



Louvain Institute of Data Analysis and  
Modeling in economics and statistics

Institute of Statistics, Biostatistics and Actuarial Sciences

---

## **LDATS2470: Statistical Machine Learning and High Dimensional Data Analysis**

Lending club database: insurance policies **classification**

---

*Authors:*

GODFRIN Siméon ([DATS2M](#)- 60921800)  
RGHIOUI Mehdi ([DATS2M](#)- 29912100)  
WARNAUTS Aymeric ([DATS2M](#)- 87031800)

*Professor:*

HAFNER Christian

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Preprocessing</b>	<b>2</b>
<b>3</b>	<b>Optimization by cross-validation</b>	<b>3</b>
<b>4</b>	<b>Linear Support Vector Machine</b>	<b>4</b>
4.1	Linear separability . . . . .	5
4.2	Results . . . . .	6
<b>5</b>	<b>Non-linear Support Vector Machine</b>	<b>7</b>
5.1	Radial Kernel SVM . . . . .	8
5.2	Polynomial Kernel SVM . . . . .	9
<b>6</b>	<b>PCA and visualisation</b>	<b>9</b>
6.1	PCA visualisation for linear Kernel ( $C = 19$ ) . . . . .	10
6.2	PCA visualisation for radial Kernel ( $C = 1987.42, \gamma = 0.0301$ ) . . . . .	11
<b>7</b>	<b>Conclusion</b>	<b>11</b>
<b>8</b>	<b>Appendices</b>	<b>12</b>
8.1	Graph and results about Kernel RBF using a Quadratic loss . . . . .	12
8.2	Polynomial Kernel results using hinge loss . . . . .	12
8.3	R Code for linear kernel grid search & CV . . . . .	12
8.4	R Code for RBF Kernel grid search & CV . . . . .	14
8.5	R Code for polynomial Kernel grid search & CV . . . . .	17
8.6	R Code for PCA projection . . . . .	19
	<b>References</b>	<b>21</b>

## 1 Introduction

For this project, we will work on the *Lending Club* loan database as this is one of the richest loan dataset that provides all loan information for loans granted by the lending club between 2007 and 2011 in the United States. This project will focus on the application of the *Support Vector Machine* which is a supervised machine learning algorithm used for classification and regression analysis. We expect the implementation of Support Vector Machine (SVM) to be useful as this method is able to construct non-linear decision boundaries between different classes of data, which can be helpful for capturing complex relationships between loan default and the various features that are used to predict it. Moreover, we will take advantage of its robustness to high-dimensional data and its ability to handle imbalanced datasets by using techniques such as cost-sensitive learning and class weighting, as we observe in the following table a non symmetric distribution of observations across the levels of our target feature *loan\_status*:

	Frequency	Percentage	Cumulative percentage
Fully Paid	32 950	85.4 %	85.5 %
Charged Off	5 627	14.6 %	100 %
Total	38 577	100 %	100 %

Table 1: *loan\_status* levels frequencies

As the choice between *hinge loss* and *quadratic loss* depends on the specific question being addressed, the desired trade-off between accuracy and robustness and the type of model used, we will explore both to find an appropriate one. However, the *hinge loss* function is typically used for classification problems with binary labels as our target feature *loan\_status* because it is a convex function that is designed to encourage the model to correctly classify data points, while also allowing an error rate. Finally, the *hinge loss* function is recognized for its robustness to outliers compared to the quadratic loss function. This characteristic enables it to effectively handle noisy data. Additionally, SVMs employing the *hinge loss* exhibit sparsity by having fewer support vectors. This sparsity not only contributes to computational efficiency but also enhances their overall performance. We defined this loss as:

$$\mathcal{L}(y, f(x)) = \max(0, 1 - yh(x))$$

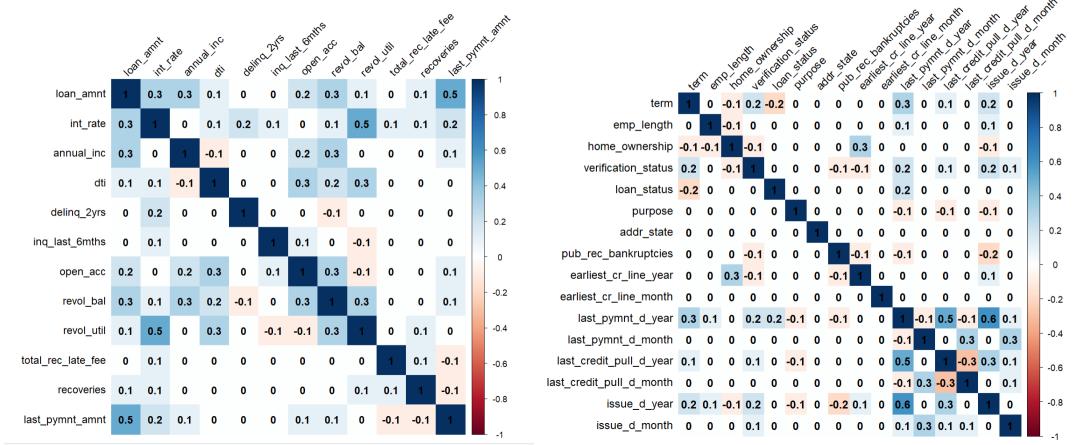
Where  $y$  is the true label,  $h(x)$  is the predicted score given by the fitted hyperplane equation as  $h(x) = w^T x + b$ . For a given data point, if the function  $h(x)$  yields a positive value, we classify it as belonging to the positive class; conversely, if the function yields a negative value, we classify it as belonging to the negative class. By multiplying  $y$  by  $h(x)$ , the hinge loss function effectively increases the penalty for incorrect predictions. When the predicted score  $h(x)$  shares the same sign as the true label  $y \in -1, 1$ , the product  $yh(x)$  is positive, indicating a correct prediction and resulting in a null loss for that data point. Conversely, when the product is negative, it signifies an incorrect prediction, and the loss is proportional to the magnitude of the product  $yh(x)$ . Larger magnitudes correspond to greater losses, reflecting the model's degree of error.

## 2 Preprocessing

For our preprocessing part we will not perform specific treatment for outliers, as for *SVM* classification of loan default, outliers may represent borrowers with unique characteristics or circumstances that could have a significant impact on their likelihood of default and thus removing these contracts could result in a biased analysis and potentially lead to incorrect predictions of loan default risk.

Moreover we will discard some features regarding correlation plots of both continuous and categorical features and features related to text and date types. However, regarding the treatment of features such as *desc*, *emp\_title* and *title* we would have used embedding transformation to gain information but this goes beyond the limits of

this course. This leaves us with non-redundant variables that we represent as a matrix of correlations to verify our cleaning:



What concerns *na* values, the strategy will be to fit linear regression to impute missing values in the loan dataset by building a model for each column with missing values and predicting these missing values using other columns as predictors.

### 3 Optimization by cross-validation

As said, cross-validation is a technique used to validate the performances of a given model in several folds of our training data, especially it's appropriate when fitting models with hyper-parameters tuning as these have to be chosen to improve performances on unseen test data. Since the unobersed testing data set is supposed to be unavailable during training, we cannot simply select the hyper-parameters that give the best performances on the training data as this can lead to over-fitting, suggesting the cutting of the training set in predetermined number of *k-folds* in order to validate our assumptions about the choice of hyper-parameters. The model is then trained on  $k - 1$  folds and evaluated on the remaining one and this process is repeated  $k$  times, with each fold allocated to the validation set exactly once. The performance of the model is then averaged across the folds to give an estimate of the model's performances on new data. In our case, **grid-search** will be the procedure for testing a wide range of different hyper-parameters. By trying out different hyper-parameter values and evaluating the model's performances using cross-validation, we can find the best combination of hyper-parameters that gives the best average performance on new data. Thus, we re-coded the algorithm to keep tracks of some key metrics explained later. Here it is in pseudo-code where we will only list the parameters that will be subject to optimization:

- *Cost* : This is the penalty parameter for misclassifying training data as a larger value of this hyper-parameter increases the penalty for misclassification, which makes the algorithm more likely to choose a decision boundary that separates the classes more accurately. Thus, a smaller value of cost allows for more misclassifications in the training data, which may result in a decision boundary that is more flexible but may overfit the data. Thus, the cost  $C$  quantify the way our SVM is over-fitting training data.
- *Gamma* : This parameter controls the width of the *Gaussian* kernel and a larger value of  $\gamma$  makes the kernel narrower, which means that each training sample has a smaller influence on the decision boundary. A smaller value of  $\gamma$  makes the *Gaussian* kernel wider, which means that each training example has a larger influence on the decision boundary. This will be used in the *RBF kernel*, but not in the *polynomial* one.
- *Coef0* : This is the independent term in the *polynomial* kernel function as it determines the offset of the decision boundary from the origin. A non-zero value of this hyperparameter shifts the decision boundary away from the origin, while a zero value places it at the origin.

- *Degree* : This is the degree of the *polynomial* kernel as it controls the complexity of the decision boundary by determining the order of the polynomial. Thus, a larger value of *degree* creates a more complex decision boundary that can fit the training data more closely, but may overfit the data. A smaller value of *degree* creates a simpler decision boundary that may not fit the training data as closely, but may generalize better to new data.

Moreover, these parameters are not independent from each others as increasing the *degree* of the *polynomial* kernel can increase the complexity of the decision boundary, which may require increasing the *cost* parameter to prevent overfitting. On the other hand, increasing the *coef0* parameter can shift the decision boundary away from the origin, which may require decreasing the *degree* parameter to prevent overfitting. This is why we will perform a grid search to find the optimal combination of parameters:

---

**Algorithm 1** Cross-validation with grid search for SVM

---

**Require:** A dataset  $\mathcal{D}\{y_i, x_i\}_{i=1}^n$ , a set of hyperparameters  $\mathcal{H} = (C_1, \gamma_1, \dots), (C_2, \gamma_2, \dots), \dots, (C_p, \gamma_p, \dots)$  for SVM, and the number of folds  $k$  for cross-validation.

- 1: Divide  $\mathcal{D}$  into  $k$  equally sized subsets  $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_k$
- 2: Initialize a variable  $best_h$  to  $\emptyset$  and a variable  $best_{score}$  to  $\infty$
- 3: **for all**  $h_i = (C_i, \gamma_i, \dots) \in \mathcal{H}$  **do**
- 4:   Initialize a variable  $score_i$  to 0
- 5:   **for**  $j = 1$  to  $k$  **do**
- 6:     Define the training set  $\mathcal{T}_j$  as  $\mathcal{D}$  excluding subset  $\mathcal{D}_j$
- 7:     Train an SVM model  $f_{\mathcal{T}_j}^{(i)}$  on  $\mathcal{T}_j$  with hyperparameters  $h_i$
- 8:     Define the validation set  $\mathcal{V}_j$  as  $\mathcal{D}_j$
- 9:     Evaluate  $f_{\mathcal{T}_j}^{(i)}$  on  $\mathcal{V}_j$  and record performance metric  $m_{ij}$
- 10:   Add  $m_{ij}$  to  $score_i$
- 11: **end for**
- 12: Calculate the average performance of  $h_i$ :  $average_i = \frac{1}{k} \sum_{j=1}^k m_{ij}$
- 13: **if**  $average_i < best_{score}$  **then**
- 14:   Update  $best_{score}$  to  $average_i$
- 15:   Update  $best_h$  to  $h_i$
- 16: **end if**
- 17: **end for**
- 18: Train a final SVM model  $f$  on  $\mathcal{D}$  with hyperparameters  $best_h$
- 19: **return**  $f$  the final model

---

For other nonlinear cases, where we have a lot more parameters and thus combinations to try, we implemented a learning rate. We initially set a  $1 \times 3$  vector of random values for each hyperparameter to work with. The algorithm updates the initial values for each parameter by selecting the best combination of hyperparameters as a central value. Then we implemented a concentration coefficient based on the current best parameter (either cost or the gamma parameters.) We end up with  $1 \times 3$  vector for each hyperparameter, all we have to do is restart the loop with those new values till the model does not improve his *detection rate* while maintaining a stable *balanced accuracy*. We will explain these metrics in what follows.

## 4 Linear Support Vector Machine

Support Vector Machine is a machine learning algorithm that can be used for classification tasks as *SVM* tries to construct a hyperplane in a high-dimensional input space (i.e in  $d$  dimensions) that separates the data into two area's that will corresponds to the classification of the *Charged\_Off* and *Fully\_Paid* labels of the *loan\_status* feature. The function of the hyperplane is  $h(x)$  as we defined :

$$h(x) = b + w_1x_1 + w_2x_2 + \dots + w_px_p$$

where  $x_1, x_2, \dots, x_p$  are the input variables,  $w_0, w_1, w_2, \dots, w_p$  are the coefficients of the hyperplane, and  $p$  is the number of input variables. All points that lie on the hyperplane have  $h(x) = 0$ .

The *SVM* tries to find the hyperplane that maximizes the margin while still correctly classifying all the training data. The margin is the distance between the hyperplane and the closest points of each class and can be either hard or soft. The hard margin represents the cases where we do not allow any misclassification. The soft margin approach allows for some degree of misclassification, and is ruled by the cost parameter and the slack variables. The slack variables, "ξ", represent in the case of a misclassified point, the distance between this point and the closest margin to his class. The value is equal to zero if the point is well classified. They embody the overlapping of classes.

In the most plausible case, our dataset is non linearly separable, here follows the objective function:

$$\min_{\mathbf{w}, b, \xi} \left( \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i^k \right)$$

The parameter  $C$  is a tuning parameter that is used to determine the penalty for misclassification of the training data points.

## 4.1 Linear separability

Linear separability is determined by the ability to find a hyperplane that can completely separate the two classes of our binary target in the feature space. Mathematically, this can be represented as follows, given a set of training data:

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)$$

where  $\mathbf{x}_i$  is a vector of features for the  $i$ -th data point, and  $y_i$  is its corresponding class label, the problem of linear separability can be formulated as finding a hyperplane  $\mathbf{w}^T \mathbf{x} + b = 0$  such that:

$$\begin{aligned} \mathbf{w}^T \mathbf{x}_i + b &\geq +1, \text{ if } y_i = +1 \\ \mathbf{w}^T \mathbf{x}_i + b &\leq -1, \text{ if } y_i = -1 \end{aligned}$$

where  $\mathbf{w}$  is the weight vector and  $b$  is the bias term.

If these constraints are satisfied, then the training data is linearly separable, and this is equivalent to making the assumption that there exists a margin  $\gamma$  such that:

$$\begin{aligned} \mathbf{w}^T \mathbf{x}_i + b &\geq +\delta^*, \text{ if } y_i = +1 \\ \mathbf{w}^T \mathbf{x}_i + b &\leq -\delta^*, \text{ if } y_i = -1 \end{aligned}$$

where the margin  $\delta^*$  can be calculated as  $\delta^* = \frac{2}{\|\mathbf{w}\|^2}$ , where  $\|\mathbf{w}\|^2$  is the Euclidean norm of  $\mathbf{w}$ . In fact, the *SVM* will choose the canonical hyperplane that yields the maximum margin among all possible separating hyper-planes, as it's the distance between the decision boundary and the closest data points from either class. Thus as we denote  $x_-$  a negative point that lies on the negative side of the hyperplane and  $x_+$  a point that lies on the positive side, we can define the margin to be maximized as:

$$\delta^* = \frac{\mathbf{w}^T \mathbf{x}_+ + b}{\|\mathbf{w}\|^2} - \frac{\mathbf{w}^T \mathbf{x}_- + b}{\|\mathbf{w}\|^2}$$

and this defines what we call *support vectors* as these are lying on the margin for both sides of the hyper-planes, then the equations of such points are respectively equal to 1, -1. Thus the optimization for support vector machine algorithms can be formulated as a constrained optimization problem, where the goal is to minimize the classification error while maximizing the margin between the two classes. Then, the optimization involves finding

a set of *Lagrange multipliers*, also known as *dual variables*, which represent the weights of the training samples in the classification problem. The utilization of *Lagrange multipliers* serves the purpose of identifying the support vectors, which are data points positioned in close proximity to the decision boundary. These support vectors hold paramount importance within the training dataset as they provide crucial information, exerting the highest influence on the optimal hyperplane. This is because the decision boundary is constructed as a linear combination of the support vectors and their corresponding *Lagrange multipliers*, further emphasizing their significance. The explanations above have been retrieved from the slides of the *Statistical Machine Learning and High Dimensional Data Analysis* lectures ([Hafner, 2021](#)) of the *LSBA* faculty at *UCLouvain*.

## 4.2 Results

For the linear support vector machine, the cost  $C$  will be our regularization parameter that controls the trade-off between achieving a low training error and a low testing error as we already said, a smaller value of this metric leads to a wider margin, but more misclassifications, while a larger value of  $C$  leads to a narrower margin, but fewer misclassifications. Moreover, as we want to compare the performances of models with different cost parameters we will introduce some metrics such as the *accuracy* which is the most straightforward metric for evaluating the performance of a classification model because it simply measures the percentage of correctly classified instances out of all instances in the dataset. We will make reference to the *precision* and *recall* percentages which are the proportion of true positive predictions out of all positive predictions made by the model and out of all actual positive instances in the dataset respectively. Finally, to summarize we will compute the *F1 score* that is a weighted average of precision and recall, with equal weighting given to both measures and that provides a single measure of the overall performance of a model. We recall these formula's that will be computed to assess from our classification percentages:

$$\text{Sensitivity} = \frac{TP}{TP + FN} \quad (1)$$

$$\text{Recall} = \frac{TP}{TP + FN} \quad (2)$$

$$\text{Precision} = \frac{TP}{TP + FP} \quad (3)$$

$$\text{F1-score} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (4)$$

$$\text{Balanced Accuracy} = \frac{1}{2} \left( \frac{TP}{TP + FN} + \frac{TN}{TN + FP} \right) \quad (5)$$

$$\text{Detection Rate} = \frac{TP}{TP + FP + FN + TN} \quad (6)$$

The first thing we want to verify : are our dataset linearly separable ? In order to find an answer, we compute the accuracy of our linear SVM given that we choose an hard margin, by choosing a cost value  $C = \infty$ . A perfectly linearly separable dataset would have an accuracy of 100%. Here are our results:

- *bal accuracy*: 75.07%
- *precision*: 29.05%
- *F1-score*: 43.45%
- *margin*:  $6.829 \times 10^{-8}$
- *detec rate*: 12.62%
- *support V*: 3979

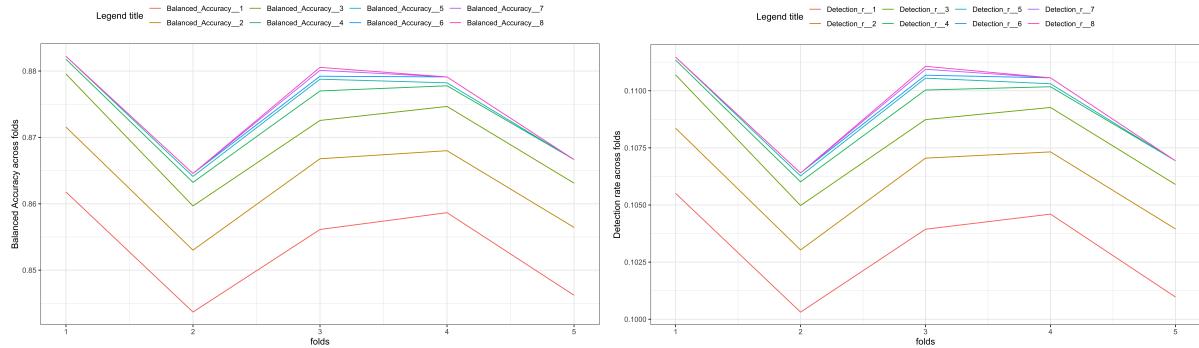
Therefore our dataset is not perfectly linearly separable as we train and test the performances on the whole dataset.

Now that this assumption has been checked, we will use a soft margin and allow misclassification, as explained earlier. We report in the following table the performances of several values of  $C$  obtained by grid search and cross-validation (we have introduced these concepts in the section [Optimization by cross-validation](#)) , with the mean values of each metric across all folds for the linear *SVM*:

Model	$C$	error	dispersion	support_V	Sensitivity	F1_score	Balanced_Accuracy	Precision	Recall	Detection_r
SVM Linear 1	5	0.04283	0.004013	4285	0.7065	0.8279	<b>0.8532</b>	1	0.7066	<b>0.1030</b>
SVM Linear 2	7	0.03953	0.003975	4710	0.7263	0.841	<b>0.8631</b>	1	0.7263	<b>0.1059</b>
SVM Linear 3	9	0.03765	0.004361	4701	0.7398	0.8503	<b>0.8699</b>	1	0.7398	<b>0.1079</b>
SVM Linear 4	11	0.03632	0.004166	4627	0.7466	0.8548	<b>0.8733</b>	1	0.7466	<b>0.1089</b>
SVM Linear 5	13	0.03596	0.004016	4610	0.7480	0.8557	<b>0.8740</b>	1	0.7480	<b>0.1091</b>
SVM Linear 6	15	0.03583	0.003925	4591	0.7487	0.8562	<b>0.8743</b>	1	0.7487	<b>0.1092</b>
SVM Linear 7	17	0.03573	0.003904	4659	0.7490	0.8564	<b>0.8745</b>	1	0.7490	<b>0.1093</b>
SVM Linear 8	19	0.03570	0.003867	4611	0.7492	0.8566	<b>0.8746</b>	1	0.7492	<b>0.1093</b>
SVM Linear 9	25	0.0575	0.0036	4582	0.7492	0.8566	<b>0.8746</b>	1	0.7492	<b>0.1093</b>
SVM Linear 10	50	0.0554	0.0037	4490	0.7496	0.8570	<b>0.8748</b>	1	0.7496	<b>0.1094</b>

Table 2: SVM Linear Model Performance

We recall that the objective of this project, as we are working in actuarial setting, is to build an appropriate classifier that detect contracts that are/will be *fully paid* but as we report the *precision metric* that is always equal to 1, the main focus will be to detect as much *charged off* (i.e true negatives) observations as possible. Thus as we know that our percentage of *charged off* contracts in our test  $\mathcal{T}$  and train set  $\mathcal{D}$  are approximately around 14.6%, we will choose the *SVM* parameters that allows our *detection rate* to be the closest to this value (see the Fire Outbreak Detection System that follows the same approach ([Umoh, Udo, & Nyoho, 2019](#)) in the appendices). For the assessment and as we want to follow the propagation of these key values in our 5 folds, we display the trace plots below:



Thus looking at these plots, we can conclude that both the convergences of the *balanced accuracy* and the *detection rate* are reaching optimal values in the folds for  $C > 15$ . Thus looking at the final table above and more specifically at the error metric we will choose a cost parameter equal to 19 for our final svm in the section [PCA and visualisation](#), as we see that with a highest cost, ther error rate is increasing from 0.0357 to 0.0575.

## 5 Non-linear Support Vector Machine

The main difference between the linear and the nonlinear case is that we project the points in an enlarged feature space rather than working in the input space. Instead of trying to fit a nonlinear model, we map the problem to a new space by doing a nonlinear transformation using,  $v(x_i)$ , a suitably chosen basis function and then use a linear model in this new space.

Afterward, the Kernel trick occurs. The idea is to replace the inner product of the basis function, which comes from the Lagrange dual function resolution, by a kernel function. So instead of mapping the whole data points in the new space and working there, we directly apply the kernel function in the original space and work on the Kernel matrix which is on the input space. Thus we can compute the dot product between any pair of data points in the transformed space without explicitly computing the transformation for all data points. We worked with the Kernel matrix and her pairwise similarities. It results in a decreasing of the computational time of our algorithm. The kernel function can take various form such as the RBF or the polynomial which we are going to explain in

more details in the next sections.

Then, the procedure is the same as before. We use the more general case of the soft margin hyperplane here because we have no guarantee that the problem is linearly separable in this new space. The equation of the hyperplane is given as :

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

where  $K(x_i, x)$  is the kernel function,  $\alpha_i$  is a Lagrangian multiplier  $0 \leq \alpha_i \leq C$  and  $y_i$  is the class belonging (-1, 1).

The objective function stays as follow :

$$\min_{\mathbf{w}, b, \xi} \left( \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i^k \right)$$

## 5.1 Radial Kernel SVM

The RBF kernel corresponds to a transformation that maps the data into an infinite-dimensional space. This means that the transformed data points can have an infinite number of dimensions. It's a popular choice because it can handle non-linearly separable data and has only one hyperparameter to tune. It uses the radial basis function. His distribution graph looks like a Gaussian Distribution curve (see article of ([Czarnecki & Tabor, 2014](#)) in the appendices) and is of the following from :

$$K(x_i, x_j) = e^{-\gamma \|x_i - x_j\|^2}$$

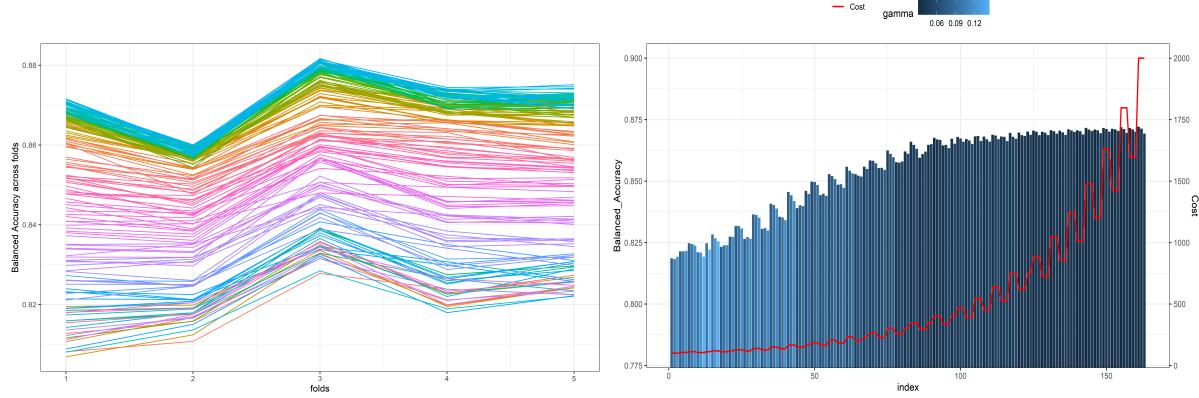
where,  $x_i$  and  $x_j$  are two data points in the feature space,  $\gamma = \frac{1}{2\sigma^2}$  is a hyperparameter that controls the width of the Gaussian density with  $\sigma^2$  the spread parameter of the Gaussian density function, and  $\|x_i - x_j\|^2$  is the squared Euclidean distance between the two data points, their similarity. The RBF kernel can lead to over-fitting if the gamma parameter is set too high. This can cause the SVM to learn noise in the data instead of the underlying patterns. On the other hand, the small gamma produce smoother decision boundaries because they produce smoother feature space mappings. But by setting it too low could results in the incapacity of capturing the underlying patterns in the data. In other words, gamma sets the influence of the points : with a high gamma the nearby points will have high influence, with a low gamma far away points will also be considered to find the decision boundary.

Model	$C$	$\gamma$	$\delta^*$	support_V	F1_score	Balanced_Accuracy	Precision	Recall	Detection_r	L_SV	R_SV	error
SVM Kernel 1	1	0.005	3.3391 e-4	6682	0.6462	<b>0.7387</b>	1	0.4773	<b>0.0696</b>	3680	3001	0.07624
SVM Kernel 2	1	1.5	2.4077e-04	19957	0.3092	<b>0.5915</b>	0.9569	0.1845	<b>0.0269</b>	15608	4353	0.1202
SVM Kernel 3	10	0.005	4.5204e-05	5880	0.7411	<b>0.7944</b>	0.9997	0.5888	<b>0.0859</b>	3796	2173	0.0600
SVM Kernel 4	10	0.3	5.7754e-05	7569	0.7481	<b>0.7921</b>	0.8899	0.5968	<b>0.0870</b>	5068	2526	0.0696
SVM Kernel 5	17.5	0.0050	2.7470e-05	5745	0.7507	<b>0.8005</b>	1.0000	0.6010	<b>0.0877</b>	3539	2100	0.0582
SVM Kernel 6	17.5	0.1000	3.2864e-05	5094	0.7645	<b>0.8167</b>	0.9513	0.6391	<b>0.0932</b>	3138	1986	0.0574
SVM Kernel 7	17.5	0.5	4.2960e-05	9492	0.6446	<b>0.7613</b>	0.7810	0.5490	<b>0.0801</b>	6671	2930	0.0883
SVM Kernel 8	25	0.1	2.4049e-05	5003	0.7690	<b>0.8215</b>	0.9420	0.6497	<b>0.0948</b>	3115	1936	0.0569
SVM Kernel 9	107.50	0.0614	6.2517e-06	4304	0.8144	<b>0.8529</b>	0.9513	0.7119	<b>0.1038</b>	2684	1683	0.0473
SVM Kernel 10	222.90	0.0393	3.1328e-06	3894	0.8360	<b>0.8660</b>	0.9672	0.7363	<b>0.1074</b>	2368	1559	0.0421
SVM Kernel 11	1150.13	0.0377	7.4137e-07	3554	0.8365	<b>0.8710</b>	0.9472	0.7491	<b>0.1093</b>	2231	1353	0.0427
SVM Kernel 12	1656.18	0.0377	5.3594e-07	3521	0.8356	<b>0.8711</b>	0.9434	0.7500	<b>0.1094</b>	2228	1313	0.0430
SVM Kernel 13	1987.42	0.0302	4.3896e-07	3387	0.8394	<b>0.8721</b>	0.9522	0.7507	<b>0.1095</b>	2136	1300	0.0419

Table 3: SVM RBF Kernel Model Performance using hinge loss

As you can see in on plot below and on our Table 3, our variable of interest, the detection rate, with SVM RBF

Kernel converges to 10, 95%. To recall it, the target is a detection rate of 14, 6%. We are quiet happy with this result since we also have a good balanced accuracy, at 87, 21% and a small error. The graphs below underline the convergence. As you can see, as the iterations goes by, we are losing precision while converging, but it is not a problem since both the balanced accuracy and the detection rate increase and simultaneously the number of support vectors drops. You will find the results of the RBF kernel using hinge loss in appendices section 8.1.



## 5.2 Polynomial Kernel SVM

As for the *Radial Kernel SVM*, the kernel trick will be helpful to build the polynomial one as once again, the polynomial kernel will be used to transform the input data from its original space to a higher-dimensional feature space, where it may be easier to separate the 2 classes. Then, the *polynomial kernel* is defined as:

$$K_d(x_i, x_j) = (\gamma x_i^T x_j + r)^d = \sum_{k=1}^d \binom{d}{k} r^{d-k} (\gamma x_i^T x_j)^k$$

where  $x_i$  and  $x_j$  are input feature vectors,  $d$  is the degree of the polynomial,  $\gamma$  is a scaling factor, and  $r$  is a coefficient and thus each pair of input vectors is represented as a polynomial of degree  $d$  with higher degrees allowing for more complex interactions between features.

This kind of kernel is often chosen to improve the classification performance on nonlinear datasets, that is non linearly separable, what we have concluded for our *loan* database. However, it's essential to note that using a *polynomial kernel* with a high degree can lead to over-fitting on the training data, as the decision boundary may become too complex and fit the noise in the data. Therefore, once again the choice of the hyperparameters should be carefully selected through grid-search and cross-validation to avoid this behaviour and maximize generalization performance. Even though we have performed a grid search with a rate covering a wide range of parameters we will only report in the table the most conclusive results. You will find these results in the appendices section [Polynomial Kernel results using hinge loss](#).

## 6 PCA and visualisation

The goal of PCA is to find a new set of features, called principal components  $z$ , which are linear combinations of the original features and that capture the maximum amount of variance in the data. The first principal component, denoted by  $z_1$ , is given by the linear combination of the variables  $x_1, x_2, \dots, x_p$ :

$$X = \begin{pmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,p} \\ x_{2,1} & x_{2,2} & \dots & x_{2,p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & \dots & x_{n,p} \end{pmatrix} \quad z_1 = \phi_{1,1}x_1 + \phi_{2,1}x_2 + \dots + \phi_{p,1}x_p$$

where  $\phi_{1,1}, \phi_{2,1}, \dots, \phi_{p,1}$  are the coefficients or loading's of the first principal component and are chosen such that the variance of  $z_1$  is maximized subject to the constraint that  $\sum_{i=1}^p \phi_{i1}^2 = 1$ .

Mathematically, PCA can be formulated as if we let  $X$  be an  $n \times p$  matrix containing  $n$  observations of  $p$  variables then by first centering the data and by subtracting the mean of each variable from each observation, then we can compute the covariance matrix  $S$  of the centered data:

$$S = \frac{1}{n-1} X^T X$$

Next, we can compute the eigenvectors and eigenvalues of  $S$  as an eigenvector of  $S$  is a non-zero vector  $v$  that satisfies the equation:

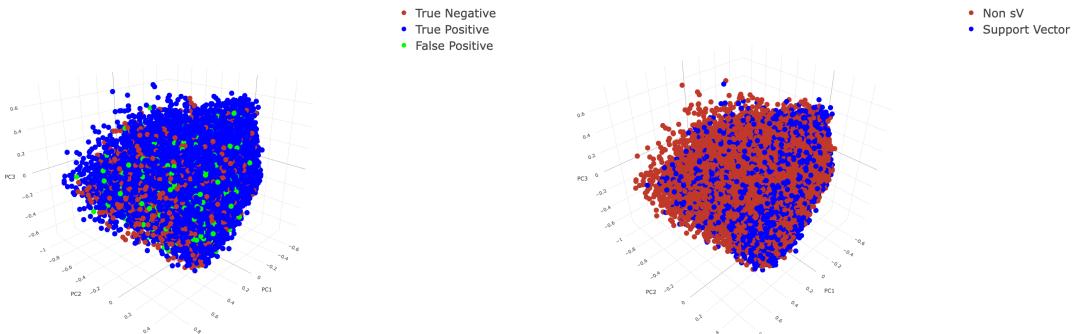
$$Sv = \lambda v$$

where  $\lambda$  is a scalar, called the eigenvalue that corresponds to the amount of variance explained by each eigenvector such that their respective eigenvectors of  $S$  form an orthonormal basis for the  $d$ -dimensional space. Then, we can use the eigenvectors of  $S$  to transform the original data into a new coordinate system as the first principal component is the eigenvector with the largest eigenvalue and thus that captures the biggest amount of variance. As we want to represent our contracts in a 3D space we will thus choose the first 3 principal components and then the transformed data matrix  $Z$  is obtained by multiplying the centered data matrix  $X$  by the matrix of 3 first eigenvectors  $V$ :

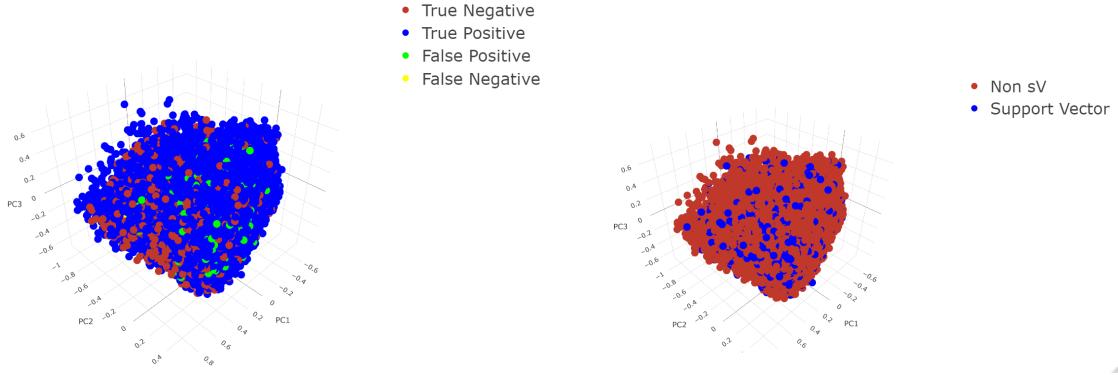
$$Z = XV$$

Thus in order to close this analysis we will display the representation of our data points with PCA and display colors corresponding to *true positive*, *true negative* that are well classified and *false positive*, *false negative* that are misclassified. On another plot we color the support vectors to assess for the distribution of points associated to the margin that are the closest to the decision boundary. But we want to add that as it has been done in the paper of (Kempfert, Wang, Chen, & Wong, 2019), it has not been considered to perform other dimensionality reduction such as the *Supervised kernel PCA* prior to our *Support vector machine* classification that may have improved our results.

## 6.1 PCA visualisation for linear Kernel ( $C = 19$ )



## 6.2 PCA visualisation for radial Kernel ( $C = 1987.42, \gamma = 0.0301$ )



## 7 Conclusion

As it's time to finish our journey through the application of support vector machines to our contract database we will choose the one that gives the highest performances and interpret our results at best. Our results showed that support vector machines are highly effective in predicting loan defaults, achieving a high *accuracy rate* and demonstrating the ability to handle complex and nonlinear relationships between the input features. This project highlights the potential of support vector machines as a valuable tool in the financial industry for risk assessment. It offers lenders an effective means of identifying high-risk loan applicants and reducing the risk of default. Our main objective was to improve the detection rate while maintaining a stable *balanced accuracy*, a high *precision* and a low *error*. We have to add that a anomaly analysis has been performed, it's reasonable to expect not perfect results since our dataset contains 4.6% of outliers and that the proportion of anomalies in the *Charged Off* level of the *loan status* target can reach 15% for some features.

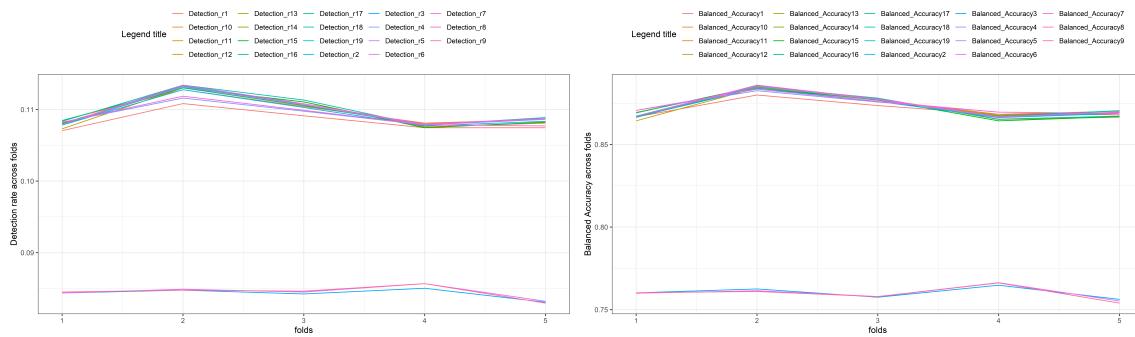
The best performing model is the support vector machine with polynomial kernel 8 ( $C = 15, d = 3, \gamma = 1, r = 1$ ) that gives performances of 87.81% and 11.13% for the *balanced accuracy* and *detection rate* respectively. We keep track of the error rate to make sure that this *svm* configuration is the one that perform better than others and as it's the lowest one we get across all kernels with *error* = 0.0336 and *dispersion* = 0.00393. The second best performing model is with RBF kernel 13 with a quadratic loss function :( $C = 1987.42, \gamma = 0.0242$ ) with a *balanced accuracy* of 87.39% and a *detection rate* of 10.98% while have a small error of 4.08%. Nonetheless every model tried performed almost similarly at optimum, even, at our surprise, the linear one. As it was asked to analyse the results for one of them, we chose to visualize the radial kernel one through PCA as the convergence across the folds was smoother. We encountered problem with the running time and we concluded, as neither expert coder nor yet data scientist, that this could be problematic in the industry. We know that the polynomial kernel with an higher cost and a quadratic loss function would perform better but given our available computational power, we decided to stop there due to lack time.

## 8 Appendices

### 8.1 Graph and results about Kernel RBF using a Quadratic loss

Model	$C$	$\gamma$	$\delta^*$	support_V	F1_score	Balanced_Accuracy	Precision	Recall	Detection_r	L_SV	R_SV	error
SVM Kernel 1	1150.13	0.005	5.8541e-07	3466	0.8515	<b>0.8713</b>	0.9971	0.7430	<b>0.1083</b>	1938	1459	0.0378
SVM Kernel 2	1150.13	0.0302	7.1611e-07	3442	0.8427	<b>0.8724</b>	0.9617	0.7500	<b>0.1094</b>	2081	1327	0.0408
SVM Kernel 3	1150.13	0.5	5.3951e-06	8210	0.6038	<b>0.7602</b>	0.6323	0.5779	<b>0.0843</b>	5720	2601	0.1106
SVM Kernel 4	1656.18	0.005	4.2305e-07	3302	0.8541	<b>0.8735</b>	0.9964	0.7475	<b>0.1090</b>	1814	1406	0.0372
SVM Kernel 5	1656.18	0.0302	5.1732e-07	3406	0.8431	<b>0.8732</b>	0.9598	0.7517	<b>0.1097</b>	2067	1292	0.0408
SVM Kernel 6	1656.18	0.5	4.9495e-06	8158	0.6018	<b>0.7602</b>	0.6259	0.5795	<b>0.0845</b>	5664	2600	0.1119
SVM Kernel 7	1987.42	0.005	3.5984e-07	3246	0.8543	<b>0.8737</b>	0.9960	0.7480	<b>0.1091</b>	1788	1375	0.0372
SVM Kernel 8	1987.42	0.0302	4.3923e-07	3379	0.8425	<b>0.8734</b>	0.9568	0.7526	<b>0.1098</b>	2087	1279	0.0410
SVM Kernel 9	1987.42	0.5	4.7890e-06	8143	0.6008	<b>0.7598</b>	0.6240	0.5793	<b>0.0845</b>	5653	2601	0.1123
SVM Kernel 10	1589.94	0.0242	5.1987e-07	3331	0.8464	<b>0.8738</b>	0.9682	0.7519	<b>0.1097</b>	2010	1293	0.0398
SVM Kernel 11	1589.94	0.0302	5.3656e-07	3407	0.8430	<b>0.8731</b>	0.9600	0.7516	<b>0.1096</b>	2063	1297	0.0408
SVM Kernel 12	1589.94	0.0362	5.5236e-07	3489	0.8380	<b>0.8717</b>	0.9494	0.7501	<b>0.1094</b>	2161	1301	0.0423
SVM Kernel 13	1987.42	0.0242	4.2535e-07	3286	0.8457	<b>0.8739</b>	0.9655	0.7524	<b>0.1098</b>	1998	1292	0.0408
SVM Kernel 14	1987.42	0.0302	4.3923e-07	3379	0.8425	<b>0.8734</b>	0.9568	0.7526	<b>0.1098</b>	2087	1279	0.0410
SVM Kernel 15	1987.42	0.0362	4.5210e-07	3488	0.8370	<b>0.8719</b>	0.9450	0.7212	<b>0.1096</b>	2168	1293	0.0427
SVM Kernel 16	2384.90	0.0242	3.6052e-07	3264	0.8453	<b>0.8738</b>	0.9649	0.7523	<b>0.1097</b>	1991	1258	0.0401
SVM Kernel 17	2384.90	0.0302	3.7254e-07	3357	0.8410	<b>0.8729</b>	0.9540	0.7521	<b>0.1097</b>	2067	1267	0.0414
SVM Kernel 18	2384.90	0.0362	3.8355e-07	3444	0.8374	<b>0.8726</b>	0.9434	0.7530	<b>0.1098</b>	2165	1280	0.0426
SVM Kernel 19	1907.92	0.0290	4.5290e-07	3361	0.8437	<b>0.8735</b>	0.9605	0.7523	<b>0.1097</b>	2052	1281	0.0406

Table 4: SVM RBF Kernel Model Performance



### 8.2 Polynomial Kernel results using hinge loss

Model	$C$	$d$	$\gamma$	$r$	$\delta^*$	support_V	F1_score	Balanced_Accuracy	Precision	Recall	Detection_r	error	dispersion
SVM Poly 1	7	2	1	1	1.032e-04	11992	0.8491	<b>0.8689</b>	1	0.7378	<b>0.1076</b>	0.03765	0.004361
SVM Poly 2	7	3	10	1	4.775e-03	11673	0.4845	<b>0.7778</b>	0.343	0.8421	<b>0.1228</b>	0.0931	0.00594
SVM Poly 3	7	4	10	1	1.209e-01	6507	0.4601	<b>0.7484</b>	0.3321	0.8171	<b>0.1192</b>	0.1047	0.00754
SVM Poly 4	11	1	1	1	6.132e-05	4644	0.8546	<b>0.8731</b>	1	0.7462	<b>0.1088</b>	0.0401	0.00469
SVM Poly 5	15	1	10	1	5.61e-05	4387	0.8568	<b>0.8748</b>	1	0.749	<b>0.1093</b>	0.03712	0.00456
SVM Poly 6	15	1	10	0	5.617e-05	4476	0.8568	<b>0.8748</b>	1	0.7496	<b>0.1093</b>	0.03668	0.00436
SVM Poly 7	15	3	0.5	1	5.4088e-05	9605	0.854	<b>0.8748</b>	0.9895	0.7510	<b>0.1096</b>	0.03564	0.00427
SVM Poly 8	15	3	1	1	6.349e-05	12776	0.8466	<b>0.8781</b>	0.9542	0.7628	<b>0.1113</b>	0.0336	0.00393

Table 5: SVM nonlinear Model Performance1

### 8.3 R Code for linear kernel grid search & CV

```

1 folds <- createFolds(data_clustering$loan_status, k = 5)
2
3
4 # Define lambda and alpha values to search over
5 Cost <- c(5,7,9,11,13,15,17,19)
6 # Initialize matrix to store cross-validation results
7 cv_results_svm <- matrix(NA, nrow = length(Cost), ncol = 8)
8 colnames(cv_results_svm) <- c("Cost", "support_V", "Sensitivity", "F1_score", "Balanced_Accuracy", "Precision", "Recall", "Detection_r")

```

```

9  row_idx <- 1
10 set.seed(87031800)
11 results_list <- list()
12
13 # Loop over lambda and alpha values
14 for (j in 1:length(Cost)) {
15   # Initialize vector to store Dev.Bern results for this combination of lambda and alpha
16   support_V <- rep(NA, 5)
17   Sensitivity <- rep(NA, 5)
18   F1_score <- rep(NA, 5)
19   Balanced_Accuracy <- rep(NA, 5)
20   Precision <- rep(NA, 5)
21   Recall <- rep(NA, 5)
22   Detection_r <- rep(NA, 5)
23
24   # Loop over folds
25   nfolds <- 0
26   for (fold in 1:5) {
27     # Split data into training and test sets for this fold
28     train_idx <- setdiff(seq_len(nrow(data_clustering)), folds[[fold]])
29     test_idx <- folds[[fold]]
30     x_train <- data_clustering[train_idx, ]
31     y_train <- data_clustering$loan_status[train_idx]
32     x_test <- data_clustering[test_idx, ]
33     y_test <- data_clustering$loan_status[test_idx]
34     print("go")
35     print(Cost[j])
36     # Fit model with elastic net regularization
37     final_model <- svm(loan_status ~ ., data = x_train, kernel = "linear", cost = Cost[j])
38     print("fin")
39     # Compute odds
40
41     y_pred <- predict(final_model, newdata = x_test)
42     confu <- confusionMatrix(data = y_pred, y_test)
43     print(confu)
44     sensi <- confu$byClass[["Sensitivity"]]
45     f1 <- confu$byClass[["F1"]]
46     balanced <- confu$byClass[["Balanced Accuracy"]]
47     prec <- confu$byClass[["Precision"]]
48     rec <- confu$byClass[["Recall"]]
49     det <- confu$byClass[["Detection Rate"]]
50
51     # Store results for this fold
52     support_V[fold] <- final_model$tot.nSV
53     print(support_V)
54     Sensitivity[fold] <- sensi
55     print(Sensitivity)
56     F1_score[fold] <- f1
57     print(F1_score)
58     Balanced_Accuracy[fold] <- balanced
59     print(Balanced_Accuracy)
60     Precision[fold] <- prec
61     print(Precision)
62     Recall[fold] <- rec
63     print(Recall)
64     Detection_r[fold] <- det
65     nfolds <- nfolds + 1
66   }
67   results_list <- cbind(results_list, list(support_V, Sensitivity, F1_score, Balanced_Accuracy, Precision, Recall, Detection_r))

```

```

68 print(results_list)
69 # Compute mean Dev.Bern loss across folds for this combination of lambda and alpha
70 mean_support_V <- mean(support_V)
71 mean_Sensitivity <- mean(Sensitivity)
72 mean_F1_score <- mean(F1_score)
73 mean_Balanced_Accuracy <- mean(Balanced_Accuracy)
74 mean_Precision <- mean(Precision)
75 mean_Recall <- mean(Recall)
76 mean_Detection_r <- mean(Detection_r)
77
78 # Store results in cv_results matrix
79 cv_results_svm[row_idx, 1] <- Cost[j]
80 cv_results_svm[row_idx, 2:8] <- c(mean_support_V, mean_Sensitivity, mean_F1_score,
  mean_Balanced_Accuracy, mean_Precision, mean_Recall, mean_Detection_r)
81 print(cv_results_svm)
82 # Increment row index
83 row_idx <- row_idx + 1
84 }

```

## 8.4 R Code for RBF Kernel grid search & CV

```

1 folds <- createFolds(data_clustering$loan_status, k = 5)
2
3 # Define lambda and alpha values to search over
4 cost <- c(1150.13, 1656.18, 1987.42)
5 gamma = c(0.005, 0.0302, 0.5)
6 previous_error = Inf
7 previous_acc = 0
8 previous_detection_r = Inf
9
10
11 decay = 0
12 conver = FALSE
13 # Initialize matrix to store cross-validation results
14 cv_results <- matrix(NA, nrow = 9, ncol = 13)
15 colnames(cv_results) <- c("Cost", "gamma", "margin",
16                           "support_V", "Sensitivity", "F1_score",
17                           "Balanced_Accuracy", "Precision", "Recall",
18                           "Detection_r", "L_SV", "R_SV",
19                           "error")
20
21 results_list <- list()
22 full_results = list()
23 #convergebce loop
24 while (conver == FALSE) {
25   print(cost)
26   print(gamma)
27   set.seed(696460)
28
29   row_idx <- 1
30 # Loop over gamma and alpha values
31   for (c in 1:length(cost)) {
32     for (gam in 1:length(gamma)){
33       # Initialize vector to store Dev.Bern results for this combination of lambda and
34       # alpha
35       support_V <- rep(NA, 5)
36       Sensitivity <- rep(NA, 5)
37       F1_score <- rep(NA, 5)
38       Balanced_Accuracy <- rep(NA, 5)
39       Precision <- rep(NA, 5)

```

```

39   Recall <- rep(NA, 5)
40   Detection_r <- rep(NA, 5)
41
42   margin = rep(NA,5)
43   L_SV = rep(NA,5)
44   R_SV = rep(NA,5)
45   error = rep(NA,5)
46
47   # Loop over folds
48   nfolds <- 0
49   for (fold in 1:5) {
50     # Split data into training and test sets for this fold
51     train_idx <- setdiff(seq_len(nrow(data_clustering)), folds[[fold]])
52     test_idx <- folds[[fold]]
53     x_train <- data_clustering[train_idx, ]
54     y_train <- data_clustering$loan_status[train_idx]
55     x_test <- data_clustering[test_idx, ]
56     y_test <- data_clustering$loan_status[test_idx]
57
58     # Fit model with elastic net regularization
59     print(c("go in",fold))
60
61     final_model <- svm(loan_status ~ ., data = x_train, kernel = "radial", cost = cost
62 [c], gamma = gamma[gam],
63           type = "C-classification", loss = "squared-hinge")
64     predictions = predict(final_model, x_test)
65     error[fold] = mean(predictions != y_test)
66
67     print("out")
68
69     # Compute odds
70
71     y_pred <- predict(final_model, newdata = x_test)
72     confu <- confusionMatrix(data = y_pred, y_test)
73
74     sensi <- confu$byClass[["Sensitivity"]]
75     f1 <- confu$byClass[["F1"]]
76     balanced <- confu$byClass[["Balanced Accuracy"]]
77     prec <- confu$byClass[["Precision"]]
78     rec <- confu$byClass[["Recall"]]
79     det <- confu$byClass[["Detection Rate"]]
80
81
82     # Store results for this fold
83     support_V[fold] <- final_model$tot.nSV
84     Sensitivity[fold] <- sensi
85     F1_score[fold] <- f1
86     Balanced_Accuracy[fold] <- balanced
87     Precision[fold] <- prec
88     Recall[fold] <- rec
89     Detection_r[fold] <- det
90
91     margin[fold] = 2/norm(final_model$coefs)
92     L_SV = final_model$nSV[1]
93     R_SV = final_model$nSV[2]
94     nfolds <- nfolds + 1
95   }
96   #end fold loop
97   print("out of fold")

```

```

98 results_list = cbind(results_list, list(support_V, Sensitivity, F1_score, Balanced_
99 Accuracy, Precision, Recall, Detection_r, margin, error, L_SV, R_SV))
100
101 print(results_list)
102 # Compute mean Dev.Bern loss across folds for this combination of lambda and alpha
103 mean_support_V <- mean(support_V)
104 mean_Sensitivity <- mean(Sensitivity)
105 mean_F1_score <- mean(F1_score)
106 mean_Balanced_Accuracy <- mean(Balanced_Accuracy)
107 mean_Precision <- mean(Precision)
108 mean_Recall <- mean(Recall)
109 mean_Detection_r <- mean(Detection_r)
110
111 mean_margin = mean(margin)
112 mean_L_SV = mean(L_SV)
113 mean_R_SV = mean(R_SV)
114 mean_error = mean(error)
115
116 # Store results in cv_results matrix
117 cv_results[row_idx, 1:2] <- c(cost[c], gamma[gam])
118 cv_results[row_idx, 3:13] <- c(mean_margin, mean_support_V, mean_Sensitivity,
119                               mean_F1_score, mean_Balanced_Accuracy, mean_Precision
120 ,
121                               mean_Recall, mean_Detection_r, mean_L_SV,
122                               mean_R_SV, mean_error)
123 print(cv_results)
124 # Increment row index
125 row_idx <- row_idx + 1
126 }
127 #####end double if c, gam
128
129 print("end double")
130 #to keep the previous retults
131 full_results = cbind(full_results, cv_results)
132 print(full_results)
133
134 #finding the optimal parameter based on the smallest detection_rate
135 #we find the index and take the minimal difference
136 opt_detec = abs(0.1406 - cv_results[,10])
137 opt_index = which.min(opt_detec)
138 best_detection_r = opt_detec[opt_index]
139
140 print(c("opt_detec"))
141 print(c("opt_index", opt_index))
142 #we assign the best hyperparameters with the opt index
143 best_cost = cv_results[opt_index,1]
144 best_gamma = cv_results[opt_index,2]
145 best_acc = cv_results[opt_index,7]
146
147
148 #update the new cost and gamma
149 # care implementation of the decay hasn't been tested
150 cost = c(best_cost - (best_cost /5), best_cost, best_cost + (best_cost/5))
151 gamma = c(best_gamma - (best_gamma /5), best_gamma, best_gamma + (best_gamma/5))
152
153 print(c(cost, gamma))
154
155 #test of convergence based on the detection_rate
156 if (best_detection_r >= previous_detection_r) {

```

```

157     print("it converged")
158     conver <- TRUE
159   }
160   previous_detection_r = best_detection_r
161
162   print(c(best_detection_r, previous_detection_r))
163   print(full_results)
164 }

```

## 8.5 R Code for polynomial Kernel grid search & CV

```

1   folds <- createFolds(data_clustering$loan_status, k = 5)
2
3 # Define lambda and alpha values to search over
4 Cost <- c(7,11,15,19,25, 50,100)
5 Degree <- c(1,2,3,4)
6 gamma <- c(0.01,0.1,0.5,1, 10)
7 Coefs <- c(0, 1)
8 # Initialize matrix to store cross-validation results
9 cv_results_svm <- matrix(NA, nrow = length(Cost)*length(Degree)*length(gamma)*length(
10   Coefs), ncol = 14)
11 colnames(cv_results_svm) <- c("Cost", "Degree", "Gamma", "Coef0", "margin", "support_V", "
12   Sensitivity", "F1_score", "Balanced_Accuracy", "Precision", "Recall", "Detection_r",
13   "Lsv", "Rsv")
14 row_idx <- 1
15 results_list <- list()
16
17 # Loop over lambda and alpha values
18 for (i in 1:length(Cost)) {
19   for (j in 1:length(Degree)) {
20     for (z in 1:length(gamma)){
21       for (a in 1:length(Coefs)) {
22         set.seed(87031800)
23         # Initialize vector to store Dev.Bern results for this combination of lambda and
24         # alpha
25         support_V <- rep(NA, 5)
26         Sensitivity <- rep(NA,5)
27         F1_score <- rep(NA, 5)
28         Balanced_Accuracy <- rep(NA,5)
29         Precision <- rep(NA, 5)
30         Recall <- rep(NA, 5)
31         Detection_r <- rep(NA, 5)
32         margin <- rep(NA,5)
33         Lsv <- rep(NA,5)
34         Rsv <- rep(NA, 5)
35
36         # Loop over folds
37         nfolds <- 0
38         for (fold in 1:5) {
39           # Split data into training and test sets for this fold
40           train_idx <- setdiff(seq_len(nrow(data_clustering)), folds[[fold]])
41           test_idx <- folds[[fold]]
42           x_train <- data_clustering[train_idx, ]
43           y_train <- data_clustering$loan_status[train_idx]
44           x_test <- data_clustering[test_idx, ]
45           y_test <- data_clustering$loan_status[test_idx]
46           print("go")
47           print(c(Cost[i], Degree[j], gamma[z], Coefs[a]))
48           # Fit model with elastic net regularization

```

```

final_model <- svm(loan_status ~ ., data = x_train, kernel = "polynomial",
cost = Cost[i], degree = Degree[j], gamma = gamma[z], coef0 = Coefs[a])
print("fin")
# Compute odds

y_pred <- predict(final_model, newdata = x_test)
confu <- confusionMatrix(data = y_pred, y_test)
print(confu)
sensi <- confu$byClass[["Sensitivity"]]
f1 <- confu$byClass[["F1"]]
balanced <- confu$byClass[["Balanced Accuracy"]]
prec <- confu$byClass[["Precision"]]
rec <- confu$byClass[["Recall"]]
det <- confu$byClass[["Detection Rate"]]

# Store results for this fold
margin[fold] <- 2/norm(final_model$coefs)
Lsv[fold] <- final_model$nSV[1]
Rsv[fold] <- final_model$nSV[2]
support_V[fold] <- final_model$tot.nSV
print(support_V)
Sensitivity[fold] <- sensi
print(Sensitivity)
F1_score[fold] <- f1
print(F1_score)
Balanced_Accuracy[fold] <- balanced
print(Balanced_Accuracy)
Precision[fold] <- prec
print(Precision)
Recall[fold] <- rec
print(Recall)
Detection_r[fold] <- det
nfolds <- nfolds + 1
}

results_list <- cbind(results_list, list(margin, support_V, Sensitivity, F1_
score, Balanced_Accuracy, Precision, Recall, Detection_r, Lsv, Rsv))
print(results_list)
# Compute mean Dev.Bern loss across folds for this combination of lambda and
alpha
mean_support_V <- mean(support_V)
mean_Sensitivity <- mean(Sensitivity)
mean_F1_score <- mean(F1_score)
mean_Balanced_Accuracy <- mean(Balanced_Accuracy)
mean_Precision <- mean(Precision)
mean_Recall <- mean(Recall)
mean_Detection_r <- mean(Detection_r)
mean_margin <- mean(margin)
mean_Lsv <- mean(Lsv)
mean_Rsv <- mean(Rsv)

# Store results in cv_results matrix
cv_results_svm[row_idx, 1:4] <- c(Cost[i], Degree[j], gamma[z], Coefs[a])
cv_results_svm[row_idx, 5:14] <- c(mean_margin, mean_support_V, mean_Sensitivity
, mean_F1_score, mean_Balanced_Accuracy, mean_Precision, mean_Recall, mean_Detection_
r, mean_Lsv, mean_Rsv)
print(cv_results_svm)
# Increment row index
row_idx <- row_idx + 1
}
}
}

```

101 }

## 8.6 R Code for PCA projection

```

1 # Train an SVM on the transformed data
2 svm <- svm(loan_status~ ., data = data_PCA, kernel = "radial", cost = 1987.42, gamma =
0.0302)
3
4 # Predict SVM output for each data point
5 pred_labels <- predict(svm, data_PCA)
6 confusionMatrix(pred_labels, data_PCA$loan_status)
7 # Create a color vector for the data points based on the predicted labels
8 colors <- rep(3, length(pred_labels)) # initialize colors vector with 3s
9 for(i in 1:length(pred_labels)) {
10   if(pred_labels[i] == "Charged Off" && y[i] == "Charged Off") {
11     colors[i] <- 1
12   } else if(pred_labels[i] == "Fully Paid" && y[i] == "Fully Paid") {
13     colors[i] <- 2
14   } else if (pred_labels[i]== "Charged Off" && y[i] == "Fully Paid" ) {
15     colors[i] = 4
16   }
17 }
18 to_comp <- table(colors)
19 to_comp
20
21
22 # PCA
23 pca <- prcomp(x)
24 pca.x <- pca$x[,1:3]
25 library(plot3D)
26 library(rgl)
27 library(plotly)
28 summary(pca)
29 library(factoextra)
30
31 png(filename = "contributions66.png", width = 10, height = 6, units = 'in', res=900) ##
      save the plot in files
32 fviz_eig(pca)
33 dev.off()
34
35 fviz_pca_var(pca,col.var = "contrib", gradient.cols = c("#00AFBB", "#E7B800", "#FC4E07"))
      ,repel = TRUE)
36 pca
37
38 plot3d(pca$rotation[,1:3],
39         texts=rownames(pca$rotation),
40         col="red", cex=0.8
41
42 text3d(pca$rotation[,1:3],
43         texts=rownames(pca$rotation),
44         col="red", cex=0.8
45
46 coords <- NULL
47 for (i in 1:nrow(pca$rotation)) {
48   coords <- rbind(coords,
49                     rbind(c(0,0,0),
50                           pca$rotation[i,1:3]))
51 }
52 lines3d(coords,
53

```

```

54     col="red",
55     lwd=1)
56
57 dev.off()
58 var <- get_pca_var(pca)
59 ind = get_pca_ind(pca)
60
61 mean(ind$cos2[,1]) # 0.3439
62 mean(ind$cos2[,2]) # 0.1778
63 mean(ind$cos2[,3]) #0.18293
64 mean(ind$cos2)
65 ind$coord
66 var$cos2
67
68
69 library("corrplot")
70 corrplot(var$cos2, is.corr=FALSE)
71 fviz_cos2(pca, choice = "var", axes = 1:3)
72 ## Plot for true positive etc...
73
74 pca.fin <- as.data.frame(pca.x)
75 pca.fin <- cbind(pca.fin, colors)
76 pca.fin$colors <- factor(pca.fin$colors, levels = c("1", "2", "3","4"), labels = c("True
    Negative", "True Positive", "False Positive","False Negative"))
77
78 fig1 <- plot_ly(pca.fin, x = ~PC1, y = ~PC2, z = ~PC3, color = ~colors, colors = c('#
    BF382A', 'blue', 'green',"yellow"))
79 fig1 <- fig1 %>% add_markers()
80 fig1 <- fig1 %>% layout(scene = list(xaxis = list(title = 'PC1'),
81                             yaxis = list(title = 'PC2'),
82                             zaxis = list(title = 'PC3')),
83                             legend = list(font = list(size = 30),
84                                           itemsizing = "constant",
85                                           itemwidth = 50,
86                                           itemheight = 20))
87
88 fig1
89
90
91 ## Plot for support vectors
92
93 supportvectors <- svm$index
94
95 pca.fin2 <- as.data.frame(pca.x)
96 colors <- rep(1, length(pca.fin2))
97 pca.fin2 <- cbind(pca.fin2, colors)
98 sv_indices <- match(supportvectors, rownames(pca.fin))
99
100 pca.fin2$colors[sv_indices] <- 2
101 pca.fin2$colors <- factor(pca.fin2$colors, levels = c("1", "2"), labels = c("Non sv", "
    Support Vector"))
102
103 fig2 <- plot_ly(pca.fin2, x = ~PC1, y = ~PC2, z = ~PC3, color = ~colors, colors = c('#
    BF382A', 'blue'))
104 fig2 <- fig2%>% add_markers()
105 fig2 <- fig2 %>% layout(scene = list(xaxis = list(title = 'PC1'),
106                             yaxis = list(title = 'PC2'),
107                             zaxis = list(title = 'PC3')),
108                             legend = list(font = list(size = 30),
109                                           itemsizing = "constant",
110                                           itemwidth = 50,

```

```
111     itemheight = 20))  
112 fig2
```

## References

- Czarnecki, W. M., & Tabor, J. (2014). Cluster based RBF kernel for support vector machines. *CoRR*, *abs/1408.2869*. Retrieved from <http://arxiv.org/abs/1408.2869>
- Hafner, C. (2021). Statistical machine learning and high dimensional data analysis. *Slides K-means & Support Vector Machines*, 22(43), 1–40. Retrieved from <https://moodle.uclouvain.be/course/view.php?id=6116>
- Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The elements of statistical learning: data mining, inference and prediction* (Second ed.). Springer. Retrieved from <http://www-stat.stanford.edu/~tibs/ElemStatLearn/>
- Kempfert, K. C., Wang, Y., Chen, C., & Wong, S. W. K. (2019). *A comparison study on nonlinear dimension reduction methods with kernel variations: Visualization, optimization and classification*. Retrieved from <https://arxiv.org/abs/1910.02114>
- Umoh, U. A., Udo, E., & Nyoho, E. E. (2019). Support vector machine-based fire outbreak detection system. *CoRR*, *abs/1906.05655*. Retrieved from <http://arxiv.org/abs/1906.05655>