

**Prediction of Survival in Hepatocellular Carcinoma Patients over a Year**  
**How effective are Support Vector Machine and Neural Networks in Combination**  
**with Dimensionality Reduction Techniques?**

Jana Louise Bensing

Master Data Science for Society and Business, Constructor University

Data Mining

Prof. Wilhelm

May 31, 2023

## **Executive Summary**

### **Goals**

The goal of this research is three-fold. Firstly, we want to identify the most influential factors that may influence survival prognosis using RFE and Boruta. Secondly, we wanted to evaluate if SVM and NN can be utilized to effectively predict survival over the time-span of one year in patients diagnosed with HCC. Thirdly, we wanted to combine SVM and NN with the dimensionality reduction techniques of PCA and LDA, since previous research has shown that this may increase accuracy of prediction.

### **Data Background**

The Hepatocellular Carcinoma Dataset (HCC) was downloaded from Kaggle and originally collected by Santos, Abreu, García-Laencina, Carvalho and Simão (2015). It contains 49 risk factors for HCC as recognized by the European Organization for Research and Treatment of Cancer. The original dataset had 165 observations, but was later synthetically increased by Santos et al., (2015) to 204 observations, which is the data used in this study.

### **Approach**

The data was preprocessed and the most influential factors explored and visualized with Boruta and RFE. PCA visualization assisted in gaining understanding of class distributions. Then eight models were fitted, five of them SVM models and three NN models. The five SVM models were SVM (radial), SVM (polynomial), SVM (cross-validation), SVM-PCA, SVM-LDA. The tree NN models were NN, NN-PCA and NN-LDA. Models were evaluated using accuracy, sensitivity and specificity.

### **Results**

The four most influential factors for predicting survival over one year were Albumin, ALP, Hemoglobin and PS. An SVM model only including these four predictors achieved 72.58% accuracy. Boruta analysis classified 15 variables as confirmed and 5 as tentative. The SVM model including these 20 predictors achieved an accuracy of 77.42%. The NN models outperformed all SVM models with exception of the SVM-LDA model, which performed as well as the NN-PCA model and better than the NN-LDA model. The best performing model was the NN model, which achieved an accuracy of 93.54% in predicting survival over one year. Further dimensionality reduction did decrease accuracy in combination with SVM and NN with exception of the SVM-LDA model.

Primary liver cancer also called Hepatocellular Carcinoma (HCC) is currently the fifth most-common form of cancer worldwide and the third deadliest cancer. In most cases HCC develops in people that have struggled with liver disease for many years, but it may also occur without pre-existing liver diseases. In more than 90% of primary liver cancer it is HCC. In the United States around six new cases per 100.000 people are diagnosed every year, with men older than 60 having a heightened risk to develop HCC compared to women and younger men. In men HCC is currently the second leading cause of cancer deaths after lung cancer. (*What is hepatocellular carcinoma (HCC)?*)

Over the years research has identified multiple epidemiological and lifestyle risk factors that seem to heighten the risk of HCC significantly. Being infected with either viral hepatitis B (HBV) or hepatitis C (HCV) has been shown to increase the likelihood of developing HCC, especially if these chronic infections are not treated sufficiently. Further high alcohol consumption has been consistently linked to heightened risk of cirrhosis and alcoholic fatty liver disease, thus directly increases the risk of developing HCC. Also, smoking has been found to increase the risk for HCC. (HB; *Hepatocellular carcinoma (HCC)*)

On top of that it has been found that so called metabolic risk factors seem to become the main cause of HCC diagnoses worldwide. These metabolic risk factors include metabolic syndrome, obesity, type 2 diabetes and non-alcoholic fatty liver disease (NAFLD) or so-called lifestyle diseases, which are more prevalent in wealthy countries. Hence why decreasing the prevalence of obesity and diabetes is now recognized as an important step in decreasing HCC rates. (HB)

As with all cancers early detection is key to ensure that effective treatment is possible. HCC can be diagnosed among other techniques with MRI's, CT's, ultrasounds, biopsies, tests of the liver function or bone scans. However, especially in the beginning growth of the HCC tumor is very slow and years can pass, before the patient experiences any symptoms. For this reason, people, who are at risk for developing HCC should go to check-ups regularly. (*Hepatocellular carcinoma (HCC)*)

Treatment options depend heavily on the size and the spread of the HCC and range from surgical intervention that remove the carcinoma from the liver over liver transplantation up to chemotherapy. A new approach that has improved survival time for patients with advanced HCC tumors is systemic therapy, which prevents the tumor from growing. This treatment may increase the time of survival beyond the one-year mark after diagnosis with advanced HCC. (*What is hepatocellular carcinoma (HCC)?*)

One of the growing fields of generally cancer research is how Machine Learning (ML) and Deep Learning (DL) techniques can be utilized to integrate medical information and predict survival of patients. In order to better understand how the presence or absence of risk factors and the progress of the disease may affect survival chances. In turn this may also aid to further understand how HCC develops and progresses and provide patients with a more accurate estimation of their expected survival time.

Previous research has found that the algorithm Support Vector Machine and Neural Networks perform best when predicting the presence of breast cancer. (Omondiagbe, Veeramani & Sidhu, 2019) Also, it established that dimensionality reduction methods such as principal component analysis (PCA) and linear discriminant analysis (LDA) in combination with SVM and NN led to higher accuracy when predicting breast cancer.

Thus, this research aims to investigate if firstly, SVM and NN can also be utilized to predict survival in HCC patients over the span of one year. Secondly, it will explore if dimensionality reduction techniques combined with SVM and NN lead to higher accuracy when predicting survival. Lastly, this research will aim to identify the most influential factors that may determine survival prognosis.

### **Background of Data**

The Hepatocellular Carcinoma Dataset (HCC) was downloaded from Kaggle. (<https://www.kaggle.com/datasets/mrsantos/hcc-dataset?select=hcc-data-complete-balanced.csv>) The original dataset contains 49 predictors with 23 being quantitative and 26 being qualitative, the outcome variable and 165 observations each representing a real patient. The predictors were included based on the current clinical practice guidelines of the European Association for the Study of the Liver - European Organization for Research and Treatment of Cancer (EASL-EORTC). (EASL-EORTC, 2012) The outcome variable was binary with 0 standing for deceased during a one year time-span and 1 standing for survival of the patient during a one year time-span.

This data was collected at the University Hospital in Portugal and was donated to Kaggle by Santos, Abreu, García-Laencina, Carvalho and Simão (2015), who explored how the cluster-based oversampling method in combination with NN and logistic regression increases survival prediction in HCC patients. Santos et al. (2015) found that cluster-based oversampling did increase accuracy in all models. They achieved a maximum accuracy of 75.2% for their NN model with clustered and augmented data. This is why it was decided to

conduct this analysis with the updated Santos et al. (2015) dataset instead of the original dataset.

The updated dataset (hcc-data-complete-balanced.csv) was synthetically enlarged with the oversampling method and had 204 observations, meaning that 39 observations were artificially created. Moreover, all missing values of the original dataset were eliminated and estimated with synthetic data. Due to his advanced pre-processing by Santos et al. (2015) the updated dataset was also perfectly balanced in terms of the outcome variable. Thus, presenting a great basis for analysis.

## **Data Preprocessing**

### **Complete Analysis**

The dataset was important and the necessary packages for analysis were loaded. The used packages were dplyr, corrplot, caret, tidyverse, caTools, e17071, Boruta, keras and Mass. The original data frame had 34 integer columns and 16 columns containing strings, despite expressing numerical values. Thus, all character columns were firstly mutated to factor variables and secondly to integer variables. This was done, so that analysis could be performed on all variables.

It was necessary to first recode them as a factor before converting them to integer values, because without doing so all values were just recognized as missing values. After that all variables were integer values. It was decided to simplify analysis by operating only with one datatype with exception of the outcome variable class, which later was recoded as a factor with two levels 0 = D for deceased and 1 = S for survived.

The dataset was split into quantitative and qualitative variables to gain better overview of variables. The quantitative variables were: Age, Grams\_day, Packs\_year, PS, INR, AFP, Hemoglobin, MCV, Leucocytes, Platelets, Albumin, Total\_Bil, ALT, AST, GGT, ALP, TP, Creatinine, Major\_Dim, Dir\_Bil, Iron, Sat and Ferritin and can be seen in Graph 1. The qualitative variables were: Gender, Symptoms, Alcohol, HBsAg, HBeAg, HCVAb, Cirrhosis, Endemic, Smoking, Diabetes, Obesity, Hemochro, AHT, CRI, HIV, NASH, Varices, Spleno, PHT, PVT, Metastasis, Hallmark, Encephalopathy, Ascites, Nodule, Class and can be seen in Graph 2.

To prepare for analysis the whole dataset was scaled and centered to implement a common range of numbers and appropriately transform the qualitative variables, which were in almost all cases a 0 or 1. Then the data was analyzed for correlations above the standard cut-off value of 0.7. Column 10 (Smoking) and 26 (Packs\_year) and column 38 (ALT) and 39

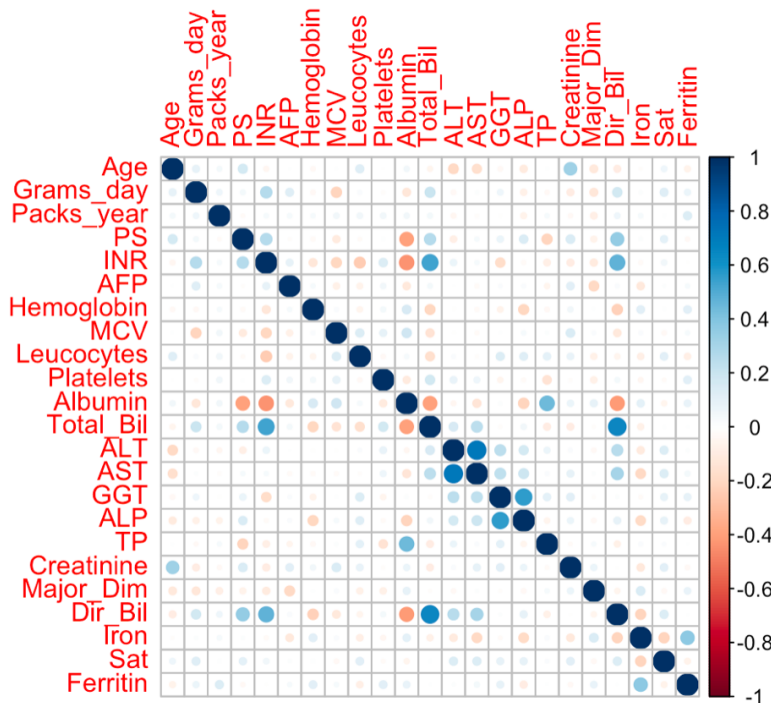
(AST) were flagged. The variables Packs\_year and AST were removed from the data. Then the data was checked for any missing values, which were not present.

### Specific Models

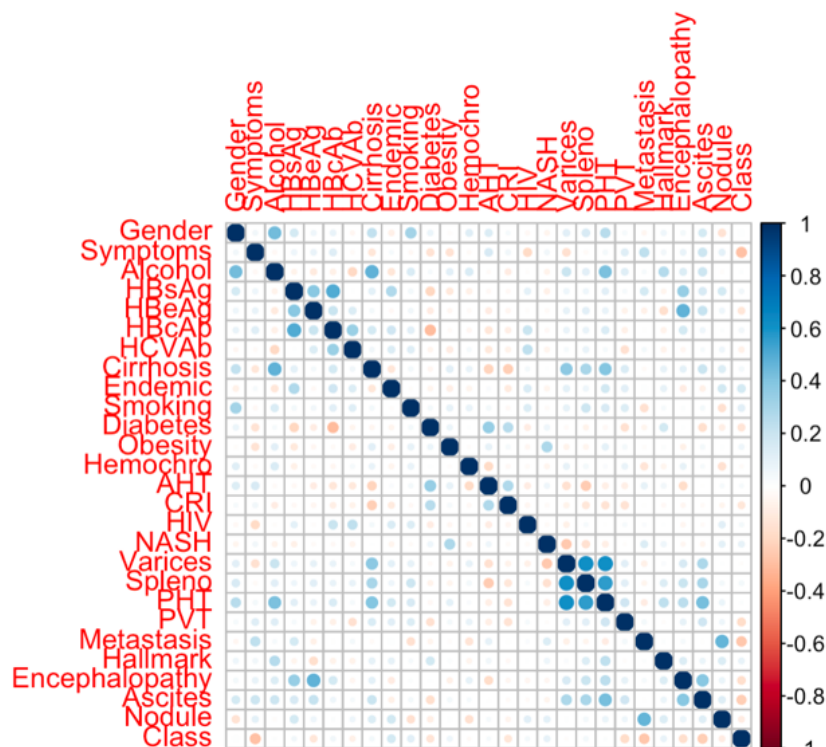
Since this study conducted analysis with two different algorithms namely, SVM and NN as well as them combined with the dimension reduction techniques PCA and LDA preprocessing and the utilized datasets varied. However, all algorithms were operated with the same split ratio of 0.7 with 70% being train and 30% being test data and survival being defined as the positive class. For the simple SVM and NN model, the SVM-RFE and the SVM-Boruta model the dataset, which excluded the two highly correlated features was used.

For the SVM-PCA, SVM-LDA, NN-PCA and NN-LDA only the scaled and centered data was used, since PCA and LDA already transform and exclude unnecessary variables. Thus, PCA and LDA can be seen as a form of preprocessing in our analysis. For the NN models the split data was transformed into the matrix format and a validation split of 0.2 was specified.

**Graph 1**



*Note:* This shows all qualitative variables and their correlations

**Graph 2**

*Note:* This shows all quantitative variables and their correlations.

### Data Exploration

Since all predictors in this dataset are confirmed by the EASL-EORTC to be risk factors for HCC, it is likely that many of them also influence the prediction of survival. Thus, it was decided to explore the data using two feature elimination techniques in order to identify most influential features when specifically predicting survival. However, since it is expected that all of the included predictors do also influence survival, because they influence development of HCC all predictors will be included in later analysis.

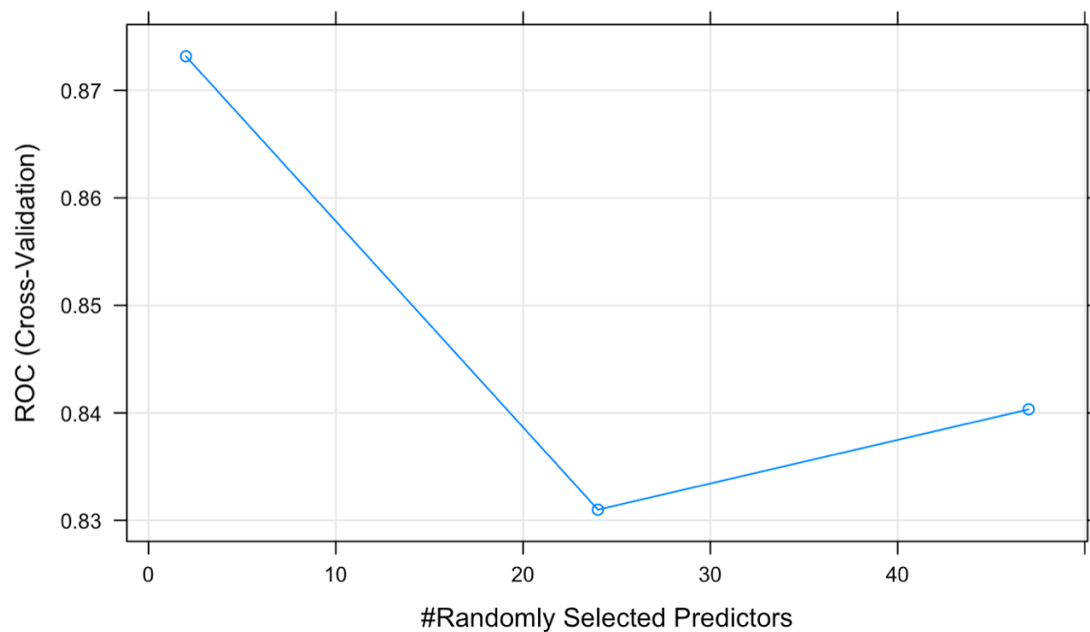
The methods used were recursive feature elimination technique (RFE) and the Boruta package in combination with the SVM algorithm. RFE uses the random-forest algorithm to check for combinations of different features. Each combination is rated in accuracy and ranked. The most highly ranked features are chosen. (Omondiagbe et al., 2019)

The Boruta package is also built on the random-forest algorithm, but evaluates the importance of features based on significance. It classifies predictors into three categories namely, rejected, tentative and confirmed. Where rejected means the predictor is not important, tentative means the algorithm is unsure how to classify the feature and confirmed that the predictor is important for predicting the outcome. (Team, 2018)

### SVM-RFE

The random forest algorithm was initiated with 15-times cross validation and achieved 79.03% in accuracy predicting survival with all included predictors. As indicated by Graph 3 four predictors were identified to be the most important when predicting survival. The four identified predictors were Albumin, ALP, Hemoglobin and PS. The SVM model fitted only with these four predictors achieved 72.58% accuracy, 70.97% sensitivity and 74.19% specificity.

**Graph 3**



*Note:* This graph shows the optimal number of randomly selected predictors chosen by a random forest model to predict diagnosis.

### SVM-Boruta

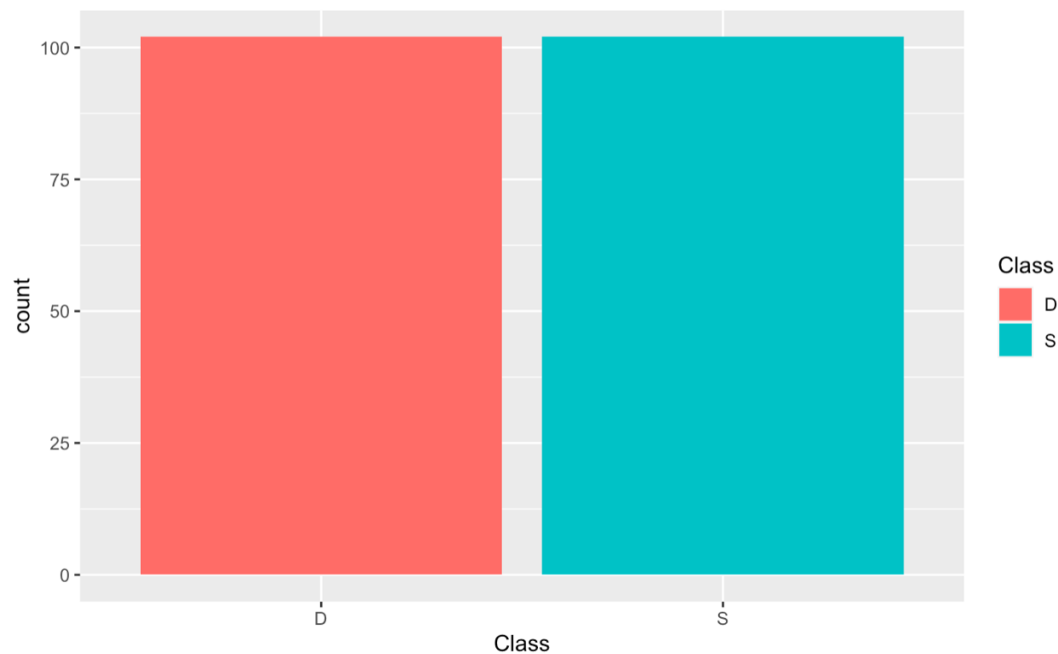
The Boruta analysis was run with 200 repetitions and in total identified 15 variables as confirmed, 5 variables as tentative and 28 as rejected. The confirmed variables were Symptoms, Endemic, Metastasis, Age, PS, Ascites, INR, Hemoglobin, Albumin, Total\_Bil, GGT, ALP, Creatinine, Dir\_Bil and Sat. The tentative variables were Encephalopathy, ALT, Platelets, Major\_Dim and Iron. It was decided to include all tentative and confirmed variables as predictors for the SVM mode. The SVM model including these 20 predictors achieved an accuracy of 77.42%, sensitivity of 77.42% and a specificity of 77.42%.



## Visualizations

### *Class Distribution*

**Graph 4**

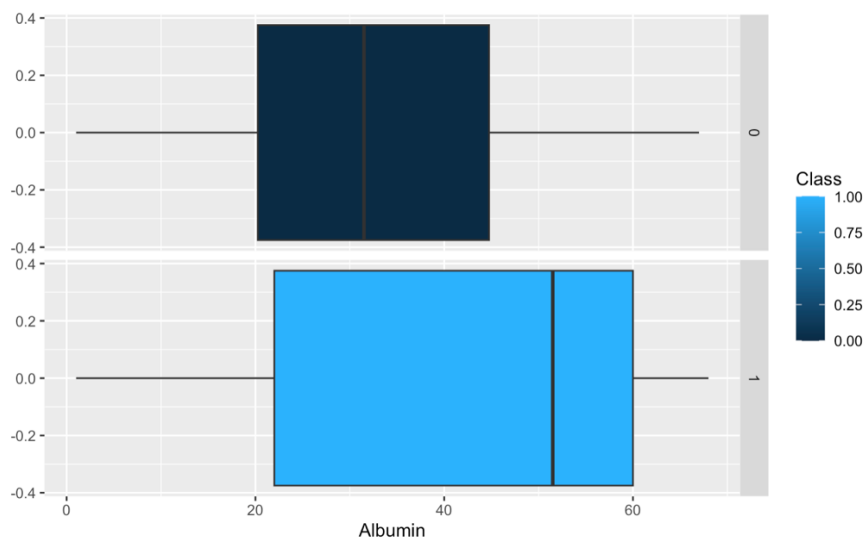


*Note:* This graph shows that both classes are equally present in the dataset.

## Four most important Predictors

### *1. Albumin*

**Graph 5**

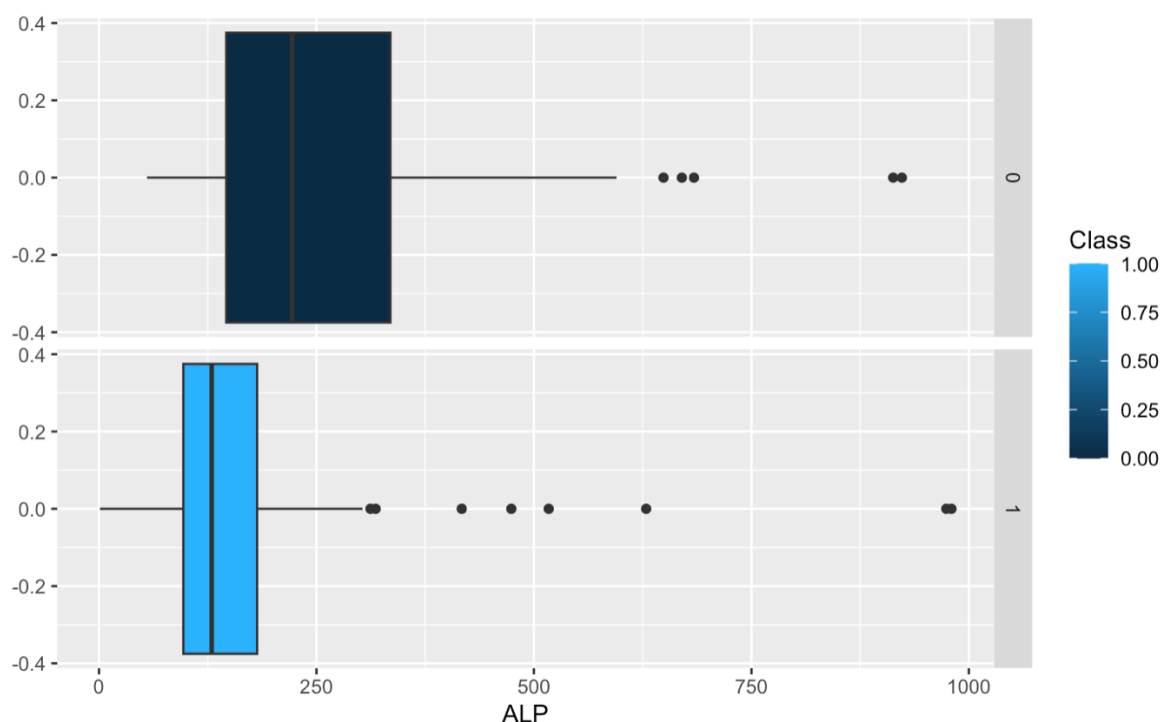


*Note:* This graph shows the distributions of Albumin levels for patients of both classes.

*Interpretation:* Albumin is a protein made by the liver. Too low albumin levels can be an indication of a liver disease including HCC. As Graph 5 shows the patients that survived had on average much higher albumin levels compared to deceased patients as indicated by the higher median. However, there seems to be an overlap in the distributions of albumin levels across both classes. (U.S. National Library of Medicine)

## 2. ALP

**Graph 6**

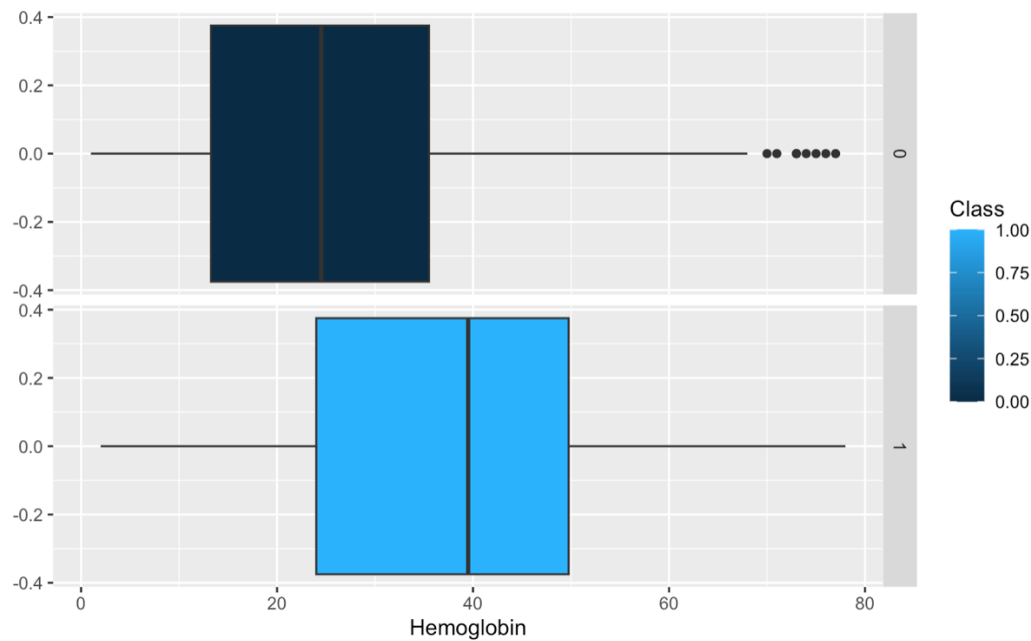


*Note:* This graph shows the distributions of ALP levels for patients of both classes.

*Interpretation:* Alkaline phosphatase (ALP) is a protein found in all body tissue including the liver. Too high ALP levels can be an indication of a liver disease including HCC. As Graph 6 shows the patients that survived had on average lower ALP levels compared to deceased patients as indicated by the lower median. There seems to be a clearer separation in distributions across both classes compared to albumin levels. (*Alkaline phosphatase*)

### 3.Hemoglobin

**Graph 7**

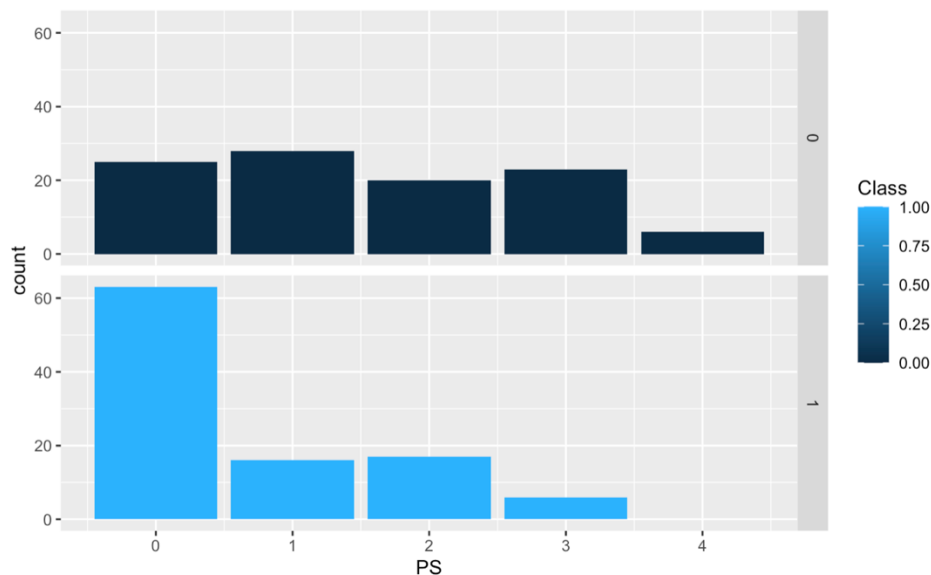


Note: This graph shows the distributions of Hemoglobin levels for patients of both classes.

*Interpretation:* Hemoglobin is a protein in the red blood cells, which carries oxygen. Too low Hemoglobin levels can be an indication of many chronic diseases including HCC. As Graph 7 shows the patients that survived had on average higher Hemoglobin levels compared to deceased patients as indicated by the higher median. However, there seems to be an overlap in the distributions of Hemoglobin levels across both classes. (Mayo Foundation for Medical Education and Research, 2022)

#### 4. PS

**Graph 8**



*Note:* This graph shows the distributions of pathology for patients of both classes.

*Interpretation:* Pathology in this data indicated the level of how far the condition of the patient deviated from normal levels. With 0 indicating that there was no deviation from normal levels at point of measurement, meaning the patient was cancer free and 4 indicating that the cancer had a highly progressed level at point of measurement. As can be seen in Graph 8 there was no patient, who survived the span of one year if identified with level 4 pathology. The majority of patients, who survived over the span of one year were cancer free at the time of measurement.

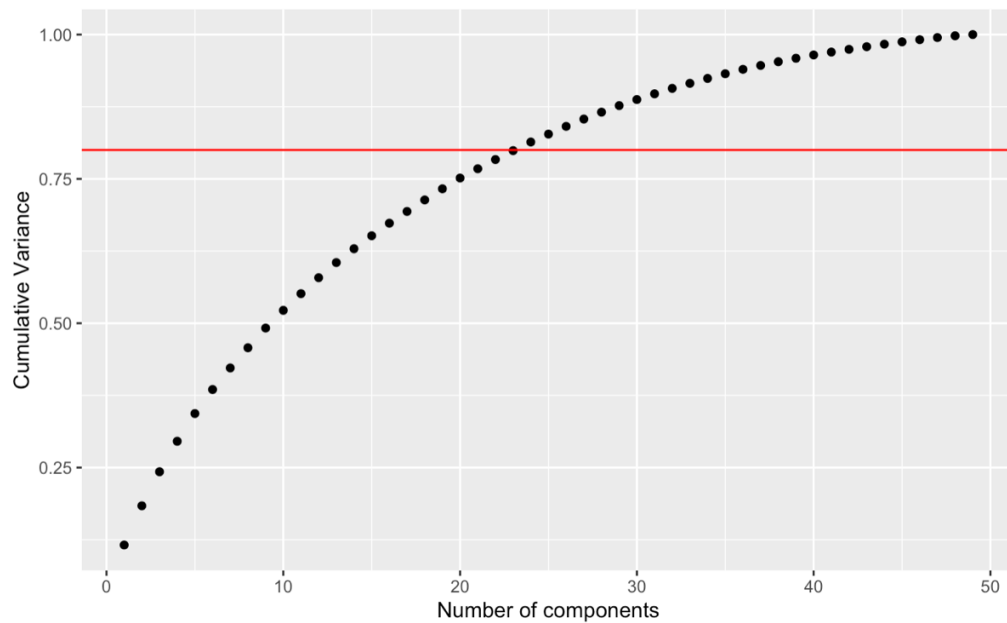
#### PCA

The PCA was conducted on the whole dataset. Cumulative the first ten principal components did only explain 52.23% in variance. As can be seen in Graph 9 the cumulative variance did increase only marginally with each component. Over 20 components were necessary to reach a cumulative variance of 80%. For further analysis 23 components were included. However, since the variance only increases marginally with each component this could be evidence, that this particular dataset may not improve prediction by utilizing dimensionality reduction.

Further Graph 9 illustrates the distribution of the two classes across the four most influential components. As expected, based on Graph 6, 7, 8 and 9 there seems to be extreme

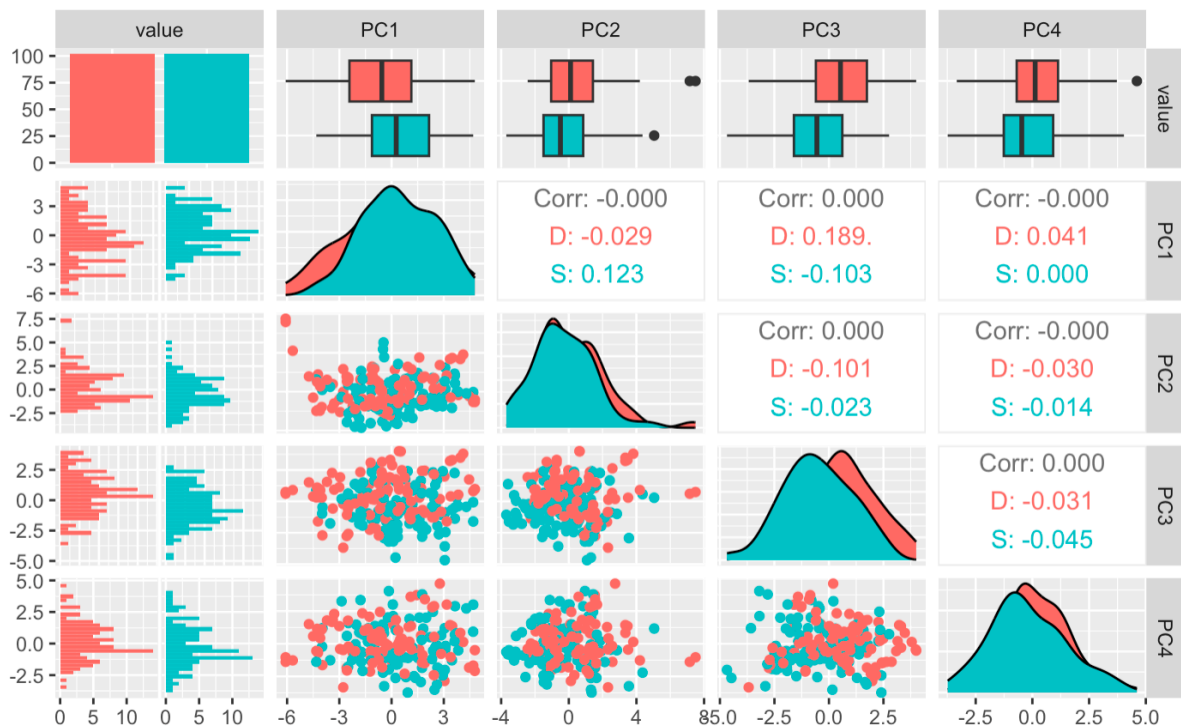
overlap in the distributions. Based on this it can be concluded that the classes are not linearly separable.

**Graph 8**



*Note:* This graph shows the cumulative variance for the number of components.

**Graph 9**



*Note:* This graph shows the distributions of the first four principal component.

### Data Analysis

Since there is evidence that algorithms that can combine linear and non-linear techniques perform better on complex problems such as cancer prediction the SVM and NN algorithms were chosen for this research. Further indication that linear separation might not be possible was given, by the four most influential predictors, since all of them seemed to overlap in their distribution of values. Additionally, SVM and NN were combined with PCA and LDA, because research indicated that this further benefits accuracy of models. (Omondiagbe et al, 2019)

SVM fits a hyperplane in a n-dimensional space, which aims at separating two classes by maximizing their distance. The function by which the hyperplane is fitted can be changed by defining the kernel as either linear, radial or polynomial. Especially, radial and polynomial kernels are beneficial if the data is not linearly separable into two classes, because they add dimensions that may make separation possible. <https://towardsdatascience.com/svm-and-kernel-svm-fed02bef1200> Hence why it was decided to only implement either radial or polynomial kernel.

NN consist of an input layer, one or multiple hidden layers and one output layer. Each layer can be run with a different activation function and thus can combine linear and non-linear processing of data. Additionally, there is the option to integrate different optimizer and loss functions that also perform well on non-convex data. (reference) For this research across all NN models the ADAM optimizer and the binary cross entropy loss function were implemented. The networks were run with Epoch = 20 and Batch = 5 and a validation split of 0.2. Input and hidden layers ran with RELU activation, while the output layer ran with the sigmoid function.

### Results

#### SVM

Three SVM models were fitted. The first model used a radial kernel and achieved an accuracy of 75.81% a sensitivity of 70.97% and a specificity of 80.65%. The second model used a polynomial kernel and achieved an accuracy of 79.03% a sensitivity of 93.55% and a specificity of 64.52%. The third model used a linear kernel and performed a 10-times cross-validation and iterated over the sigma values of 0.01, 0.015 and 0.2 and cost values of 0.001, 0.01, 0.1, 1, 5, 10 and 100. The final values used in the model were sigma = 0.01 and cost = 100. It also achieved an accuracy of 79.03% as the polynomial model, but with a sensitivity of 83.87% and a specificity of 74.19%.

**SVM-PCA**

The model combining SVM with PCA achieved only 62.9% of accuracy, 77.2% sensitivity and 48.39% specificity.

**SVM-LDA**

The model combining SVM with LDA achieved the highest performance out of all SVM models with an accuracy of 85.48%, a sensitivity of 87.1% and a specificity of 83.87%. Thus, it was also the most balanced model out of all SVM models.

**NN**

The NN model was designed with two hidden layers. There were 47 neurons in the input layer congruent with all predictors of the data, 30 neurons in the first hidden layer, 16 neurons in the second hidden layer and one output neuron. The best NN model achieved an accuracy of 93.55% with a sensitivity of 100% and specificity of 87.1%. Thus, it was the best performing NN model and best overall model.

**NN-PCA**

The model combining NN and PCA was designed with two hidden layers. There were 23 neurons in the input layer congruent with all predictors of the data, 13 neurons in the first hidden layer, 9 neurons in the second hidden layer and one output neuron. The best NN-PCA model achieved an accuracy of 85.45%, a sensitivity of 100% and a specificity of 70.97%.

**NN-LDA**

The model combining NN and LDA was designed with one hidden layer. There were two input neurons, two neurons in the hidden layer and one output neuron. The best NN-LDA model achieved an accuracy of 83.87%, a sensitivity of 87.1% and a specificity of 80.65%.

**Discussion**

As identified by the RFE model the four most influential factors for predicting survival over one year were Albumin, ALP, Hemoglobin and PS. An SVM model only including these four predictors achieved 72.58% accuracy, 70.97% sensitivity and 74.19% specificity. As shown by the visualizations of these factors as well as the distributions of the

two classes across the first four principal components, distributions of for both classes seem to overlap heavily. This was further evidence that the classes were not linearly separable. Thus, confirmed the choice of algorithms, namely SVM and NN was appropriate to predict survival, because they have to the ability to combine linear and non-linear properties for prediction.

Further, Boruta analysis classified 15 variables as confirmed and 5 as tentative. The 15 confirmed variables did include Albumin, ALP, Hemoglobin and PS. The SVM model including these 20 predictors achieved an accuracy of 77.42% as well as a sensitivity and specificity of 77.42%. Notably, the Boruta analysis did not include confirmed risk factors for HCC such as alcohol, cirrhosis, obesity or smoking, meaning that despite being risk factors for HCC they are not as important in predicting survival of patients that are already diagnosed with HCC.

The best performing SVM model was the SVM-LDA model which achieved 85.48% accuracy, a sensitivity of 87.1% and a specificity of 83.87%. The best performing NN model as well as overall best performing model was the NN model with 93.55% accuracy, a sensitivity of 100% and specificity of 87.1%. However, with exception of the SVM-LDA the NN-PCA and NN-LDA model did also outperform all other SVM models with 85.45% accuracy and 83.87% accuracy respectively.

This shows that NN is more suitable in effectively predicting survival of HCC patients over the span of one year compared to the SVM algorithms. Additionally, our results indicate that when predicting survival dimensionality reduction techniques seemed to decrease accuracy of both algorithms. However, there is some evidence that SVM performance may benefit from LDA.

An explanation for this finding could be that the predictors included in this data have already been established by research to be the most important risk factors for HCC. Thus, the data has already been highly curated and targeted towards identifying HCC risk. Meaning that no random information and barely highly correlated information was included in this data.

Despite that most common risk factors for HCC such as cirrhosis, alcohol consumption, smoking, obesity or infection with HBV or HBC were not identified as main predictors of survival. In fact, the most important predictors were Albumin, ALP, Hemoglobin and PS, which shows us that common risk factors for HCC are not necessarily as influential in predicting survival of HCC patients.



Hence a practical application of this research could be to especially focus on collecting data that includes the most influential data when predicting survival instead of collecting all 49 risk factors for HCC. A shorter questionnaire could possibly increase participation rates due to being less time consuming and being perceived as less invasive by patients, without overly reducing predictive power. In turn this could solve the issue of the original dataset, which had many missing values and not many observations, which had to be compensated by synthetic data.

### **Conclusion**

Based on this research NN seems to be the superior algorithm when predicting survival of HCC patients over the timespan of one year. Further SVM performance may increase when combined with LDA, but still performs lower than an NN model. The four most influential factors for predicting survival over one year were Albumin, ALP, Hemoglobin and PS. Thus, future research could benefit from focusing on NN models when predicting survival as well as targeting questionnaires more towards predictors of survival of diagnosed HCC patients instead of risk factors for HCC.

### **Limitations**

The generalizability and applicability of these findings is limited due to two reasons. Firstly, only one dataset was analyzed, meaning our sample is very small. Secondly, the data used in this research was partly synthetically reduced. Thus, we cannot exclude that values have been either wrongly imputed or that artificial observations may not reflect data patterns of real patients. Additionally, this dataset had very curated predictors, meaning it is possible that other datasets, which aim at predicting survival of HCC patients may still benefit from dimensionality reduction.

### **Appendix**

The Appendix containing the code can be found under the following link:

<https://github.com/janabensing/Data-Mining-JLB.git>

### References

- Activation functions in neural networks [12 types & use cases]*. V7. (n.d.). Retrieved 31<sup>st</sup> of May, 2023 from <https://www.v7labs.com/blog/neural-networks-activation-functions>
- Alkaline phosphatase (ALP): What it is, causes & treatment*. Cleveland Clinic. professional, C. C. medical. (n.d.-a). Retrieved 31<sup>st</sup> of May, 2023 from [https://my.clevelandclinic.org/health/diagnostics/22029-alkaline-phosphatase-alp#:~:text=High%20alkaline%20phosphatase%20\(ALP\)%20levels,ALP%20than%20bone%20disorders%20do.](https://my.clevelandclinic.org/health/diagnostics/22029-alkaline-phosphatase-alp#:~:text=High%20alkaline%20phosphatase%20(ALP)%20levels,ALP%20than%20bone%20disorders%20do.)
- EASL–EORTC Clinical Practice Guidelines: Management of hepatocellular carcinoma. (2012). *European Journal of Cancer*, 48(5), 599–641. <https://doi.org/10.1016/j.ejca.2011.12.021>
- HB;, M. K. J.-S. (n.d.). *Epidemiology of hepatocellular carcinoma*. Hepatology (Baltimore, Md.). Retrieved 31<sup>st</sup> of May, 2023 from <https://pubmed.ncbi.nlm.nih.gov/32319693/>
- Hepatocellular carcinoma (HCC): Causes, symptoms, treatments & prognosis*. professional, C. C. medical. (n.d.-b). Cleveland Clinic. Retrieved 31<sup>st</sup> of May, 2023 from <https://my.clevelandclinic.org/health/diseases/21709-hepatocellular-carcinoma-hcc>
- Mayo Foundation for Medical Education and Research. (2022, February 11). *Hemoglobin test*. Mayo Clinic. Retrieved 31<sup>st</sup> of May, 2023 from <https://www.mayoclinic.org/tests-procedures/hemoglobin-test/about/pac-20385075>
- Omondiagbe, D. A., Veeramani, S., & Sidhu, A. S. (2019). Machine learning classification techniques for breast cancer diagnosis. *IOP Conference Series: Materials Science and Engineering*, 495, 012033. <https://doi.org/10.1088/1757-899x/495/1/012033>
- Santos, M. S., Abreu, P. H., García-Laencina, P. J., Simão, A., & Carvalho, A. (2015). A new cluster-based oversampling method for improving survival prediction of hepatocellular carcinoma patients. *Journal of Biomedical Informatics*, 58, 49–59. <https://doi.org/10.1016/j.jbi.2015.09.012>
- SVM and kernel SVM. learn about SVM or support ... - towards data science. (n.d.). Retrieved 31<sup>st</sup> of May, 2023 from <https://towardsdatascience.com/svm-and-kernel-svm-fed02bef1200>
- Team, D. (2018, March 7). *Boruta feature selection in R*. DataCamp. Retrieved 31<sup>st</sup> of May, 2023 from <https://www.datacamp.com/tutorial/feature-selection-R-boruta>
- U.S. National Library of Medicine. (n.d.). *Albumin blood test: Medlineplus medical test*. MedlinePlus. Retrieved 31<sup>st</sup> of May, 2023 from <https://medlineplus.gov/lab-tests/albumin-blood-test/>

*What is hepatocellular carcinoma (HCC)?*. Mount Sinai Health System. (n.d.). Retrieved 31<sup>st</sup> of May, 2023 from

<https://www.mountsinai.org/care/cancer/services/liver/hepatocellular-carcinoma#:~:text=Primary%20liver%20cancer%E2%80%94cancer%20that,liver%20disease%20for%20many%20years.> ;

<https://www.ncbi.nlm.nih.gov/books/NBK559177/>