

The George Washington University

Titanic Indi-Report

The Columbian College of Art and Science

Tan, Yueyan

2018-8-12

Directory

Introduction.....	1
Description of individual work	1
background information.....	1
My portion of the work.....	4
Calculation	11
References	11

Introduction

The sinking of the RMS Titanic is one of the most infamous shipwrecks in history. On April 15, 1912, during her maiden voyage, the Titanic sank after colliding with an iceberg, killing 1502 out of 2224 passengers and crew. This sensational tragedy shocked the international community and led to better safety regulations for ships.

One of the reasons that the shipwreck led to such loss of life was that there were not enough lifeboats for the passengers and crew. Although there was some element of luck involved in surviving the sinking, some groups of people were more likely to survive than others, such as women, children, and the upper-class.

In our project, we found out the connection between passengers' features and the probability of death. We used decision tree for finding out the most important features that affect a passenger's death in a shipwreck. Secondly, SVM algorithm were used to find a specific decision boundary to classify and finally, we compare the performances of prediction by using these three algorithms.

And our group shared works are shown in the following table:

ITEM	Zhengting Zhu	Yueyan Tan	Jinyu Hu
Coding	SVM & Naïve Bayes	DATA & Decision Tree	
Report	Part 4	Part 5 Result	Part 1,2,3
PPT	Part 4	Part 5 Result	Part 1,2,3

Description of individual work

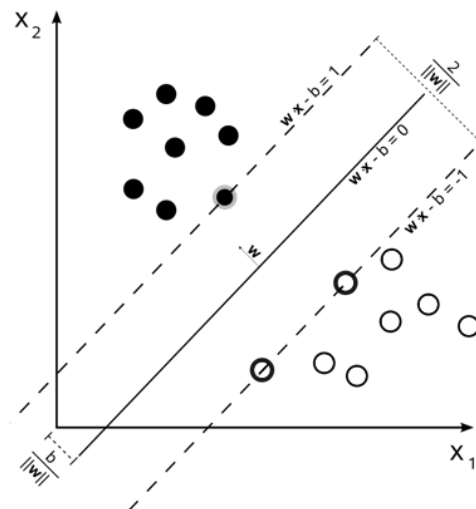
background information

In this project, we use three models listing below:

Support Vector Machines (SVM):

In machine learning field, support vector machines (SVM) are supervised learning models that analyzing data for classification and regression.

A Support Vector Machine is a discriminative classifier formally defined by a separating hyperplane. In other words, given labeled training data, the algorithm outputs an optimal hyperplane which categorizes new

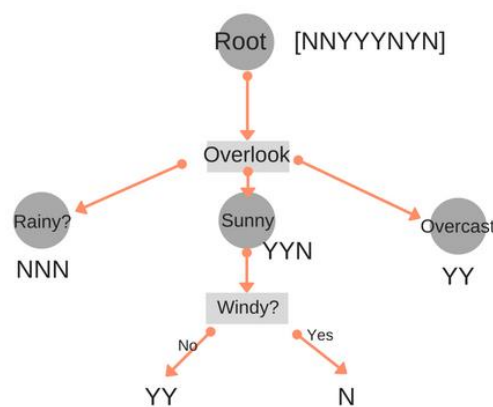


examples. The easiest model picture looks like this:

Decision Tree:

In Data Science field, a decision tree is support tool that uses a tree-like graph to show the decisions. A decision tree specially starts with a single node, and then branches into possible outcomes. Each of those outcomes leads to additional nodes, which branch into other possibilities. It is one approach to display an algorithm that only contains conditional control statements in classification analysis.

One possible model looks like this:



Some important factors of Decisions Tree are:

$$\text{Entropy} = \sum -p_i \log_2 p_i$$

$$\text{Gini} = 1 - \sum (p_i^2)$$

Information Gain (IG) = Entropy (parent) – Average Entropy (children)

Naïve Bayes

In machine learning field, naive Bayes classifiers are a group of simple probabilistic classifiers based on Bayes' theorem with independence assumptions between the features. It is a powerful algorithm used for real time prediction, text classification, recommendation system, etc.

Key mathematic equation list below:

$$P(C(Class)|X(Features)) = \frac{P(X|C) \times P(C)}{P(X)}$$

$P(C(Class)|X(Features))$: Posterior Probability of class (C) given predictor (X)

$P(X|C)$: Likelihood - the conditional probability of the predictor

$P(C)$: Prior Probability of the Class

$P(X)$: Total probability

My portion of the work

My work in Coding

```
1 import numpy as np
2 import pandas as pd
3 from sklearn.tree import DecisionTreeClassifier
4 from sklearn.metrics import confusion_matrix
5 from sklearn import tree
6 from sklearn.model_selection import train_test_split
7 from sklearn.svm import SVC
8 from sklearn.metrics import accuracy_score, roc_curve, roc_auc_score
9 from sklearn.metrics import classification_report
10 import seaborn as sns
11 import matplotlib.pyplot as plt
12 import warnings
13 from pydotplus import graph_from_dot_data
14 import webbrowser
15 from sklearn.naive_bayes import GaussianNB
16 warnings.filterwarnings("ignore")
17 import os
18 os.environ["PATH"] += os.pathsep + r'C:\Users\yueya\Anaconda3\Library\bin\graphviz/'
19
20 #importing data
21 tt = pd.read_csv(r'data.csv', sep=',', header=0)
22 ###-----
23 ###-----
24 #data preprocessing
25 # look at first few rows
26 tt.head()
27
28 # replace missing characters as NaN
29 tt.replace('?', np.NaN, inplace=True)
30
31 # check the structure of mushroom_data
32 tt.info()
33
34 # check the null values in each column
35 print(tt.isnull().sum())
36
37 # check the summary of the mushroom_data
38 tt.describe(include='all')
39
40 # replace categorical mushroom_data with the most frequent value in that column
41 tt = tt.apply(lambda x: x.fillna(x.value_counts().index[0]))
42
43 # drop Cabin, ticket and name of passenger
44 tt.drop(['Cabin', 'Name', 'Ticket', 'Embarked'], axis=1, inplace=True)
45
46 # again check the null values in each column
47 print(tt.isnull().sum())
48
49 #Encoder features
50 X_data = pd.get_dummies(tt.iloc[:, 1:])
51 X = X_data.values
52 Y = tt.values[:, 0]
53 Y = Y.astype('int')
54 ###-----
```

```

54  """-----
55
56  # split the dataset into train and test
57  X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.3, random_state=100)
58
59  # data visualization
60  ntrain=tt.shape[0]
61  import seaborn as sns
62  train=tt[:ntrain]
63  sns.pairplot(train, vars=['Age', 'Sex', 'Pclass'], hue='Survived',)
64  plt.title("data visualization")
65  plt.show()
66
67  # data standardize
68  from sklearn.preprocessing import StandardScaler
69  sc = StandardScaler()
70  sc.fit(X_train)
71  X_train = sc.transform(X_train)
72  X_test = sc.transform(X_test)
73
74  """-----
75
76  # Decision Tree
77
78  # perform training with giniIndex.
79  # creating the classifier object
80  clf_gini = DecisionTreeClassifier(criterion="gini", random_state=100, max_depth=3, min_samples_leaf=5)
81
82  # performing training
83  clf_gini.fit(X_train, y_train)
84
85  # perform training with entropy.
86  # Decision tree with entropy
87  clf_entropy = DecisionTreeClassifier(criterion="entropy", random_state=100, max_depth=3, min_samples_leaf=5)
88  clf_entropy.fit(X_train, y_train)
89  #prediction
90  y_pred_gini = clf_gini.predict(X_test)
91
92  y_pred_entropy = clf_entropy.predict(X_test)
93
94  # calculate metrics gini model
95  print("\n")
96  print("Results Using Gini Index: \n")
97  print("Classification Report: ")
98  print(classification_report(y_test, y_pred_gini))
99  print("\n")
100  print("Accuracy : ", accuracy_score(y_test, y_pred_gini) * 100)
101  print("\n")
102  print('-'*80 + '\n')
103
104  # calculate metrics entropy model
105  print("\n")
106  print("Results Using Entropy: \n")
107  print("Classification Report: ")
108  print(classification_report(y_test, y_pred_entropy))
109  print("\n")
110  print("Accuracy : ", accuracy_score(y_test, y_pred_entropy) * 100)
111  print('-'*80 + '\n')

```

```

113 # confusion matrix for gini model
114 conf_matrix = confusion_matrix(y_test, y_pred_gini)
115 class_names = tt.Survived.unique()
116 df_cm = pd.DataFrame(conf_matrix, index=class_names, columns=class_names)
117
118 plt.figure(figsize=(5,5))
119 hm = sns.heatmap(df_cm, cbar=False, annot=True, square=True, fmt='d', annot_kws={'size': 20}, yticklabels=df_cm.columns, xticklabels=df_cm.columns)
120 hm.yaxis.set_ticklabels(hm.yaxis.get_ticklabels(), rotation=0, ha='right', fontsize=20)
121 hm.xaxis.set_ticklabels(hm.xaxis.get_ticklabels(), rotation=0, ha='right', fontsize=20)
122 plt.ylabel('True label', fontsize=20)
123 plt.xlabel('Predicted label', fontsize=20)
124 plt.title("Decision Tree - Gini")
125 plt.tight_layout()
126 plt.show()
127
128 # confusion matrix for entropy model
129 conf_matrix = confusion_matrix(y_test, y_pred_entropy)
130 class_names = tt.Survived.unique()
131 df_cm = pd.DataFrame(conf_matrix, index=class_names, columns=class_names)
132
133 plt.figure(figsize=(5,5))
134 hm = sns.heatmap(df_cm, cbar=False, annot=True, square=True, fmt='d', annot_kws={'size': 20}, yticklabels=df_cm.columns, xticklabels=df_cm.columns)
135 hm.yaxis.set_ticklabels(hm.yaxis.get_ticklabels(), rotation=0, ha='right', fontsize=20)
136 hm.xaxis.set_ticklabels(hm.xaxis.get_ticklabels(), rotation=0, ha='right', fontsize=20)
137 plt.ylabel('True label', fontsize=20)
138 plt.xlabel('Predicted label', fontsize=20)
139 plt.title("Decision Tree - Entropy")
140 plt.tight_layout()
141 plt.show()
142
143 # display decision tree
144 dot_data = tree.export_graphviz(clf_gini, filled=True, rounded=True, class_names='survived', feature_names=tt.iloc[:, 0:].columns, out_file=None)
145
146 graph = graph_from_dot_data(dot_data)
147 graph.write_pdf("decision_tree_gini.pdf")
148 webbrowser.open_new(r'decision_tree_gini.pdf')
149
150 dot_data = tree.export_graphviz(clf_entropy, filled=True, rounded=True, class_names='survived', feature_names=tt.iloc[:, 0:].columns, out_file=None)
151
152 graph = graph_from_dot_data(dot_data)
153 graph.write_pdf("decision_tree_entropy.pdf")
154 webbrowser.open_new(r'decision_tree_entropy.pdf')
155
156 #%%

```

My work in Report

Results

Decision Tree

For the first part of our project, we used Decision-Tree algorithm for figuring out the degree of importance of each feature. We began with Decision Tree for the following reasons:

1. It is an algorithm with widely used and could be well-visualized, which would be much easier for the readers from different background to understand our project results.
2. Decision Tree is simple to understand and to interpret, which also benefit the readers.
3. Decision Tree requires little data preparation. Other techniques such as SVM often require data normalization, dummy variables need to be created and blank values to be removed.
4. It is able to handle both numerical and categorical data. Other techniques like logistic regression and other statistical algorithms are usually specialized in analyzing datasets that have only one type of variable and the data we used (The Titanic data set) includes both categorical as well as numerical data, which means decision tree would be a good tool for the project.
5. It is possible to validate a model using statistical tests. That makes it possible to account for the reliability of the model, such that we could fit the ROC for the prediction which would be much simple for readers to get our results.

Firstly, we got the decision tree by using entropy for calculating, as is shown in figure1:

Before getting the result, we may guess that the most important features for affecting the possibility of survival would be gender or age, while according to the result, we could find out that the most important feature for deciding whether the passenger could survival from the shipwreck is “parch” (number of parents / children aboard the Titanic), which indicates that passengers with parents or children together would have much stronger desire of survival or higher rate or survival than others when facing the big disaster.

Moving down, we could see that “sibsp” (number of siblings / spouses aboard the Titanic) and “pclass” (Ticket class) are also important features that have effect on possibility of the survival. The reason of “sibsp” maybe the same with the feature “parch”. While “Pclass” also have such importance may because the rich or upper class always have priority over resources which is a harsh reality in our society.

Then, we got the decision tree by using Gini for calculating, as is shown in figure2:

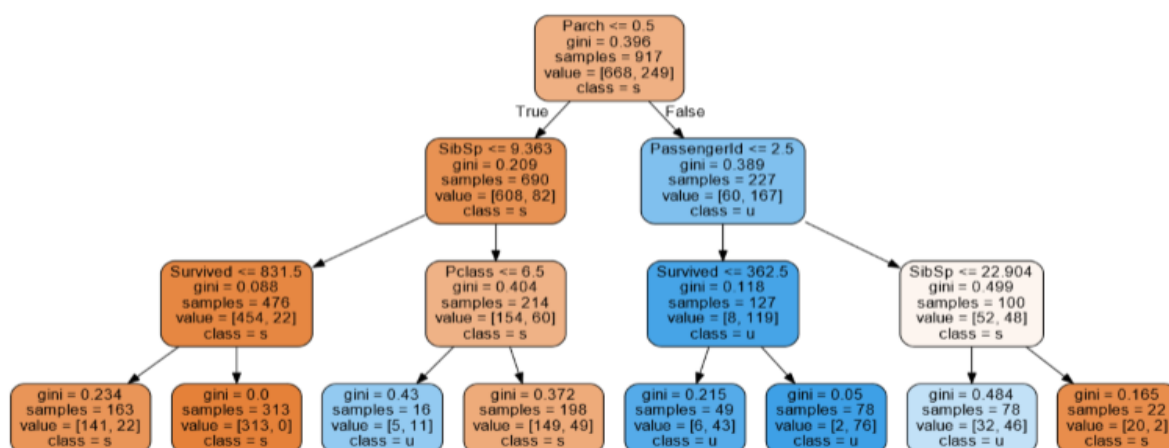


Figure 2 Decision Tree-GINI

Firstly, we got the decision tree by using entropy for calculating, as is shown in figure1:

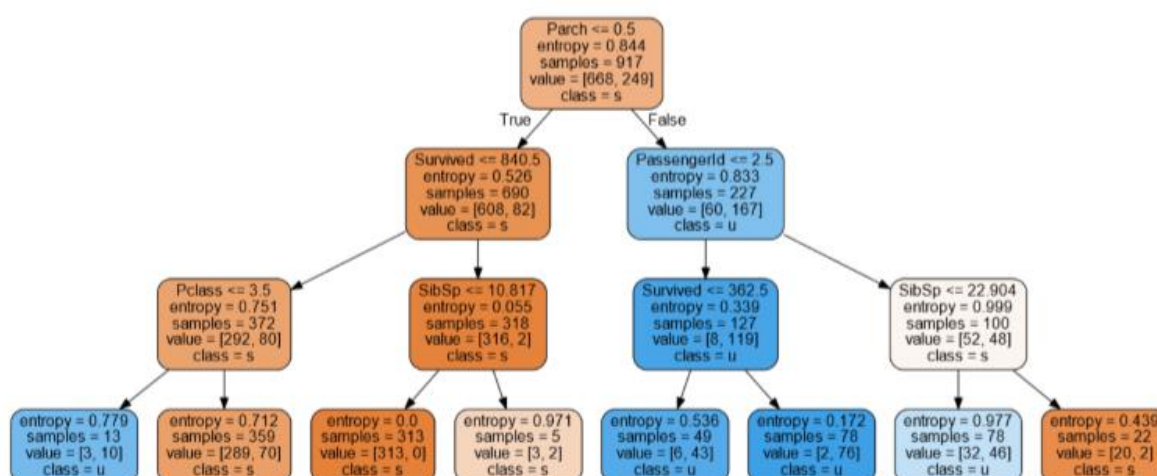


Figure 1 Decision Tree-Entropy

Additionally, we also predict in both methods as is shown in Figure 3 and Figure 4:

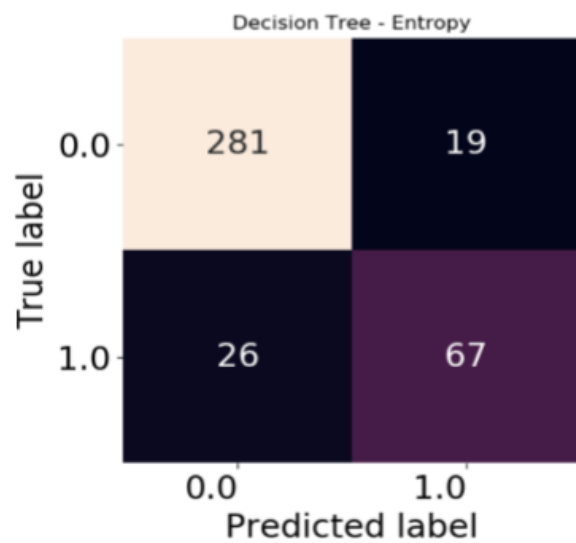


Figure 3 Entropy

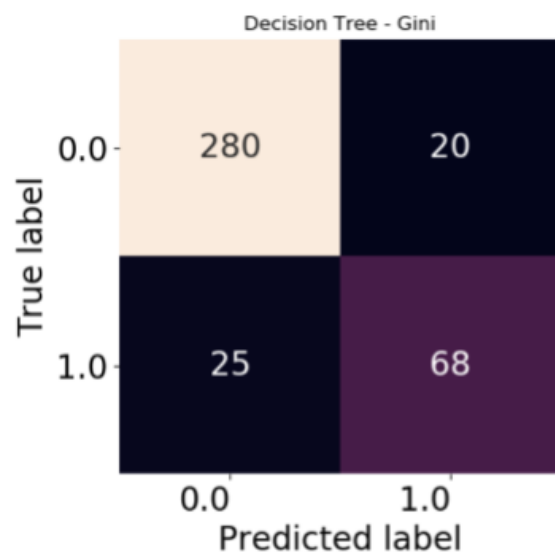


Figure 4 GINI

From the Figure 2, Figure 3 and Figure 4, it is easy to find out that Figure 2 is almost the same as figure 1 we got from entropy which means that both of these two calculations provide quite similar information and would not have much differences in the results. And for the performance of prediction with this two method, they were also the same (The prediction accuracy of both GINI and Entropy are 88.55%). While in the real life's application, more researchers would like to use Gini for calculation, then what is the differences between these two calculations ? After summarizing, we find that they are different in the following ways:

1. Entropy mostly used on categorical data set (discrete data sets).
2. Gini mostly used on contagious data set.
3. Entropy is slower to compute, while Gini is faster.

Support Vector Machine

For the second part of our project, Support Vector Machine were used for getting the specific decision boundary between two classes (survived and died), which may be helpful for discovering the effective features for that possibility of survival in shipwreck.

A Support Vector Machine (SVM) is a discriminative classifier formally defined by a separating hyperplane. In other words, given labeled training data (supervised learning), the algorithm outputs an optimal hyperplane which categorizes new examples. The reason we used SVM in this part because of the following reasons:

1. It has a regularization parameter, which makes the user think about avoiding overfitting. Secondly it uses the kernel trick, so you can build in expert knowledge about the problem via engineering the kernel.
2. An SVM is defined by a convex optimization problem (no local minima) for which there are efficient methods (e.g. SMO).
3. It is an approximation to a bound on the test error rate, and there is a substantial body of theory behind it which suggests it should be a good idea.

In this part of our project, the specific decision boundary would be helpful for classify correctly and we could get the equation of the decision boundary. Besides, SVM also used for predict the possibility of survival of a passenger in the shipwreck, therefore the results of SVM and Decision Tree could be used for compare the performance of different algorithms.

According to the results we got from SVM, as shown in Figure 5 and Figure 6:

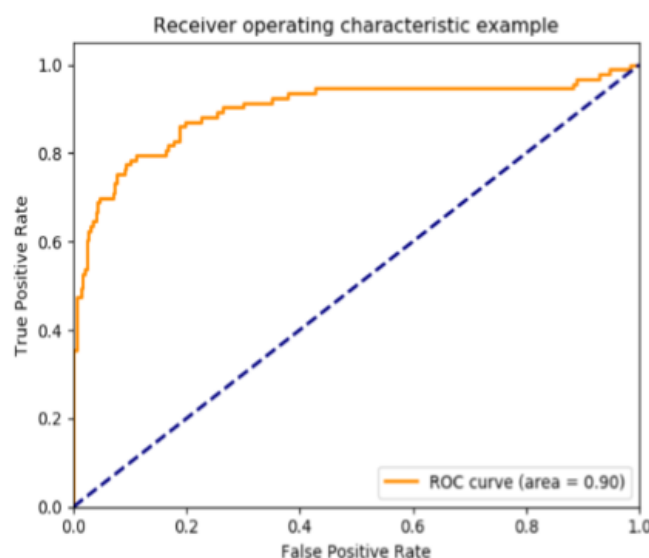


Figure 5 ROC Curve

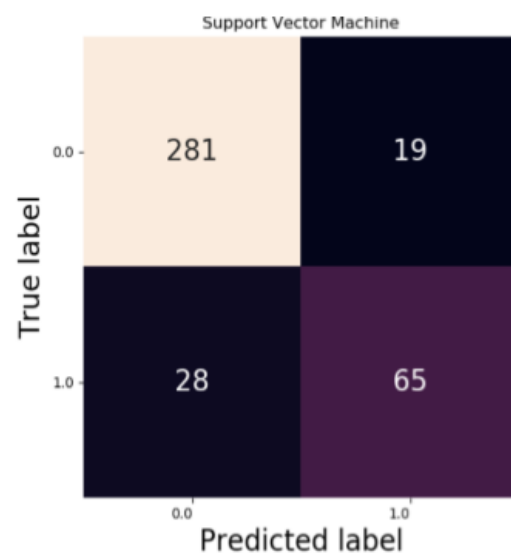


Figure 6 SVM

According to the Figure 5, the AUC=0.9 which means SVM is a good classifier in this case. The AUC value is equivalent to the probability that a randomly chosen positive example is ranked higher than a randomly chosen negative example. ROC Curve were used in this part is because when the distribution of positive and negative samples in the test set changes, the ROC curves can remain unchanged. Class imbalance is often found in the actual data set, that is, the negative sample is much more than the positive sample (or the opposite), and the distribution of the positive and negative samples in the test data may change over time.

From Figure 6, we could indicate that the prediction accuracy of SVM is the same as Decision Tree (Entropy), which is about 87.78%, therefore, we may draw the conclusion that there is no significant difference between the performances of Decision Tree and SVM algorithm in predicting the possibility of survival in shipwreck (Titanic case).

Naïve Bayes

For the last part of research, we intended to find out that whether Naïve Bayes algorithm would perform better in prediction than Decision Tree or SVM, Since Naïve Bayes has lots of advantages than other Machine Learning algorithms:

1. NB is a very simple algorithm which is easy to implement and fast.
2. If the NB conditional independence assumption holds, then it will converge quicker than discriminative models like logistic regression.
3. Even if the NB assumption doesn't hold, it works great in practice.
4. NB only needs less training data.
5. NB is highly scalable. It scales linearly with the number of predictors and data points.
6. It can be used for both binary and multi-class classification problems.
7. It can make probabilistic predictions.
8. It can handle continuous and discrete data.
9. It is not sensitive to irrelevant features.

From the results we got from SVM, as shown in Figure 7:

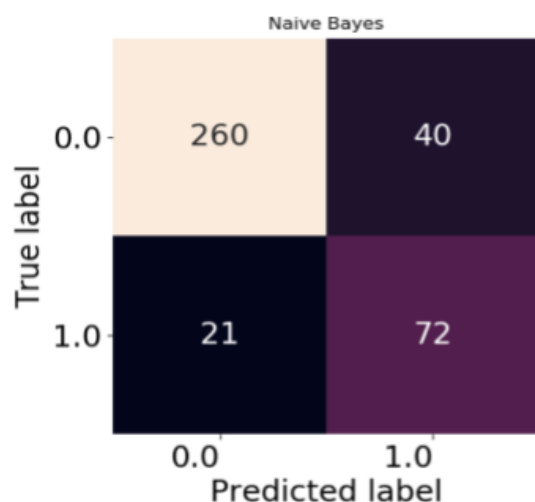


Figure 7 Naïve Bayes

According to the results we got from python in Figure 7, we could find out that the false positive rate is much higher than the other two algorithms (accuracy of prediction is only 84.47%), which means in real life of future, rescue workers would spend much time for rescuing trapped person

who may have died in the shipwreck, and it would absolutely waste much rescuing time and money if researchers doing such prediction by using Naïve Bayes algorithm. Taking all these above into consideration, we may draw the conclusion that although Naïve Bayes has lots of advantages, it is not an appropriate algorithm for prediction in this case.

Calculation

Calculation= $(154-21)/(154+15)*100\%=0.78$

References

- [1] Sapatinas, T. (n.d.). The Elements of Statistical Learning. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*. Oxford, UK: Blackwell Publishing. doi:10.1111/j.1467-985X.2004.298_11.x
- [2] Durbin, R., Miall, C., & Mitchison, G. (1989). The Computing neuron . Wokingham, England ;: Addison-Wesley Pub. Company.
- [3] Raschka, S. (2015). Python Machine Learning. Birmingham: Packt Publishing Ltd.
- [4].<https://www.kaggle.com/ldfreeman3/a-data-science-framework-to-achieve-99-accuracy/notebook>
- [5] <https://www.kaggle.com/c/titanic>