

Model #101: Credit Card Default Model

Performance Validation Guide

Mimi Trinh

Section 1: Production Model

Using the results of the first two models of random forest and gradient boosting, we identify a pool of interesting predictors to use in the logistic regression model. Specifically, we remove the four insignificant demographic variables MARRIAGE, AGE, SEX, EDUCATION and leave the seven payment-related predictors remain in the predictor pool to develop a logistic regression model. Then among these seven variables, we use the stepwise automatic variable selection method to arrive at the optimal logistic regression model.

Figure 1 – Logistic Regression Model Summary Result

```
Call:
glm(formula = DEFAULT ~ LIMIT_BAL + AVG_BILL_AMT + AVG_PAY_AMT +
     AVG_UTIL + MAX_BILL_AMT + MAX_PAY_AMT + MAX_DLQ, family = binomial(),
     data = model_train)

Deviance Residuals:
    Min       1Q   Median       3Q      Max
-1.4782  -0.7867  -0.6500  -0.2307   4.9009

Coefficients:
              Estimate Std. Error z value Pr(>|z|)
(Intercept) -8.342e-01  5.612e-02 -14.865  < 2e-16 ***
LIMIT_BAL   -1.960e-06  2.622e-07  -7.473  7.82e-14 ***
AVG_BILL_AMT  9.209e-06  1.826e-06   5.043  4.59e-07 ***
AVG_PAY_AMT  -1.787e-04  1.638e-05 -10.905  < 2e-16 ***
AVG_UTIL      2.838e-01  8.895e-02   3.190   0.00142 **
MAX_BILL_AMT -8.600e-06  1.719e-06  -5.003  5.65e-07 ***
MAX_PAY_AMT   2.970e-05  3.125e-06   9.504  < 2e-16 ***
MAX_DLQ      -4.876e-06  2.102e-06  -2.319   0.02037 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 16205  on 15179  degrees of freedom
Residual deviance: 15464  on 15172  degrees of freedom
AIC: 15480

Number of Fisher Scoring iterations: 6
```

The result above shows that all seven predictors are statistically significant at 95% confidence level, each with p-value less than 0.05 alpha. Thus, the stepwise automatic variable selection algorithm indicates that all variables in the model are significant. Among these seven predictors, LIMIT_BAL, AVG_PAY_AMT, MAX_BILL_AMT, MAX_DLQ have negative coefficients, which mean that they have a negative correlation with the dependent variable. In other words, the lower the limit balance and the lower the average payment amount and the lower of the maximum bill and the lower the maximum delinquency value, the higher the probability of default on payment. The other three predictors AVG_BILL_AMT, AVG_UTIL, MAX_PAY_AMT have positive coefficient, which mean that these variables have a positive correlation with the response variable. In other words, the higher the average billing amount and the higher the utilization rate and the higher the maximum payment amount, the higher the chance of default on payment.

Section 2: Model Development Performance

Figure 2 – Model Performance on Train Dataset

	No Default Actual	Default Actual	Sum
No Default Predicted	7443	1332	8775
Default Predicted	4314	2091	6405
Sum	11757	3423	15180

The classification table above is used to calculate the following performance metrics for the train dataset.

- $TPR = 2091 / 3423 = 61.09\%$
- $FPR = 4314 / 11,757 = 36.69\%$
- $Accuracy = (7443 + 2091) / 15,180 = 62.81\%$

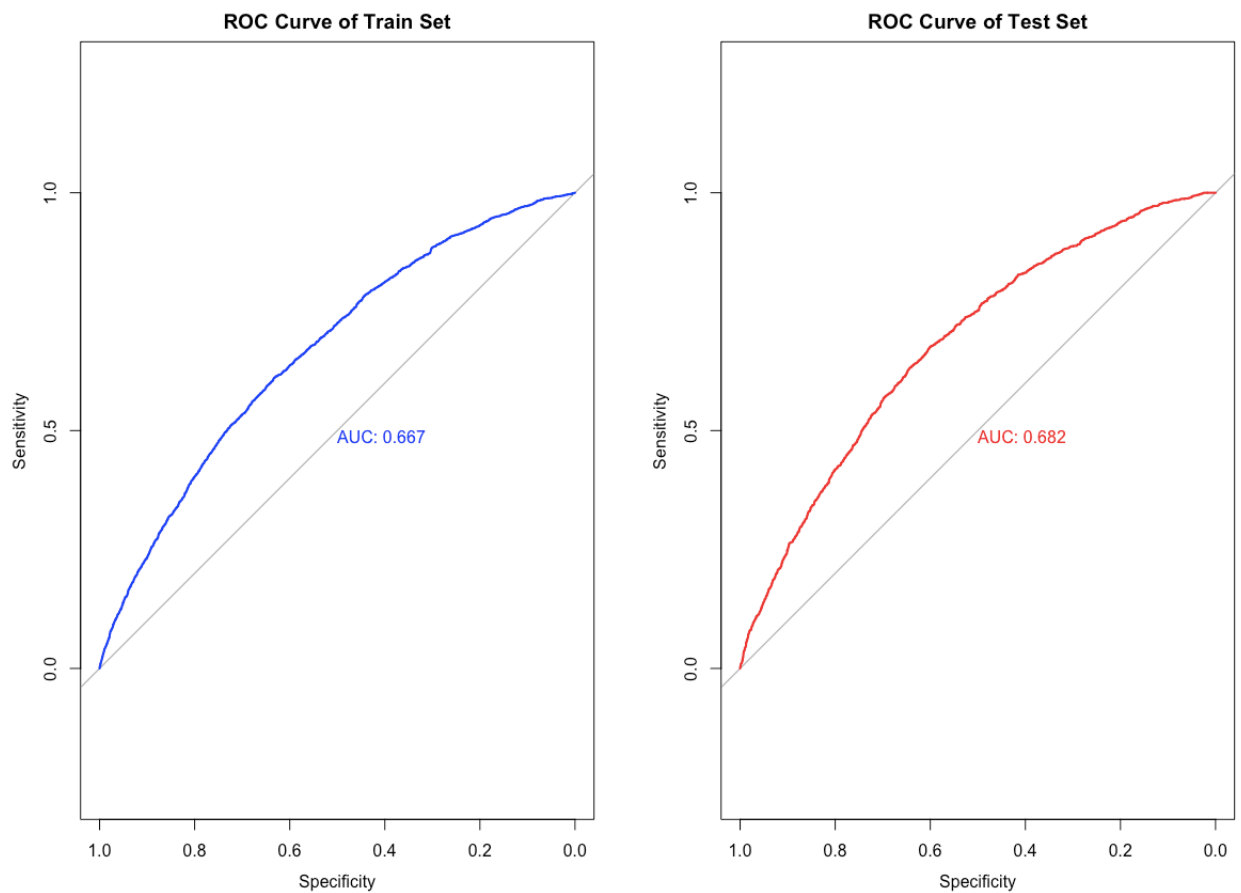
Figure 3 – Model Performance on Test Dataset

	No Default	Actual Default	Actual	Sum
No Default Predicted		3470	506	3976
Default Predicted		2296	1051	3347
Sum		5766	1557	7323

The classification matrix above is used to calculate the following performance metrics for the test dataset.

- $TPR = 1051 / 1557 = 67.5\%$
- $FPR = 2296 / 5766 = 39.82\%$
- $Accuracy = (3470 + 1051) / 7323 = 61.74\%$

Figure 4 – ROC Curve and AUC on both Train and Test Dataset



The two ROC curves and AUC above for train and test sets are very similar to one another. Thus, we can conclude that there's no overfitting issue in the logistic regression model. In general, the

higher the AUC, the better the model. The AUC value ranges from 0.5 to indicate no diagnostic ability to 1.0 to indicate perfect diagnostic ability. In this scenario, the AUC on the train set is 0.667 whereas the AUC on the test is 0.682, which mean that the logistic regression model doesn't successfully capture the response variable. Thus, though this is a good model to start, there is still much work needed to be done to improve the model. Moreover, since the AUC numbers for both train and test sets are similar to one another, we can conclude that the model doesn't have an overfitting issue.

Figure 5 – Lift Chart for Train Dataset

Train Set		Goods		Bads				
Decile	Obs	Target (Y=1)	NonTarget (Y=0)	Target Density	NonTarget Density	Target CDF	NonTarget CDF	KS Stat
1	759	358	401	10.5%	3.4%	10.5%	3.4%	7.0%
2	759	294	465	8.6%	4.0%	19.0%	7.4%	11.7%
3	759	260	499	7.6%	4.2%	26.6%	11.6%	15.0%
4	759	238	521	7.0%	4.4%	33.6%	16.0%	17.6%
5	759	247	512	7.2%	4.4%	40.8%	20.4%	20.4%
6	759	227	532	6.6%	4.5%	47.4%	24.9%	22.5%
7	759	194	565	5.7%	4.8%	53.1%	29.7%	23.4%
8	759	184	575	5.4%	4.9%	58.5%	34.6%	23.9%
9	759	165	594	4.8%	5.1%	63.3%	39.7%	23.6%
10	759	157	602	4.6%	5.1%	67.9%	44.8%	23.1%
11	759	155	604	4.5%	5.1%	72.4%	49.9%	22.5%
12	759	167	592	4.9%	5.0%	77.3%	55.0%	22.3%
13	759	141	618	4.1%	5.3%	81.4%	60.2%	21.2%
14	759	125	634	3.7%	5.4%	85.1%	65.6%	19.5%
15	759	136	623	4.0%	5.3%	89.0%	70.9%	18.1%
16	759	89	670	2.6%	5.7%	91.6%	76.6%	15.0%
17	759	97	662	2.8%	5.6%	94.5%	82.2%	12.2%
18	759	76	683	2.2%	5.8%	96.7%	88.0%	8.6%
19	759	72	687	2.1%	5.8%	98.8%	93.9%	4.9%
20	759	41	718	1.2%	6.1%	100.0%	100.0%	0.0%
Totals	15,180	3423	11,757	67.9%	100.0%			

Figure 6 – Lift Chart for Test Dataset

Test Set		Goods		Bads				
Decile	Obs	Target (Y=1)	NonTarget (Y=0)	Target Density	NonTarget Density	Target CDF	NonTarget CDF	KS Stat
1	367	170	197	10.9%	3.4%	10.9%	3.4%	7.5%
2	366	134	232	8.6%	4.0%	19.5%	7.4%	12.1%
3	366	127	239	8.2%	4.1%	27.7%	11.6%	16.1%
4	366	119	247	7.6%	4.3%	35.3%	15.9%	19.5%
5	366	106	260	6.8%	4.5%	42.1%	20.4%	21.8%
6	366	101	265	6.5%	4.6%	48.6%	25.0%	23.6%
7	366	104	262	6.7%	4.5%	55.3%	29.5%	25.8%
8	366	91	275	5.8%	4.8%	61.1%	34.3%	26.9%
9	366	82	284	5.3%	4.9%	66.4%	39.2%	27.2%
10	366	68	298	4.4%	5.2%	70.8%	44.4%	26.4%
11	367	66	301	4.2%	5.2%	75.0%	49.6%	25.4%
12	366	67	299	4.3%	5.2%	79.3%	54.8%	24.5%
13	366	60	306	3.9%	5.3%	83.2%	60.1%	23.1%
14	366	54	312	3.5%	5.4%	86.6%	65.5%	21.1%
15	366	39	327	2.5%	5.7%	89.1%	71.2%	18.0%
16	366	49	317	3.1%	5.5%	92.3%	76.7%	15.6%
17	366	40	326	2.6%	5.7%	94.9%	82.3%	12.5%
18	366	39	327	2.5%	5.7%	97.4%	88.0%	9.4%
19	366	22	344	1.4%	6.0%	98.8%	94.0%	4.8%
20	367	19	348	1.2%	6.0%	100.0%	100.0%	0.0%
Totals	7,323	1557	5,766	70.8%	100.0%			

Using the semi-deciles or half-deciles with 20 groups, the two lift charts above are provided, one for train set and one for test set. From the lift charts above, the Kolmogorov-Smirnov (KS) statistics are calculated: 23.9 for train set and 27.2 for test set. Similar to the AUC numbers, the two KS statistics for train and test sets are very close to each other, which means that the predictive model has no overfitting issue.

Section 3: Performance Monitoring Plan

According to Lyn C. Thomas in the book “Consumer Credit Models – Pricing, Profit, and Portfolios,” the rule of thumb is that KS statistics of 0.4 (or 40 as the metric used in this project)

suggest good discrimination (p112). In this project, the KS statistics of 23.9 and 27.2 on train and test sets are quite low below the threshold. To be more specific, this performance monitoring plan uses the table below outlining the metric threshold for the KS statistics that determine each RAG (red, amber, green) status. The KS statistics in this study uses an absolute change instead of a percentage change to define the performance status.

Figure 7 – KS Statistics Threshold for RAG Status

Performance Status	KS Statistics Threshold
Red	0 – 30
Amber	31 – 60
Green	61 – 100

In the RAG status, red category means that the model needs redevelopment. Amber means that the model needs to be re-validated in three months. Green means that the model is performing as expected. Therefore, using the table above, the KS statistics for the logistic regression predictive model falls in the red category, which means that the logistic regression model needs to be redeveloped before being used in operations.

Section 4: Performance Monitoring Results

Figure 8 – Lift Chart for Validate Dataset

Validate Set		Goods		Bads				
Decile	Obs	Target (Y=1)	NonTarget (Y=0)	Target Density	NonTarget Density	Target CDF	NonTarget CDF	KS Stat
1	375	163	212	9.8%	3.6%	9.8%	3.6%	6.2%
2	375	149	226	9.0%	3.9%	18.8%	7.5%	11.3%
3	375	142	233	8.6%	4.0%	27.4%	11.5%	15.9%
4	375	143	232	8.6%	4.0%	36.1%	15.5%	20.6%
5	374	113	261	6.8%	4.5%	42.9%	19.9%	22.9%
6	375	105	270	6.3%	4.6%	49.2%	24.6%	24.7%
7	375	84	291	5.1%	5.0%	54.3%	29.5%	24.8%
8	375	85	290	5.1%	5.0%	59.4%	34.5%	24.9%
9	375	82	293	5.0%	5.0%	64.4%	39.5%	24.9%
10	374	76	298	4.6%	5.1%	69.0%	44.6%	24.3%
11	375	75	300	4.5%	5.1%	73.5%	49.8%	23.7%
12	375	70	305	4.2%	5.2%	77.7%	55.0%	22.7%
13	375	59	316	3.6%	5.4%	81.3%	60.4%	20.9%
14	375	74	301	4.5%	5.2%	85.7%	65.5%	20.2%
15	374	59	315	3.6%	5.4%	89.3%	70.9%	18.4%
16	375	44	331	2.7%	5.7%	92.0%	76.6%	15.4%
17	375	38	337	2.3%	5.8%	94.3%	82.4%	11.9%
18	375	42	333	2.5%	5.7%	96.8%	88.1%	8.7%
19	375	31	344	1.9%	5.9%	98.7%	94.0%	4.7%
20	375	22	353	1.3%	6.0%	100.0%	100.0%	0.0%
Totals	7,497	1656	5,841	69.0%	100.0%			

Using the lift chart above with semi-deciles or half-deciles with 20 groups, the KS statistics for validate dataset is 24.9, which is very close to the KS statistics for train and test datasets. Thus, the model doesn't have overfitting issue. All three KS statistics for train, test, and validate datasets fall within the red category in the RAG status since they are all below the threshold of 30. Thus, the model needs to be redeveloped. As a result, below are the recommendations to model governance, and specifically what future researchers can do to improve the model, which is also reflected in the model development guide.

1. Consider more relevant predictors. From the results of this project, demographics variables such as sex, marriage, age don't have significant impact on the response

variable. Thus, data scientists should consider additional payment-related predictors such as FICO score, payment method, etc.

2. Try different modeling techniques such as neural networks. Though the logistic regression model performs moderately, it still needs to be redeveloped.
3. Consider options such as zero-based Poisson and zero-based negative binomial approaches along with the logistic regression model to address the imbalance of default vs. no default in response variable.