

Video Link: <https://www.loom.com/share/da1857f462c34d839be9a2c6b88656ca>

---

- Find the correlation between 'survived' (target column) and 'sex' column for the Titanic use case in class. Do you think we should keep this feature?

The screenshot shows a Jupyter Notebook with a Python script for analyzing the Titanic dataset. The script imports pandas, reads the 'train.csv' file, drops null values, and calculates the correlation between 'survived' and 'sex'. The output shows a correlation of approximately 0.54.

```
1 import pandas as pd
2
3 #Reading the dataset and dropping null values
4 train_df= pd.read_csv('train.csv')
5 print(train_df.columns.values)
6 print(train_df['Survived'].value_counts(dropna=False))
7
8 #Remove Survived from the train set
9 X_train= train_df.drop("Survived",axis=1)
10 Y_train= train_df["Survived"]
11
12 #Training on sex column and categorizing sex
13 train_df['Sex'] = train_df['Sex'].map({'female': 1, 'male': 0}).astype(int)
14
15 #calculating correlation between survived number and sex of passengers
16 print(train_df['Survived'].corr(train_df['Sex']))
17
18 #corelation between sex and survived
19 relation = train_df[["Survived","Sex"]].groupby(['Sex'], as_index = False).mean()
20 print (relation)
21
22
```

Run output:

```
C:\ProgramData\Anaconda3\python.exe "C:/Users/1d630534/Documents/Anurag/2020/Python with DL/ICP4/1_correlation_titanic.py"
['PassengerId' 'Survived' 'Pclass' 'Name' 'Sex' 'Age' 'SibSp' 'Parch' 'Ticket' 'Fare' 'Cabin' 'Embarked']
0    549
1    342
Name: Survived, dtype: int64
0.543351380657755
Sex Survived
0    0    0.188908
1    1    0.742038
Process finished with exit code 0
```

- Implement Naïve Bayes method using scikit-learn library.
  - Use dataset available in <https://umkc.box.com/s/anjic6c8g6034ptm0hgii6fhcu919kx8x>
  - Use train\_test\_split to create training and testing part
  - Evaluate the model on testing part using score and

The screenshot shows a Jupyter Notebook with a Python script for implementing a Naïve Bayes model. The script imports necessary libraries, reads the 'glass.csv' file, splits the data into training and testing sets, fits a Gaussian Naïve Bayes model, and evaluates its performance. The output shows a model accuracy of 46.0.

```
1 # Importing libraries as needed
2 import pandas as pd
3 from sklearn.model_selection import train_test_split
4 from sklearn.naive_bayes import GaussianNB
5 from sklearn.metrics import accuracy_score
6 from sklearn.metrics import classification_report
7
8 # Reading the provided data set
9 train_df = pd.read_csv('glass.csv')
10 X = train_df.drop("Type", axis=1)
11 Y = train_df["Type"]
12
13 # Defining the training and test data sets
14 X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.3, random_state=0)
15
16 # Fitting Naive bayes model
17 model = GaussianNB()
18
19 # Predicting the results of the model on the test data
20 Y_prediction = model.fit(X_train, y_train).predict(X_test)
21 acc_model = round(model.score(X_test, y_test) * 100)
22
23 # Computing the error rate of the model fit
24 print("Model accuracy is:", acc_model)
25 print(classification_report(y_test, Y_prediction))
```

Run output:

```
C:\ProgramData\Anaconda3\python.exe "C:/Users/1d630534/Documents/Anurag/2020/Python with DL/ICP4/2_naive.py"
Model accuracy is: 46.0
precision    recall  f1-score   support

1         0.39         0.86         0.54         21
2         0.50         0.12         0.19         26
3         0.00         0.00         0.00          7
5         0.00         0.00         0.00          2
6         0.67         1.00         0.80          2
7         0.88         1.00         0.93          7

micro avg         0.46         0.46         0.46         65
macro avg         0.41         0.50         0.41         65
weighted avg         0.44         0.46         0.37         65
Process finished with exit code 0
```

### 3. Implement linear SVM method using scikit library

- Use the same dataset above
- Use train\_test\_split to create training and testing part
- Evaluate the model on testing part using score and

```
3 SVM.py x
1 # Importing Libraries as needed
2 import pandas as pd
3 from sklearn.model_selection import train_test_split
4 from sklearn.metrics import classification_report
5 from sklearn.svm import SVC
6
7 # Reading the provided data set
8 # Defining the training and test
9 train_df = pd.read_csv('glass.csv')
10 X = train_df.drop("Type", axis=1)
11 Y = train_df["Type"]
12 data sets
13 X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.3, random_state=0)
14
15 # Fitting Naive bayes model
16 svc = SVC()
17 svc.fit(X_train, y_train)
18 Y_prediction= svc.predict(X_test)
19 acc_svc = round(svc.score(X_train, y_train) * 100,2)
20
21 # Computing the error rate of the model fit
22 print("SVM accuracy is:", acc_svc)
23 print(classification_report(y_test, Y_prediction))
24
```

Run: 3\_SVM x

"avoid this warning.", FutureWarning)  
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\metrics\cl  
UndefinedMetricWarning: Precision and F-score are ill-defined  
labels with no predicted samples.  
'precision', 'predicted', average, warn\_for)  
SVM accuracy is: 77.18

	precision	recall	f1-score	support
1	0.52	0.76	0.62	21
2	0.70	0.62	0.65	26
3	0.00	0.00	0.00	7
5	1.00	1.00	1.00	2
6	0.00	0.00	0.00	2
7	0.88	1.00	0.93	7
micro avg	0.63	0.63	0.63	65
macro avg	0.51	0.56	0.53	65
weighted avg	0.57	0.63	0.59	65

Process finished with exit code 0

Which algorithm you got better accuracy? Can you justify why?

SVM got better accuracy in provided dataset. SVMs are more likely to perform better as they can handle non-linearities in the data. Naive Bayes performs best when the features are independent of each other.

-----