

GERMAN_LOGISTIC_MODEL

The functional form of the Logistic Regression Model is given by

$$P(Y=1) = \frac{e^z}{1 + e^z}$$

$$\text{Where } Z = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \dots + \beta_{11} X_{11} + \beta_{12} X_{12} + \beta_{13} X_{13}$$

Where

$Y = 1$: Credit Classification as Bad

$Y = 0$: Credit Classification as Good

X_1 = CHK_ACCT

X_2 = Duration

X_3 = Credit History

X_4 = Credit Amount

X_5 = Balance in Savings A/C

X_6 = Employment

X_7 = Install_rate

X_8 = Marital status

X_9 = Present Resident

X_{10} = Age

X_{11} = Other installment

X_{12} = Num_Credits

X_{13} = Job

Model 1 Summary

$$P(Y=1) = \frac{e^Z}{1 + e^Z}$$

Where

$$Z = -0.0754 + 0.0104 * \text{Duration} + 0.0002 * \text{CreditAmount} + 0.3399 * \text{Install_rate} + \dots + 0.1159 * \text{Job_skilled} - 0.3676 * \text{Job_Unemployed} - 0.0018 * \text{Job_Unskilled}$$

The model summary suggest that as per walds test, only 7 features are statistically significant at a significance value of $\alpha=0.05$. The significant variables are

Credit Amount

Install_rate

Age

CHK_ACCT_no-account

CHK_ACCT_over-200DM

Credit History_ critical

Credit History_ delay

Model 2 Summary

The new Model with significant variables are

$$P(Y=1) = \ln [e^Z / (1 + e^Z)]$$

$$Z = -0.4385 + 0.0002 * \text{CreditAmount} + 0.2987 * \text{Install_rate} - 0.0267 * \text{Age} - 1.7346 * \text{CHK_ACCT_no-account-} \\ 1.1121 * \text{CHK_ACCT_over-200DM} - 0.6778 * \text{Credit History_critical} - 0.4645 * \text{Credit History_delay}$$

Some observations from the Model output are

- The probability of being bad customer increases as credit amount and install rate increases.
- The probability of being bad customer decreases as Age decreases.

Confusion matrix

	Bad credit	Good credit
Bad credit	24 (True Negative)	37 (False Positive)
Good credit	30 (False Negative)	149 (True Positive)

In the confusion matrix the columns represent the predicted class while the rows represent the actual class.

Accuracies

$$\text{Sensitivity} = \text{TP} / (\text{TP} + \text{FN}) = 149 / (149 + 30) = 0.832$$

$$\text{Specificity} = \text{TN} / (\text{TN} + \text{FP}) = 24 / (24 + 37) = 0.393$$

$$\text{Precision} = \text{TP} / (\text{TP} + \text{FP}) = 149 / (149 + 37) = 0.801$$

ROC Curve Summary

The Receiver Operating Curve can be used to understand the overall performance of Logistic Regression Model. It is a plot between sensitivity(True Positive Rate) on the vertical axis and 1-Specificity (false positive rate) on the horizontal axis.

The diagonal line in the figure represents the case of not using a model. The sensitivity and Specificity are likely to change when the cut-off probability is changed. The line above the diagonal line in the figure captures how sensitivity and 1-Specificity changes when the cut-off probability is changed. Model with higher AUC is preferred. The AUC for the model obtained is 0.72 validates the model.