adj_amt = 54 | -0.03 |
| account_age_days = 3908 | -0.014 |
| transaction_amt = 1947 | -0.011 |
| billing_state = 5 | +0 |
| currency = 1 | -0.004 |
| transaction_type = 9 | -0.001 |
| historic_velocity = 4631 | -0.005 |
| signature_image = 9 | -0.006 |
| transaction_env = 16 | -0.011 |
| + all other factors | +0 |
| prediction | 0 |

**Break Down profile**
workflow

| | |
|---|---|
| intercept | 0.051 |
| currency = 2 | +0.145 |
| transaction_adj_amt = 65 | -0.056 |
| account_age_days = 3606 | -0.093 |
| transaction_amt = 2035 | -0.02 |
| cvv = 25 | +0.004 |
| historic_velocity = 4160 | -0.005 |
| transaction_env = 7 | -0.003 |
| transaction_type = 4 | -0.002 |
| billing_state = 23 | -0.003 |
| signature_image = 4 | -0.019 |
| + all other factors | +0 |
| prediction | 0 |

- *Top 10 FP and FN, Shapley and Breakdown*
  - The results are not entirely the same but there are some similarities. For example, both transaction-adjusted amount and account age play an important part in both Shapley and Breakdown's FPs and FNs.

- However, transaction types seem to play a more important role in false negatives from Breakdown than from Shapley.
- Challenge: different methods will yield different results, they will not be exactly the same, except for some really important variables. It is harder to detect the impact of some less impacting variables.