

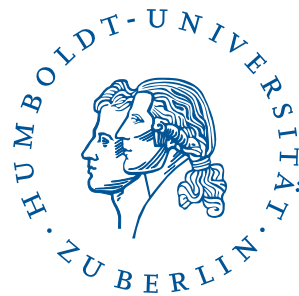
Estimating Probabilities of Default With Support Vector Machines

M.Sc. Thesis in Statistics

Rouslan Arthur Moro

Research advisors:

Prof. Dr. Wolfgang Härdle and PD Dr. Dorothea Schäfer



2006

1. COMPANY RATING METHODOLOGY

Application of statistical techniques to corporate bankruptcy started in the 60's. The first technique introduced was discriminant analysis (DA) for univariate (Beaver, 1966) and multivariate models (Altman, 1968). After DA the logit and probit models were introduced in (Martin, 1977) and (Ohlson, 1980). Nowadays these models are widely used in practice, e.g. they are at the core of the rating solutions at most European central banks. The solution in the traditional framework is a linear function (a hyperplane in a multidimensional feature space) separating successful and failing companies. A company score is computed as a value of that function. In the case of the probit and logit models the score can be directly transformed into a probability of default (PD), which denotes the probability with which a company can go bankrupt within a certain period. The major disadvantages of these popular approaches is the linearity of the solution and, in the case of logit and probit models, the prespecified form of the link function between PDs and the linear combination of predictors (Figure 1.1).

In Figure 1.1 successful and failing companies are denoted with black triangles and white quadrangles respectively. There is an equal number of companies of both classes in the sample. Following the DA and logit classification rule, which give virtually the same result, we are more likely to find a failing company above and to the right from the straight line. This may lead to a conclusion that companies with significantly negative values of op-

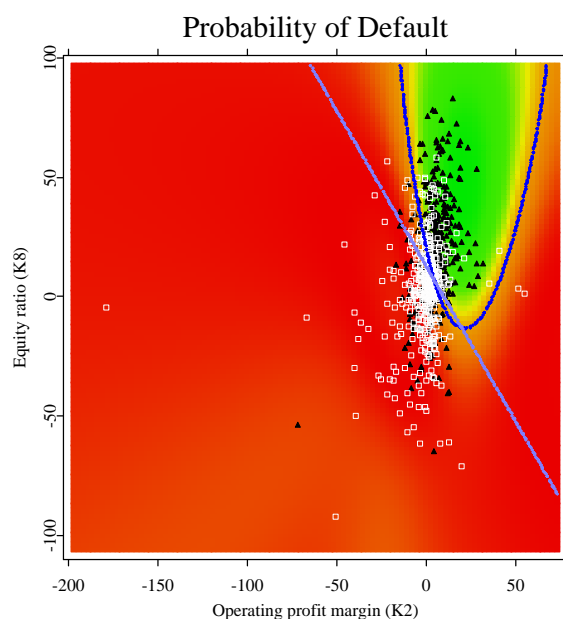


Fig. 1.1: A classification example. The boundary between the classes of solvent and insolvent companies was estimated using DA and logit regression (two indistinguishable linear boundaries) and an SVM (a non-linear boundary).

erating profit margin and equity ratio can be classified as successful. This, for example, allows for companies with liabilities much greater than total assets to be classified as successful. Such a situation is avoided by using a non-linear classification method, such as the support vector machine (SVM), which produces a non-linear boundary.

Following a traditional approach we would expect a monotonic relationship between predictors and PDs, like the falling relation for the interest coverage ratio (Figure 1.2). However, in reality this dependence is often non-monotonic as for such important indicators as the company size or net income change. In the latter case companies that grow too fast or too slow have a higher probability of default. That is the reason for contemplating non-linear techniques as alternatives. Two prominent examples are recursive partitioning (Frydman,

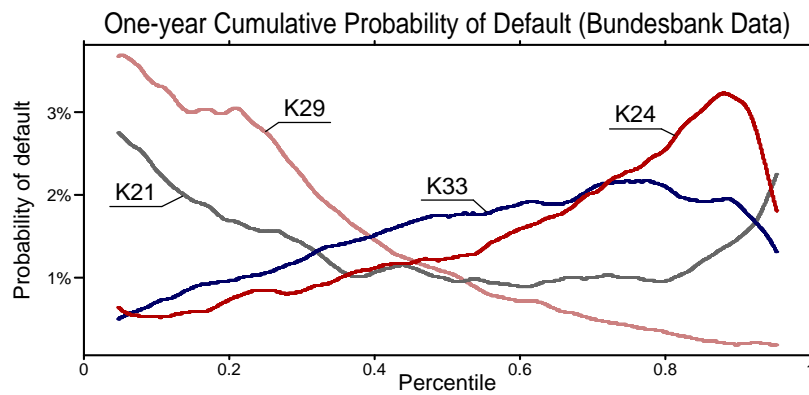


Fig. 1.2: One year cumulative PDs evaluated for several financial ratios on the German Bundesbank data. The ratios are net income change, K21 (gray), net interest ratio, K24 (red), interest coverage ratio, K29 (pink) and logarithm of total assets, K33 (blue). The k -nearest-neighbours procedure was used with the size of the window being around 8% of all observations. The total number of observations is 553500.

Altman & Kao, 1985) and neural networks (Tam & Kiang, 1992). Despite the strength of the two approaches they have visible drawbacks: orthogonal division of the data space in recursive partitioning that is usually not justified and heuristic model specification in neural networks.

Recursive partitioning, also known as classification and regression trees performs classification by orthogonally dividing the data space. At each step only a division (split) along one of the axes is possible. The axis is chosen such, that a split along it reduces the variance in each of the subspaces and maximises the variance between them. Entropy based criteria can also be used. The visible drawback is the orthogonal division itself which imposes severe restrictions on the smoothness of the classifying function and may not adequately capture the correlation structure between the variables. Orthogonal division means that the separating hyperplane can only consist of orthogonal segments parallel to the coordinate grid, whereas the boundary

between the classes has a smoothly changing gradient.

The neural network (NN) is a network of linear classifiers (neurons) that are connected with one another in a prespecified way. The outputs of some of the neurons are inputs for others. The performance of a NN greatly depends on its structure that must be adapted for solving different problems. The network must be designed manually that requires a substantial experience from the operator. Moreover, NNs mostly do not provide a global solution but only a local one. This feature, as well as too much heuristics create many obstacles on the way of using NNs at the rating departments of banks.

We would like to have a model that is able to select a classifying function based on very general criteria. The SVM is a statistical technique that in many applications, such as optical character recognition and medical diagnostics, showed very good performance. It has a flexible solution and is controlled by adjusting only few parameters. Its overall good performance and flexibility make the SVM a suitable candidate (Härdle, Moro & Schäfer, 2004).

Within a rating methodology each company is described by a set of variables x , such as financial ratios. Financial ratios, such as debt ratio (leverage) or interest coverage (earnings before interest and taxes to interest) characterise different sides of company operation. They are constructed on the basis of balance sheets and income statements. For example, the Bundesbank uses 32 ratios (predictors) computed using the company statements from its corporate bankruptcy data base. The predictors and basic statistics are given in Table 4.1. The whole Bundesbank data base covers the period 1987–2005 and consists of 553500 anonymised statements of solvent and insolvent companies. Most companies appear in the database several times in different years.

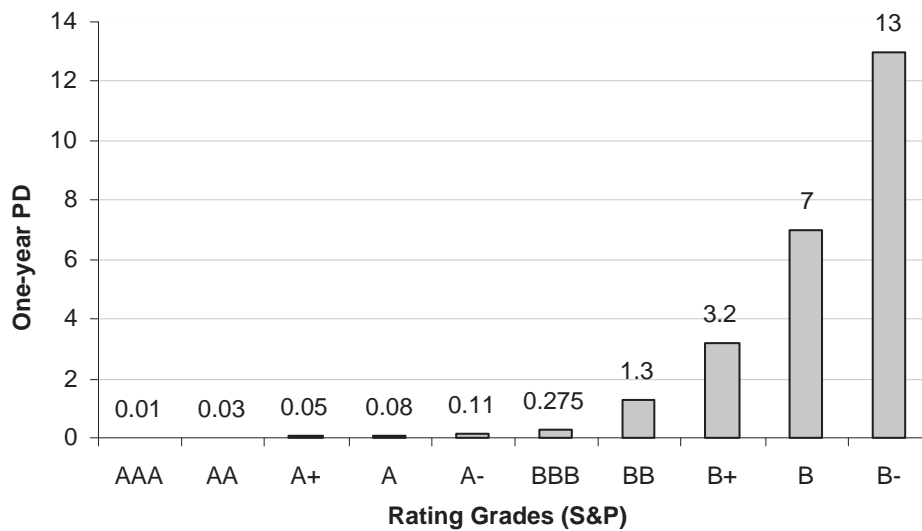


Fig. 1.3: One year probabilities of default for different rating grades (Füser, 2002).

The class y of a company can be either $y = -1$ ('successful') or $y = 1$ ('bankrupt'). Initially, an unknown classifier function $f : x \rightarrow y$ is estimated on a training set of companies (x_i, y_i) , $i = 1, \dots, n$. The training set represents the data for companies which are known to have survived or gone bankrupt. In order to obtain PDs from the estimated scores f , rating practitioners usually rely on prespecified rating classes (i.e. BBB, C, AA, etc.). A certain range of scores and PDs belong to each rating class. The ranges are computed on the basis of historical data. To derive a PD for a newly scored company its score f is compared with the historical values of f 's for each class. Basing on the similarity of the scores a company is assigned to one particular class. The PD of this class becomes the PD of the company.

Company bond ratings play an important role in determining the cost of debt refinancing since they reflect the probability of defaulting on the debt (Figure 1.3).

2. THE SVM APPROACH

The SVM (Vapnik, 1995) is a regression (and classification) technique that is based on margin maximisation (Figure 2.1) between two data classes. The margin is the distance between the hyperplanes bounding each class where in a linear perfectly separable case no observation may lie. The classifier function used by the linear SVM is a hyperplane symmetrically surrounded with a margin zone. It can be shown (Härdle, Moro & Schäfer, 2004) that by maximising the margin one reduces the complexity of such a classifier. By applying kernel techniques the SVM can be extended to learn non-linear classifying functions (Figure 2.2).

In Figure 2.1 misclassifications are unavoidable when using linear classifying functions (linearly non-separable case). To account for misclassifications the penalty ξ_i is introduced, which is related to the distance from the hyperplane bounding observations of the same class to observation i . $\xi_i > 0$ if a misclassification occurs. All observations satisfy the following two constraints:

$$y_i(x_i^\top w + b) \geq 1 - \xi_i, \quad (2.1)$$

$$\xi_i \geq 0. \quad (2.2)$$

With the normalisation of w , b and ξ_i as in (2.1) the margin equals to $2/\|w\|$. The convex objective function to be minimised given the constraints

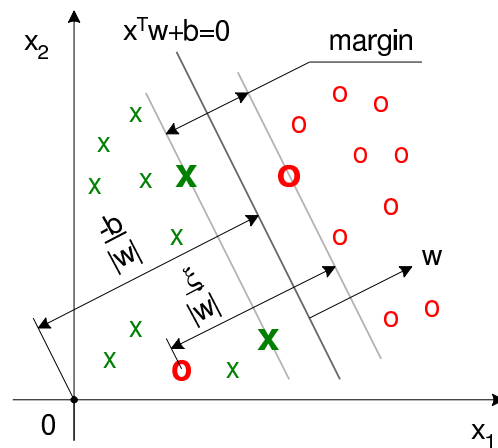


Fig. 2.1: The separating hyperplane $x^\top w + b = 0$ and the margin in a non-separable case. The observations marked with bold crosses and zeros are support vectors. The hyperplanes bounding the margin zone equidistant from the separating hyperplane are represented as $x^\top w + b = 1$ and $x^\top w + b = -1$.

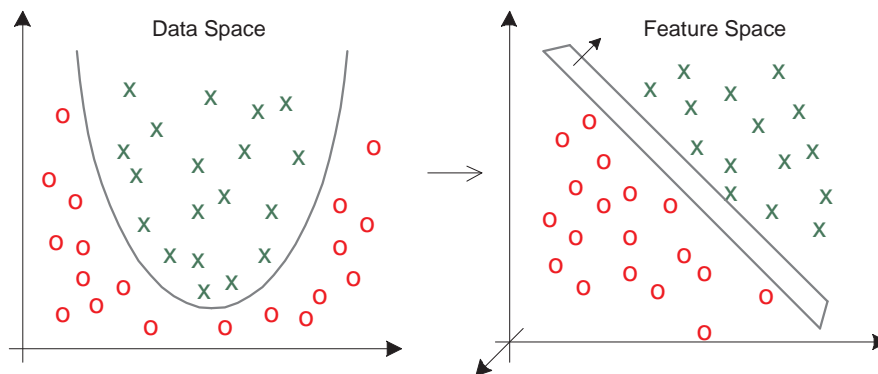


Fig. 2.2: Mapping from a two-dimensional data space into a three-dimensional space of features $\mathbb{R}^2 \mapsto \mathbb{R}^3$ using a quadratic kernel function $K(x_i, x_j) = (x_i^\top x_j)^2$. The three features correspond to the three components of a quadratic form: $\tilde{x}_1 = x_1^2$, $\tilde{x}_2 = \sqrt{2}x_1x_2$ and $\tilde{x}_3 = x_2^2$, thus, the transformation is $\Psi(x_1, x_2) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)^\top$. The data separable in the data space with a quadratic function will be separable in the feature space with a linear function. A non-linear SVM in the data space is equivalent to a linear SVM in the feature space. The number of features will grow fast with the dimension of the data d and the degree of the polynomial kernel p , which equals 2 in our example, making the closed-form representation of Ψ such as here practically impossible

(2.1) and (2.2) is:

$$\frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i. \quad (2.3)$$

The constrained optimisation problem is:

$$\min_{w_k, b, \xi_i} \max_{\alpha_i \geq 0, \mu_i \geq 0} L_P, \quad (2.4)$$

for all $i = 1, \dots, n$ and $k = 1, \dots, d$. Here L_P is the Lagrange functional for the primal problem (2.4):

$$L_P = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i \{y_i (x_i^\top w + b) - 1 + \xi_i\} - \sum_{i=1}^n \mu_i \xi_i. \quad (2.5)$$

The Karush-Kuhn-Tucker (KKT) first order optimality conditions (Gale, Kuhn & Tucker, 1951) that must hold for all $i = 1, \dots, n$ are:

$$\nabla_w = w - \sum_{i=1}^n \alpha_i y_i x_i = 0; \quad (2.6)$$

$$\frac{\partial L_P}{\partial b} = \sum_{i=1}^n \alpha_i y_i = 0; \quad (2.7)$$

$$\frac{\partial L_P}{\partial \xi_i} = C - \alpha_i - \mu_i = 0; \quad (2.8)$$

$$\alpha_i \geq 0; \quad (2.9)$$

$$\alpha_i \{1 - \xi_i - y_i (w^\top x_i)\} = 0; \quad (2.10)$$

$$\mu_i \geq 0; \quad (2.11)$$

$$\mu_i \xi_i = 0. \quad (2.12)$$

After substituting the KKT conditions into (2.5) we can obtain the La-

grangian L_D for the dual problem:

$$\min_{\alpha_i \geq 0, \delta_i \geq 0, \gamma_i \geq 0, \beta \geq 0} L_D, \quad (2.13)$$

where L_D is:

$$L_D = \frac{1}{2} w(\alpha)^\top w(\alpha) - \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \delta_i \alpha_i + \sum_{i=1}^n \gamma_i (\alpha_i - C) - \beta \sum_{i=1}^n \alpha_i y_i. \quad (2.14)$$

α_i , δ_i , γ_i and β are Lagrange multipliers for all $i = 1, \dots, n$. The function $w(\alpha)^\top w(\alpha)$ is a scalar product in some Hilbert space (hence the notation).

For a linear SVM:

$$w(\alpha)^\top w(\alpha) = \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^\top x_j. \quad (2.15)$$

For obtaining non-linear classifying functions in the data space a more general form is applicable:

$$w(\alpha)^\top w(\alpha) = \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(x_i, x_j). \quad (2.16)$$

The parameter C called capacity is related to the width of the margin zone.

The smaller the C is, the bigger margins are possible.

The function $K(x_i, x_j)$ is called a kernel function. Since it has a closed form representation, the kernel is a convenient way of mapping low dimensional data into a highly dimensional (often infinitely dimensional) space of features. It must satisfy the Mercer conditions (Mercer, 1909), i.e. be symmetric and semipositive definite or, in other words, represent a scalar product in some Hilbert space (Weyl, 1928).

In our study we applied an SVM with an anisotropic Gaussian kernel

$$K(x_i, x_j) = \exp \left\{ -(x_i - x_j)^\top r^{-2} \Sigma^{-1} (x_i - x_j) / 2 \right\}, \quad (2.17)$$

where r is a coefficient and Σ is a variance-covariance matrix. The coefficient r is related to the complexity of classifying functions: the higher the r is, the lower is the complexity. If kernel functions allow for sufficiently rich feature spaces, the performances of SVMs are comparable in terms of out-of-sample forecasting accuracy (Vapnik, 1995).

3. COMPANY SCORE EVALUATION

The company score is computed as:

$$f(x) = x^T w + b, \quad (3.1)$$

where $w = \sum_{i=1}^n \alpha_i y_i x_i$ and $b = \frac{1}{2}(x_+ + x_-)^T w$; x_+ and x_- are the observations from the opposite classes for which constraint (2.1) becomes equality. By substituting the scalar product with a kernel function we will derive a non-linear score function:

$$f(x) = \sum_{i=1}^n K(x_i, x) \alpha_i y_i + b. \quad (3.2)$$

The non-parametric score function (3.2) does not have a compact closed form representation. This may necessitate the use of graphical tools for its visualisation.

4. VARIABLE SELECTION

In this section we describe the procedure and the graphical tools for selecting the variables of the SVM model used in forecasts. We have two most important criteria of model accuracy: the accuracy ratio (AR), which will be used here as a criterion for model selection, (Figure 4.1) and the percentage of correctly classified out-of-sample observations. Higher values indicate better model accuracy.

We start model selection from the simplest, i.e. univariate models and then pick up the one with the highest AR. The problem that arises is how to determine the variable which provides the highest AR across possible data samples. For a parametric model we would need to estimate the distribution of the coefficients at the variables and, hence, their confidence intervals. This approach, however, is practically irrelevant for non-parametric models.

Instead we can compare goodness of models with respect to some accuracy measure, in our case AR. Firstly we will estimate the distributions of AR for different models. This can be done using bootstrapping (Horowitz, 2001). We randomly select training and validation sets as subsamples of 500 solvent and 500 insolvent companies each. We used the 50/50 ratio since this is the worst case with the minimum AR. The two sets are not overlapping, i.e. do not contain common observations. For each of these sets we apply the SVM with parameters that provide the highest AR for bivariate models (Figure 4.2) and estimate ARs. Then we perform a Monte Carlo experiment: repeat the

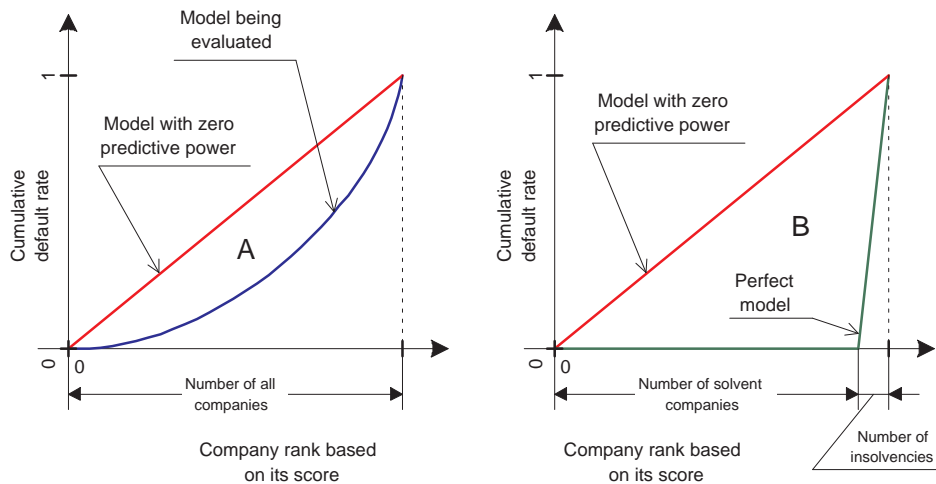


Fig. 4.1: The power curves for a perfect (green), random (red) and some real (blue) classification models. The AR is the ratio of two areas A/B . It lies between 0 for a random model with no predictive power and 1 for a perfect model.

generation of subsamples and computing of ARs 100 times. Each time we will record the ARs and then estimate their distribution.

At the end of this procedure we obtain an empirically estimated distribution of AR on bootstrapped subsamples. The median AR provides a robust measure to compare different variables as predictors. The same approach can be used for comparing SVM with DA and logit regression in terms of predictive power. We compute AR for the same subsamples with the SVM, DA and logit models. The median improvements in AR for the SVM over DA and the SVM over the logistic regression are also reported below (Figure 4.6).

We will start this procedure with all univariate models with 33 variables K1-K9, K11-K33 as they are denoted at the Bundesbank and variable K10, which is a standard normal random variable used as a reference (Table 4.1). For each model the resulting distribution of ARs will be represented as box plots (Figure 4.3). The red line depicts medians. The box within each box plot

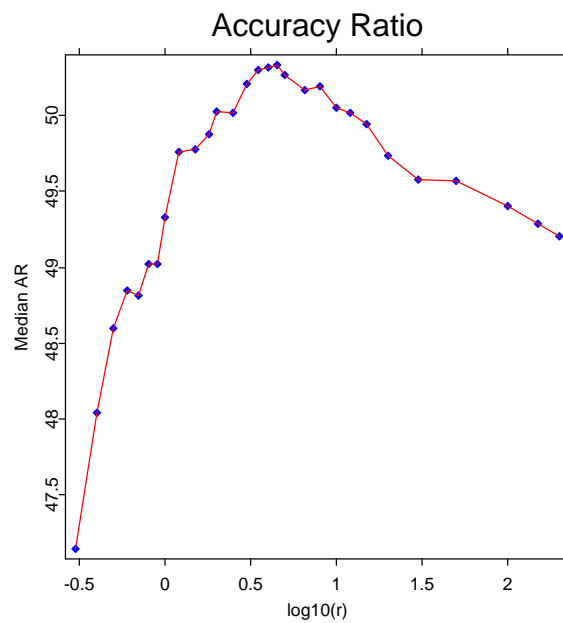


Fig. 4.2: The relationship between an accuracy measure (AR) and the coefficient r in the SVM formulation. Higher r 's correspond to less complex models. The median ARs were estimated on 100 bootstrapped subsamples of 500 solvent and 500 insolvent companies both in the training and validation sets. A bivariate SVM with the variables K5 and K29 was used. We will be using $r = 4$ in all SVMs used in this chapter.

shows the interquartile range (IQR), while the whiskers span to the distance of $3/2$ IQR in each direction from the median. Outliers beyond that range are denoted with circles.

Basing on Figure 4.3 we can conclude that variables K5 (Debt Cover) and K29 (Interest Coverage Ratio) provide the highest median AR around 50%. We can also notice that variables K12, K26 and K28 have a very low accuracy: their median ARs do not exceed 11.5%. The model based on random variable K10 has AR equal zero, in other words, it has no predictive power whatsoever. For the next step we will select variable K5 that was included in the best univariate model.

For bivariate models we will select the best predictor from the univariate models (K5) and one of the rest that delivers the highest AR (K29) (Figure 4.4). This procedure will be repeated for each new variable added. The AR is growing until the model has eight variables, then it slowly declines. Median ARs for the models with eight variables are shown in Figure 4.5.

We have also conducted experiments with subsamples of the size of 5000 observations. The change of median was extremely small (one–two orders of magnitude smaller than the interquartile range). The interquartile range got narrower as it was expected, i.e. the difference between models with bigger samples is only more statistically significant. Thus, proving that if the difference is significant on a sample of 1000 observations, it can be guaranteed that this will remain so for bigger samples.

The SVM based on variables K5, K29, K7, K33, K18, K21, K24, K33 and K9 attains the highest median AR of around 60.0%. For comparison we plot an improvement in AR for the SVM vs. DA and logit regression on the same 100 subsamples. The data used in the DA and logit models were

Tab. 4.1: Summary Statistics. q_α is an α quantile. IQR is the interquartile range.

Var.	Name	Group	$q_{0.01}$	Median	$q_{0.99}$	IQR
K1	Pre-tax profit margin	Profitability	-26.9	2.3	78.5	5.9
K2	Operating profit margin	Profitability	-24.6	3.8	64.8	6.3
K3	Cash flow ratio	Liquidity	-22.6	5.0	120.7	9.4
K4	Capital recovery ratio	Liquidity	-24.4	11.0	85.1	17.1
K5	Debt cover	Liquidity	-42.0	17.1	507.8	34.8
K6	Days receivable	Activity	0.0	31.1	184.0	32.7
K7	Days payable	Activity	0.0	23.2	248.2	33.2
K8	Equity ratio	Financing	0.3	14.2	82.0	21.4
K9	Equity ratio (adj.)	Financing	0.5	19.3	86.0	26.2
K10	Random Variable	Test	-2.3	0.0	2.3	1.4
K11	Net income ratio	Profitability	-29.2	2.3	76.5	5.9
K12	Leverage ratio	Leverage	0.0	0.0	164.3	4.1
K13	Debt ratio	Liquidity	-54.8	1.0	80.5	21.6
K14	Liquidity ratio	Liquidity	0.0	2.0	47.9	7.1
K15	Liquidity 1	Liquidity	0.0	3.8	184.4	14.8
K16	Liquidity 2	Liquidity	2.7	63.5	503.2	58.3
K17	Liquidity 3	Liquidity	8.4	116.9	696.2	60.8
K18	Short term debt ratio	Financing	2.4	47.8	95.3	38.4
K19	Inventories ratio	Investment	0.0	28.0	83.3	34.3
K20	Fixed assets ownership r.	Leverage	1.1	60.6	3750.0	110.3
K21	Net income change	Growth	-50.6	3.9	165.6	20.1
K22	Own funds yield	Profitability	-510.5	32.7	1998.5	81.9
K23	Capital yield	Profitability	-16.7	8.4	63.1	11.0
K24	Net interest ratio	Cost struct.	-3.7	1.1	36.0	1.9
K25	Own funds/pension prov. r.	Financing	0.4	17.6	84.0	25.4
K26	Tangible asset growth	Growth	0.0	24.2	108.5	32.6
K27	Own funds/provisions ratio	Financing	1.7	24.7	89.6	30.0
K28	Tangible asset retirement	Growth	1.0	21.8	77.8	18.1
K29	Interest coverage ratio	Cost struct.	-1338.6	159.0	34350.0	563.2
K30	Cash flow ratio	Liquidity	-14.1	5.2	116.4	8.9
K31	Days of inventories	Activity	0.0	42.9	342.0	55.8
K32	Current liabilities ratio	Financing	0.3	58.4	98.5	48.4
K33	Log of total assets	Other	4.9	7.9	13.0	2.1

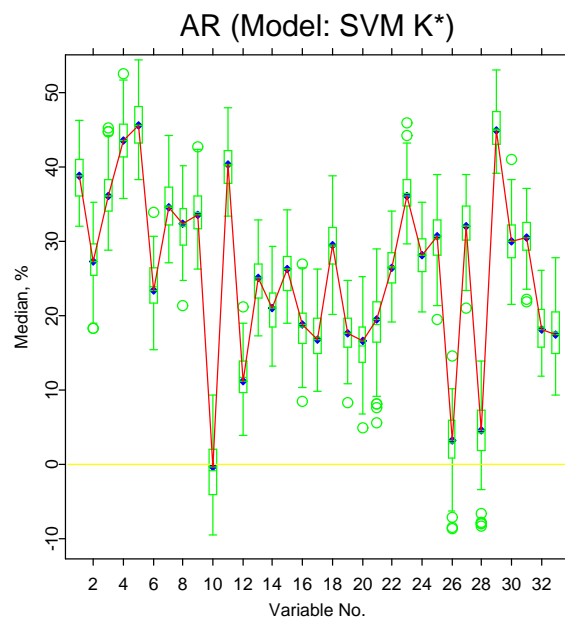


Fig. 4.3: Accuracy ratios for univariate SVM models. Box-plots are estimated basing on 100 random subsamples. The AR for the model containing only random variable K10 is zero.

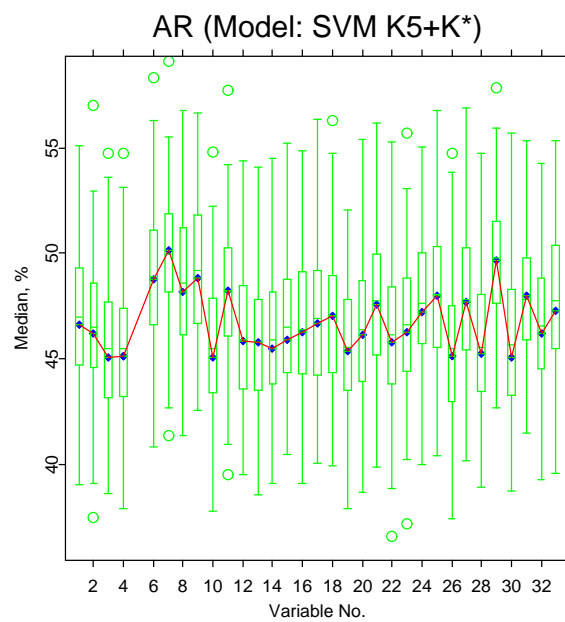


Fig. 4.4: Accuracy ratios for bivariate SVM models. Each model includes variable K5 and one of the remaining. Box-plots are estimated basing on 100 random subsamples.

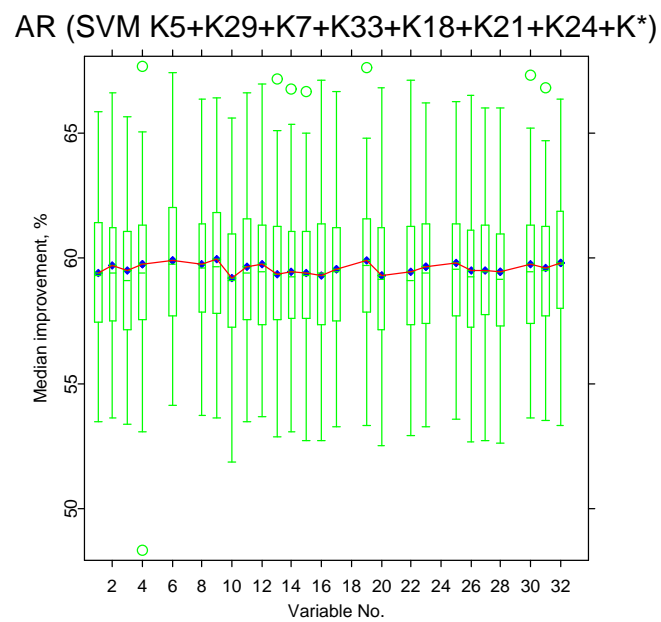


Fig. 4.5: Accuracy ratios for SVM models with eight variables. Each model includes variables K5, K29, K7, K33, K18, K21, K24 and one of the remaining. Box-plots are estimated basing on 100 random subsamples.

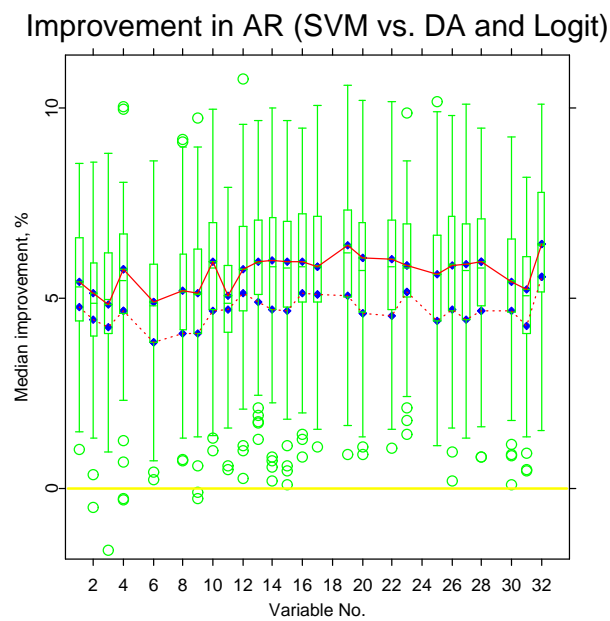


Fig. 4.6: Median improvement in AR. SVM vs. DA (the upper line) and SVM vs. logit regression (the lower line). Box-plots are estimated basing on 100 random subsamples for the case of DA. Each model includes variables K5, K29, K7, K33, K18, K21, K24 and one of the remaining

processed as following: if $x_i < q_{0.05}(x_i)$ then $x_i = q_{0.05}(x)$ and if $x_i > q_{0.95}(x_i)$ then $x_i = q_{0.95}(x_i)$; $i = 1, 2, \dots, 8$; $q_\alpha(x_i)$ is an α quantile of x_i . Thus, the DA and logit regression applied were *robust versions* not sensitive to outliers. Without such a procedure the improvement would be much higher.

Figure 4.6 represents the absolute improvement for SVM over robust DA (upper line) and SVM over robust logit regression (lower line). We can see that for all models containing variables K5, K29, K7, K33, K18, K21, K24 and one of the remaining variables the median AR was always higher for the SVM. Thus, the SVM model is always dominating in accuracy DA and logit regression with regard to AR. In terms of the percentage of correctly classified out-of-sample observations a similar result is achieved.

5. CONVERSION OF SCORES INTO PDS

There is another way to look at the company score. It defines the distance between companies in terms of the distance to the boundary between the classes. The lower is the score, the farther is a company from the class of bankrupt companies, therefore, we can assume, the lower PD it must have. This means that the dependence between scores and PDs is assumed to be monotonous. This is the only kind of dependence that was assumed in all rating models mentioned in this chapter and the only one we use for PD calibration.

The conversion procedure consists of the estimation of PDs for the observation of the training set with a subsequent monotonisation (step one and two) and the computation of a PD for some new company (step three).

Step one is the estimation of PDs for the companies of the training set. We used kernel techniques to preliminary evaluate PDs for observation i from the training set, $i = 1, 2, \dots, n$:

$$\widetilde{PD}(x_i) = \frac{\sum_{j=1}^n K_h(x_i, x_j) I_{\{y_j=1\}}}{\sum_{j=1}^n K_h(x_i, x_j)} \quad (5.1)$$

Here a k -nearest-neighbour Gaussian kernel was used. h is the kernel bandwidth.

The preliminary PDs evaluated in this way are not necessarily a monotonical function of the score. The monotonisation of \widetilde{PD}_i , $i = 1, 2, \dots, n$

is achieved at step two using the Pool Adjacent Violator (PAV) algorithm ((Barlow, Bartholomew, Bremmer & Brunk, 1972) and (Mammen, 1991)). As a result we obtain monotonised probabilities of default $PD(x_i)$ for the observations of the training set.

Finally, at step three the PDs are computed for any observation described with x as an interpolation between two PDs of the neighbouring, in terms of the score, observations from the training set, x_i and x_{i-1} , $i = 2, 3, \dots, n$:

$$PD(x) = PD(x_i) + \frac{f(x) - f(x_{i-1})}{f(x_i) - f(x_{i-1})} \{PD(x_i) - PD(x_{i-1})\}. \quad (5.2)$$

If the score for an observation x lies beyond the range of scores for the training set, then $PD(x)$ equals to the score of the first neighbouring observation of the training set.

Figure 5.1 is an example of the cumulative PD curve (power curve) and estimated PDs for a subsample of 200 companies. The PD curve has a plateau area for the observations with a high score. Default probabilities can change from 15% to 80% depending on the score.

We will be following a common in finance convention and use the red colour to highlight negative information and green and blue to convey positive information. Therefore, we would like to code PDs with colours ranging from red for the highest PD to blue-green for the most solvent company.

The graphs that show the data and PDs in the dimensions of variables K33 and K29 for different complexities of the SVM are represented in Figures 5.2–5.4. The three figures correspond consequently to three SVMs with high, average and high complexity. The outliers that lie beyond the 5% and 95% quantiles are plotted at the rand. The contour lines separating the rating classes can also be added to the graph as illustrated by Figure 5.5.

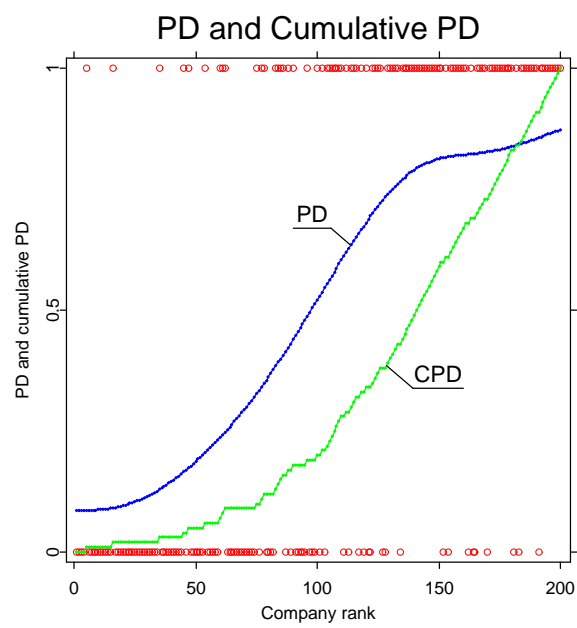


Fig. 5.1: PD (blue line) and cumulative PD (green line) estimated with the SVM for a subsample of 200 observation from the Bundesbank data. The variables were included into the model that achieved the highest AR: K5, K29, K7, K33, K18, K21, K24 and K9. The higher is the score, the higher is the rank of a company.

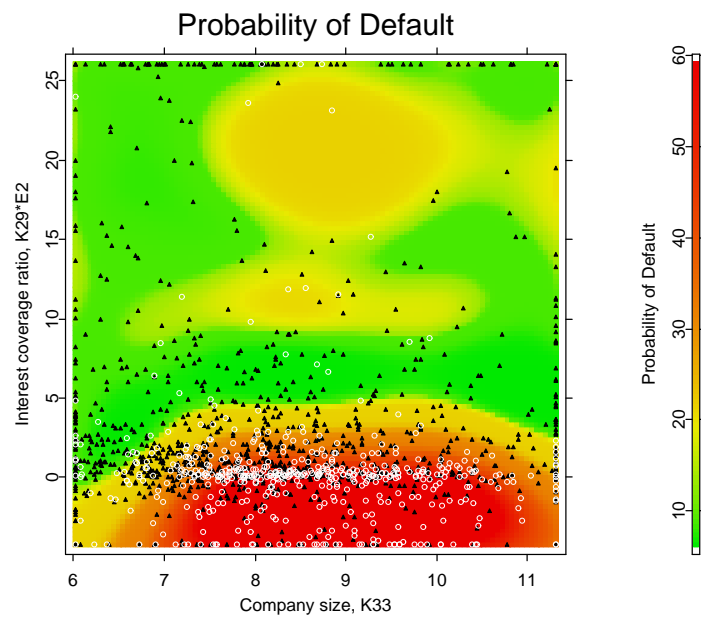


Fig. 5.2: Probability of default estimated for a random subsample of 500 failing and 500 surviving companies plotted for the variables K33 and K29. An SVM of high complexity with the radial basis kernel $0.5\Sigma^{1/2}$ was used.

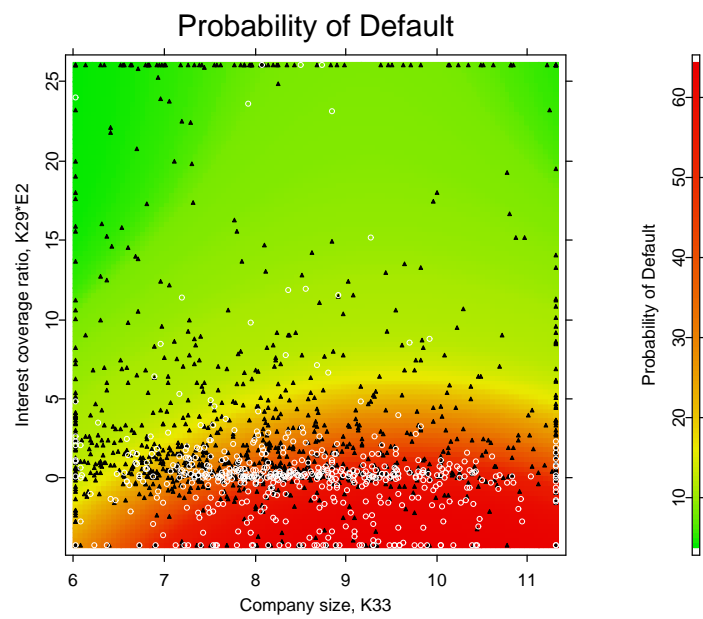


Fig. 5.3: Probability of default estimated for a random subsample of 500 failing and 500 surviving companies plotted for the variables K33 and K29. An SVM of average complexity with the radial basis kernel $4\Sigma^{1/2}$ was used. The case of the highest out-of-sample classification accuracy.

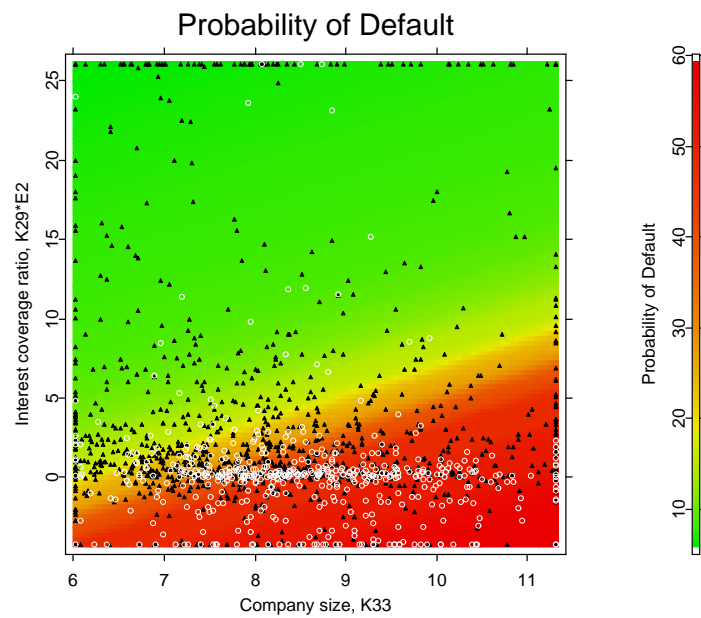


Fig. 5.4: Probability of default estimated for a random subsample of 500 failing and 500 surviving companies plotted for the variables K33 and K29. An SVM of low complexity with the radial basis kernel $100\Sigma^{1/2}$ was used.

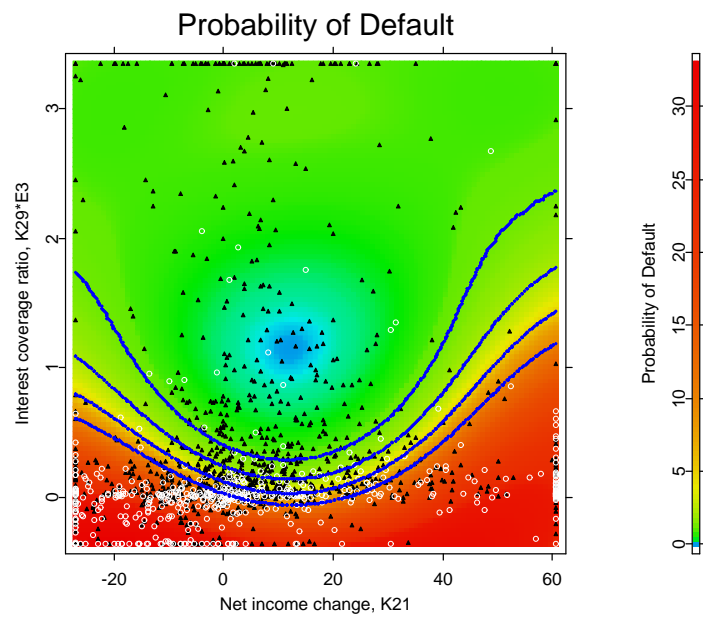


Fig. 5.5: Probability of default plotted for the variables K21 and K29. The boundaries of five risk classes are shown in blue, which correspond to the rating classes: BBB and above (investment grade), BB, B+, B, B- and lower.

6. CONCLUSION

In this paper we show that a rating model based on SVMs is dominating traditional linear parametric approaches such as DA and logit regression. The forecasting accuracy improvement is significant already for small samples. It was demonstrated how non-linear non-parametric techniques can be a basis for a rating model. The implementation of an SVM rating model and its extensive testing on the data of the German Bundesbank was performed. We believe that non-parametric techniques such as SVM will become more commonplace in the rating community since they better represent the data and provide higher forecasting accuracy.

7. ACKNOWLEDGEMENTS

I am thankful to the German Bundesbank for providing access to the unique database of the financial statements of German companies. The data analysis took place on the premises of the German Bundesbank in Frankfurt. The work was funded by the German Academic Exchange Service (DAAD) and the German Bundesbank. I also appreciate the support of the German Research Foundation (DFG) through the SFB 649 of Humboldt-Universität zu Berlin.

BIBLIOGRAPHY

- Altman, E., 1968: Financial ratios, discriminant analysis and the prediction of corporate bankruptcy. *The Journal of Finance* 589–609.
- Barlow, R. E., J. M. Bartholomew, J. M. Bremner, and H. D. Brunk, 1972: *Statistical Inference Under Order Restrictions*. John Wiley & Sons, New York, NY.
- Beaver, W., 1966: Financial ratios as predictors of failures. empirical research in accounting: Selected studies. *Journal of Accounting Research* 71–111. supplement to vol. 5.
- Frydman, H., E. Altman, and D.-L. Kao, 1985: Introducing recursive partitioning for financial classification: The case of financial distress. *The Journal of Finance*, **40**(1), 269–291.
- Füser, K., 2002: Basel II – was muß der Mittelstand tun? [http://www.ey.com/global/download.nsf/Germany/Mittelstandsrating/\\$file/Mittelstandsrating.pdf](http://www.ey.com/global/download.nsf/Germany/Mittelstandsrating/$file/Mittelstandsrating.pdf).
- Gale, D., H. W. Kuhn, and A. W. Tucker, 1951: *Linear Programming and the Theory of Games, in Activity Analysis of Production and Allocation*, T. C. Koopmans (ed.). John Wiley & Sons, New York, NY, 317–329.
- Härdle, W., R. A. Moro, and D. Schäfer, 2004: *Predicting Bankruptcy with*

-
- Support Vector Machines in Statistical Tools in Finance*, W. Härdle (ed.). Springer Verlag, Berlin.
- Härdle, W., M. Müller, S. Sperlich, and A. Werwatz, 2004: *Nonparametric and Semiparametric Models*. Springer Verlag, Berlin.
- Härdle, W. and L. Simar, 2003: *Applied Multivariate Statistical Analysis*. Springer Verlag.
- Horowitz, J. L., 2001: *The Bootstrap*, J. J. Heckman and E. E. Leamer (eds.), volume 5. Elsevier Science B.V., 3159–3228.
- Mammen, E., 1991: Estimating a smooth monotone regression function. *Annals of Statistics*, **19**, 724–740.
- Martin, D., 1977: Early warning of bank failure: A logit regression approach. *Journal of Banking and Finance*, (1), 249–276.
- Mercer, J., 1909: Functions of positive and negative type and their connection with the theory of integral equations. *Philosophical Transactions of the Royal Society of London*, **209**, 415–446.
- Ohlson, J., 1980: Financial ratios and the probabilistic prediction of bankruptcy. *Journal of Accounting Research* 109–131.
- Tam, K. and M. Kiang, 1992: Managerial application of neural networks: the case of bank failure prediction. *Management Science*, **38**(7), 926–947.
- Vapnik, V. N., 1995: *The Nature of Statistical Learning Theory*. Springer, New York.
- Weyl, H., 1928: *Gruppentheorie und Quantenmechanik*. Hirzel, Leipzig.