

Compare Support vector machine to three layer Neural Network on Titanic dataset

Rituraj Singh

Bihar, India

Rituraj.nitrkl@gmail.com

Abstract- Titanic disaster occurred 100 years ago on April 15, 1912, killing about 1500 passengers and crew members. With the use of machine learning methods and a dataset provided by Kaggle consisting of 891 rows in the train set and 418 rows in the test set, we attempt to determine the correlation between factors such as age, sex, passenger class, fare etc. to the chance of survival of the passengers. In particular, compare two different machine learning technique SVM and neural network.

1. Introduction

The field of machine learning has allowed analysts to uncover insights from historical data and past events. Titanic disaster is one of the most famous shipwrecks in the world history. 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. Using data provided on <https://www.kaggle.com/c/titanic> our goal is to apply machine-learning techniques to successfully predict which passengers survived the sinking of the Titanic. The method use in this project include SVM and three layer neural network. Various tools are used to implement these algorithms including Python, Pandas, Sklearn, TensorFlow etc.

2. Data set

The data I have used for this project is provided on the Kaggle website. Data consists 891 passenger sample for training set and their associated labels of whether or not the passenger survived. For each passenger, his/her passenger class, name, sex, age, number of siblings/spouses aboard, number of parents/children aboard, ticket number, fare, cabin and embarked are given. In Fig. 1 training dataset sample is given.

PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked	
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	S
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	NaN	S
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C123	S
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	NaN	S

Fig. 1- training data set sample

ATTRIBUTES IN TRAINING DATASET

- Survived: Survival (0 = No, 1 = Yes)
- PassengerId: Id given to each traveler on the boat
- Pclass: Ticket class. It has three possible values: 1,2,3 (first, second and third class)
- Sex : Gender of the passengers (Male or Female)
- Age: Age of the Passengers
- Sibsp : number of siblings and spouses traveling with the passenger
- Parch : number of parents and children traveling with the passenger
- Ticket: Ticket number
- Fare : Passenger fare
- Cabin : Cabin number
- Embarked: Port of Embarkation (C = Cherbourg, Q = Queenstown, S = Southampton)

3. Data Analysis

We need to explore dataset to consider potential data input for the solution. This step is very important because the quality and quantity of data determine how good our model can be. Ticket feature may not be a correlation with survival. It contains 210 duplicates values. We may drop ticket feature. Cabin feature may be dropped as it is highly incomplete or contains many null values in training. PassengerId is not correlated with survival so we may be drop it from training dataset. Out of all passengers in training dataset 38% survived. From Fig.2 we can see that there is significant difference is survival between Female (74.20%) and male (18.89%).

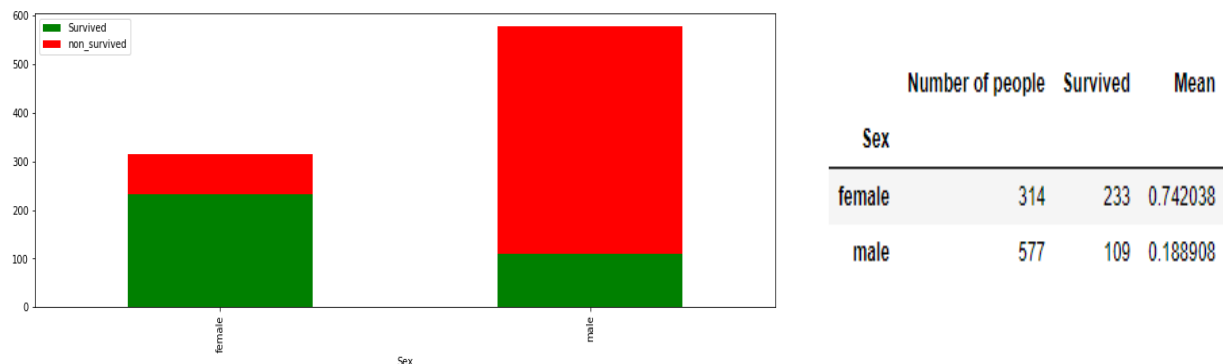


Fig.2 Sex vs Survival

Female from first and second class have more than 90% survival chance and from third class only 50% survived while male has a much higher survival rate (36.88%) from first class then from second class(15.74%) and then from third class(13.54%)(Fig.3).

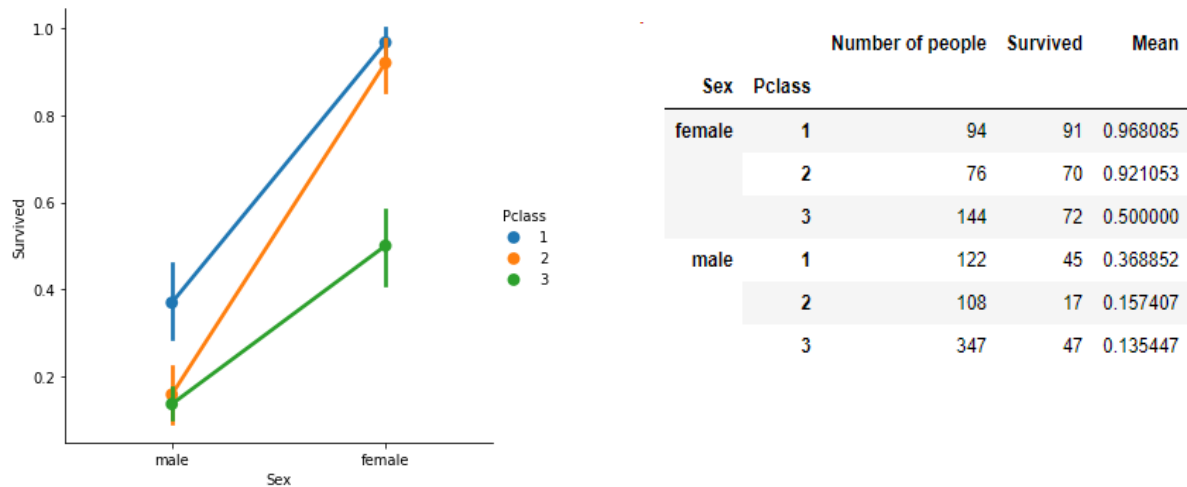


Fig.3 Pclass vs. Survival

From Fig.4 we can see that younger male of age range 5-10 year tend to survive as depicted by green histogram male of age range 20year-40year has more tend to die. Women survived more than men, as depicted by the larger female green histogram.

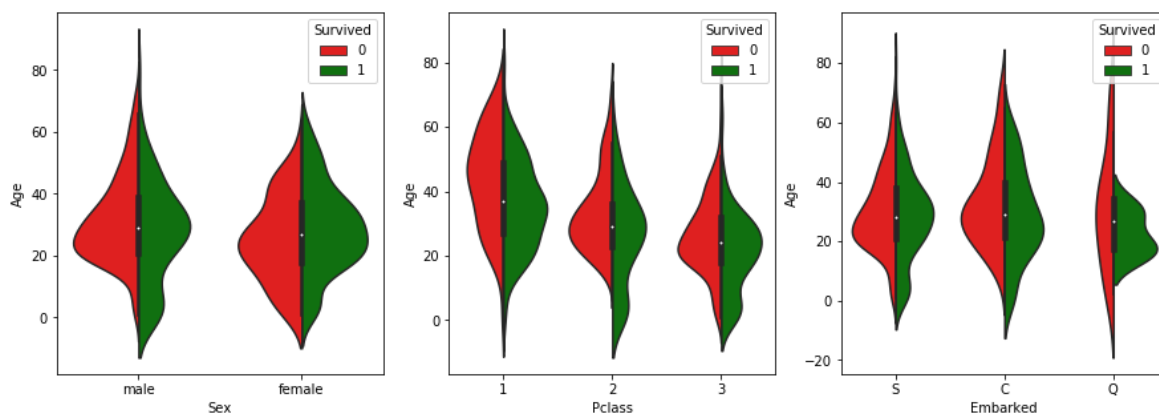
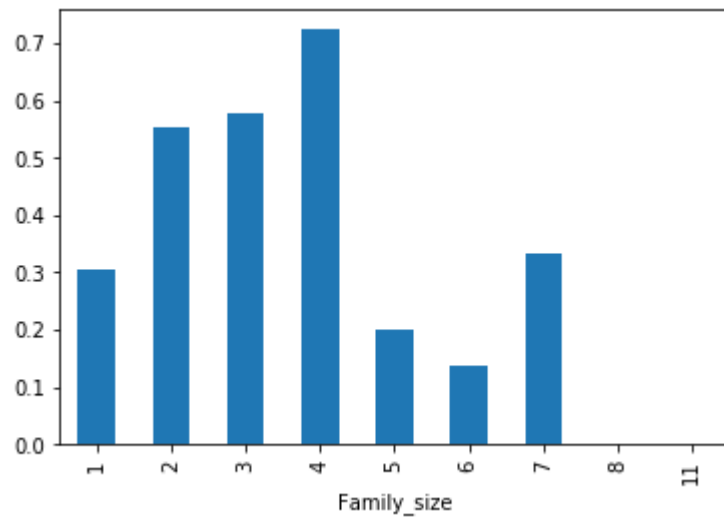


Fig.4 Histogram plot of Age vs. Survival for different features like Sex, Pclass, Embarked

4. Feature Engineering

First I have combined train and test data for feature engineering to prevent any information mismatch in train and test data set. From the name extract title that can give additional information about the social status. There are total 18 title in data, map these title with Mr., Miss., Mrs., Officer, Master and Royalty, we need to convert these feature in binary format. Sex features were mapped to 0 as male and 1 as female. Next feature is age, null value of age has been filled with mean value of given age data. Now I combined the Sibsp and Parch to get family size as we can see family size with 2, 3 and 4 members were high survival chance. Family has been converted into three group as singleton, small family and large family.



	Family_size	Survived
0	1	0.303538
1	2	0.552795
2	3	0.578431
3	4	0.724138
4	5	0.200000
5	6	0.136364
6	7	0.333333
7	8	0.000000
8	11	0.000000

Fig.5 Family Size vs. Survived

Fare features were converted into three group as low, medium and high fare range. In Embarked feature two null values were replaced by most frequent value S.

After completion of feature engineering data were separated into two part first 891 data as train data and rest 418 data as test data. Then engineered feature were selected for modeling.

5. Modeling

Models uses in this project were SVM and neural network model with training data. The following feature were selected for data modeling a.) Pclass, b.) Sex, c.) Age, d.) Embarked, e.) Title, f.) Family, g.) Fare_data. Train data was split into train and test set for modeling purpose.

The SVM model was implemented for classification on train dataset. rbf function was used as kernel in SVM model. We were able to achieve an accuracy rate of 88.88% on test data set.

The neural network build was a three-layer neural network. Two layer with relu function and third layer with sigmoid function. We were able to achieve an accuracy rate of 90.00% on test data set.

Algorithm	Accuracy
SVM	88.88%
Neural Network	87.77%

6. Conclusion

In this project machine learning algorithm SVM and neural network has been successfully Implemented. We also determined the feature that were most the most significant for the

prediction. We were observed that shows SVM higher accuracy rate than neural network model.

7. References

- a. Kaggle, Titanic: Machine Learning form Disaster [Online]. Available: <http://www.kaggle.com/>
- b. Eric Lam, Chongxuan Tang. Titanic – Machine LearningFromDisaster.AvailableFTP: s229.stanford.edu Directory: proj2012
- c. Cortes, Corinna; and Vapnik, Vladimir N.; "Support-Vector Networks", Machine Learning, 20, 1995