

Proposal for a Bankruptcy Prediction Model with Modified Definition of Bankruptcy for Slovak Companies

Peter Adamko¹

Tomas Kliestik²

Abstract

The purpose of this paper was to develop a bankruptcy prediction model for Slovak companies. The model was constructed on a modified definition of bankruptcy. One of the most important requirements when creating any prediction model is to obtain a suitable sample of data. A sample of healthy companies and a sample of bankrupt companies are usually used when bankruptcy models are created. Since 2014, in Slovakia there is the possibility to use the Register of Financial Statements to retrieve data of all Slovak companies. Unfortunately, it is not possible to clearly distinguish between healthy companies and bankrupt companies from these data. Therefore, we used the Act. 431/2002 Coll. on Accounting, as amended, to determine the status of the company. It defines company in debt: "... the value of its liabilities exceeds the value of its assets.". For the purposes of our modelling and testing the difference between assets and liabilities expresses state of company, where a negative value classifies company as being in debt (bankruptcy). To test the performance of the model we used a standard metric (AUC, Sensitivity, Confusion Matrix, RMSE, ...). We found that under certain conditions the model works very well.

Keywords: default, financial distress, default prediction model

JEL Classification:G17, C52, C53

Introduction

A suitable sample of data is one of the most important and the most demanding requirements when any prediction model is created. In case of model of bankruptcy prediction model we need a sample financial statements. Before 2014, it was not easy to get to relevant data about the Slovak enterprises. However, from 2014, there is possible to obtain data from the register of the financial statements (register, <http://www.registeruz.sk>). It was created with the goal of improving and simplifying the business environment and reducing the administrative burden of the business. The register also improves handiness and quality of the information about the accounting objects. Register allows to view and download available financial documents published in the register. From the viewpoint of mass data collection, register's interface is a big advantage for automated downloading of data from the public section. Requirement of building application to communicate with the interface is only small disadvantage or limitation.

¹ University of Zilina, Slovakia

² University of Zilina, Slovakia

Massive amount of available financial statements in the registry is big advantage. But this benefit of mass processing is balanced by a lack of clear determination of whether the company is in bankruptcy. So, we used the Act. 431/2002 Coll. on Accounting, as amended, to define the status of the company. It describes company in debt: "... *the value of its liabilities surpasses the value of its assets.*". For the purposes of our modelling and testing the difference between assets and liabilities expresses state of company, where a negative value classifies company as being in debt (bankruptcy) (Kocisova & Misankova, 2014).

1. Metrics

To evaluate the quality or effectiveness of a model we can use different metrics. There are many ways to determine the performance of the model, unfortunately, each has its drawbacks, which are necessary to be considered.

Confusion matrix (error matrix) contains information about actual and predicted classifications done by a model. It is an N x N matrix, where N is the number of classes being predicted. For our case we have a 2x2 matrix with the following cells: TP –number of true positive cases, TN –true negative, FP - false positive, FN – false negative. From the matrix can be calculated the various metrics:

Accuracy($\Sigma TP + \Sigma TN$)/ Σ Total population. Accuracy is the proportion of the total number of predictions that were correct.

Sensitivity (Recall) is the true positive rate. It is proportion of true outcomes which are correctly predicted.

Specificity is false positive rate. It is proportion of false outcomes which are correctly predicted.

Positive likelihood ratio is ratio between Sensitivity/(1-Specificity).

Negative likelihood ratio is ratio between (1-Sensitivity)/Specificity.

Positive predictive value (Precision), $TP/(TP + FP)$. It is probability that the outcome is true when it is predicted.

Negative predictive value, $TN/(TN + FN)$. It is probability that the outcome is false when it is not predicted.

$$F_{\beta} = (1 + \beta^2) \frac{\text{Precision} \cdot \text{Sensitivity}}{\beta^2 \cdot \text{Precision} + \text{Sensitivity}}$$

F₁ is the harmonic mean of precision and sensitivity.

F₂ considers recall higher than precision - penalizes high numbers of FN

F_{0.5} puts more emphasis on precision than recall - penalizes high number of FP.

Absolute MCC (Matthews correlation coefficient) $\sqrt{\frac{\chi^2}{N}}$. MCC is related to the chi-square statistic for a 2×2 contingency table. Balanced measure which can be used even if the classes are of very different sizes. It returns a value between 0 and 1, 0 being totally dissimilar, 1 being identical. The Matthews correlation coefficient is often regarded as being one of the best metrics.

Minimum per class accuracy is the worst accuracy, maximize this is doing the best to raise the lowest accuracy.

Kappa $\kappa = \frac{p_0 - p_e}{1 - p_e}$ is a level of agreement between observations and predictions where $p_0 = TP + TN$ and $p_e = \frac{(TP + FN)(TP + FP) + (FN + TN)(FP + TN)}{N}$. Kappa value less than 0.2 is considered as poor agreement, values above 0.8 are a sign of high correlation.

Other frequently used criteria are maximize(Sensitivity + Specificity), maximize(Percent Correctly Classified), minimize(distance ROC curve from the upper left point), (Predicted prevalence = Observed prevalence) and (Specificity = Sensitivity).

ROC (Receiver Operating Characteristic) is curve with points [x, y], where x = 100-Specificity and y = Sensitivity for different cut-off (threshold) points. The closer the ROC curve is to the upper left corner, the higher the overall accuracy of the model. The closer the curve comes to the diagonal line, the worse the model. An ROC curve reveals trade off between sensitivity and specificity - increasing of sensitivity imply decreasing of specificity and vice versa). The biggest benefit of using ROC curve is that it is independent of the change in proportion of outcomes.

AUC (Area Under Curve) is one of the most common metrics. It is area under curve ROC. An area of 1 characterizes a perfect model. AUC above 0.9 is excellent and AUC below 0.6 means very bad performance of a model. Like other methods, high AUC does not automatically guarantee a quality of the model. For example, there are situations where the sensitivity is in the range of only a few hundreds, and the specificity is over 90%, and the AUC is still above 0.8.

Gini coefficient is derived from the AUC. $Gini = 2 * AUC - 1$. Gini above 60% is a considered as good model.

RMSE (Root Mean Squared Error) follows an assumption that errors are unbiased and follow a normal distribution. $RMSE = \sqrt{\frac{\sum_{i=1}^N (Predicted_i - Actual_i)^2}{N}}$. RMSE is highly affected by outlier values (Cisco, S. & Kliestik, T., 2013).

2. Model

In our model we used four predictors: Working capital / Total assets, Retained earnings / Total assets, EBIT / Total assets, Equity / Total liabilities.

We had financial statements from 2014 and 2015. Data from 2014 were divided in the ratio of 70:15:15 -for training (31,443), validation (6,738) and testing (6,738). Partition was done in order to preserve the same proportion of companies in each group (18.6% in bankruptcy). All the data of 2015 (45,368) were used to test the model, 15.9% of them was classified as in bankruptcy. Comparison of training data and testing (2015) data is summarized in **Table 1**.

Table 1: Summary of data used in training and testing model

	2014 train				2015			
	X1	X2	X3	X4	X1	X2	X3	X4
SD	0.333	0.179	0.1	0.944	0.328	0.188	0.117	1.156
VAR	0.111	0.032	0.01	0.891	0.107	0.035	0.014	1.337
1stQ	0.372	-0.037	-0.011	0.042	0.395	-0.038	0	0.057
Mean	0.651	0.02	0.028	0.598	0.663	0.027	0.057	0.75
Median	0.739	0	0.017	0.253	0.754	0.002	0.036	0.307
3rdQ	0.98	0.099	0.072	0.756	0.985	0.114	0.115	0.925

Source: Authors

We decided to create a logistic regression model for several reasons. Logistic regression does not need a linear relationship between the predictors and response. It can handle all sorts of relationships, because it uses a log transformation to the predicted odds ratio. And even though multivariate normality gives a more stable solution, predictors do not need to be multivariate normal. Additionally, the residuals do not need to be multivariate normally distributed and homoscedasticity is not needed (Spuchlakova, E.& Others, 2014).

Model is created and tested in “R”.

Table 2: Coefficients of model

	Coefficients	Standardized coefficients
Intercept	-0.852224	-2.871364
Working capital / Total assets	-0.281747	-0.093763
Retained earnings / Total assets	-4.649456	-0.833883
EBIT / Total assets	-7.082977	-0.709904
Equity / Total liabilities	-2.586588	-2.441832

Source: Authors

In **Table 2** are coefficients of the new model. A result of each logistic regression model is a probability. In this case, it is probability of default. An important step is to determine a threshold. **Table 3** shows some of the most common ways how define thresholds. For example, using a value of 0.299504 (maximized F_1) we get the confusion matrices, which are listed in **Table 4**.

Table 3: Model metrics

data max	2014 test		2015	
	threshold	value	threshold	value
F_1	0.299504	0.671375	0.298979	0.661148
F_2	0.203167	0.744390	0.224545	0.728951
$F_{0.5}$	0.381561	0.653424	0.408554	0.658281
Accuracy	0.317903	0.871327	0.389177	0.889129
Precision	0.903690	0.851064	0.846382	0.833333
Sensitivity	0.000047	1.000000	0.000036	1.000000
Specifity	0.992718	0.999635	0.997583	0.999764
Abs MCC	0.299504	0.591943	0.298979	0.593478
Min per Class Accuracy	0.230440	0.824731	0.217513	0.831197
Mean per Class Accuracy	0.248878	0.830292	0.236188	0.833729

Source: Authors

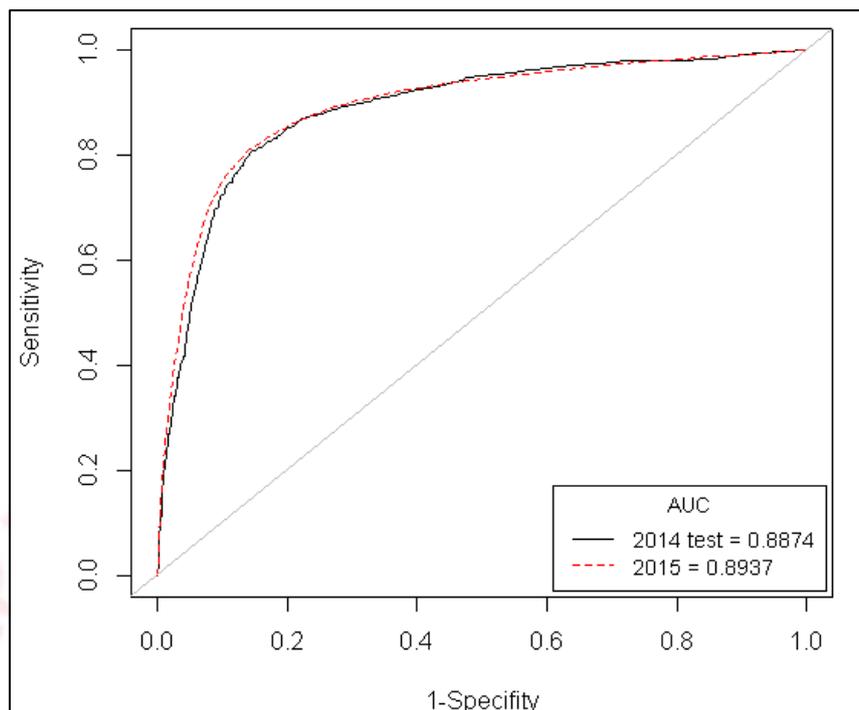
From confusion matrices in Table 4 can be seen that error rate of healthy companies is approximately 8%-10% and error rate of companies in bankruptcy is 28%-30%. This result can be considered as quite good.

Figure 1 shows the ROC curves with AUC.

Table 4: Confusion matrices for F_1 optimal threshold

		Predicted					
		2014 test			2015		
		0	1	Error rate	0	1	Error rate
Actual	0	4951	532	0.0970	35016	3095	0.081
	1	352	903	0.280	2145	5112	0.296

Source: Authors

Figure 1: ROC and AUC of model

Source: Authors

Conclusion

In the paper we present new bankruptcy prediction model for Slovak companies wherein we use modified definition of bankruptcy. The data was obtained from the register of the financial statements. New model works very well and important performance characteristics of the model are shown in **Table 5**.

Table 5: Performance characteristics of the model

	MSE	RMSE	LogLoss	Per-Class Error	AUC	Gini
2014 test	0.096	0.310	0.328	0.189	0.887	0.775
2015	0.084	0.289	0.298	0.188	0.894	0.787

Source: Authors

Acknowledgement

This research was financially supported by the Slovak Research and Development Agency – Grant NO. APVV-14-0841: Comprehensive Prediction Model of the Financial Health of Slovak Companies.

References

- Cisko, S. & Kliestik, T. (2013). *Financnymanazmentpodniku II*, EDIS Publishers, Zilina, Slovakia.
- Kocisova, K. & Misankova, M. (2014). Prediction of Default by the Use of Merton's Model and Black and Cox Model, 4th international conference on applied social science (ICASS 2014), *Advances in Education Research*, Singapore, Singapore, Vol. 51, pp. 563-568.
- Spuchlakova, E.& Others (2014). Credit Risk Measurement. *2nd International Conference on Economics and Social Science (ICESS)*, *Advances in Education Research*, Vol. 61, pp. 75-79.

