Naïve Bayes

Nguyen Thanh Phat

B2005853

Contents

Naive Bayes: Dataset Golf	1
·	
•	
-	
-	
	Each case of Probability: Using likelihood for each case: Result Naive Bayes Numerical features: Dataset Golf Each case of Probability: Using likelihood for each case: Result

1. Naive Bayes: Dataset Golf

Outlook	Temp	Humidity	Windy	Play
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No

Given dataset Golf with 4 attributes Outlook, Temp, Humidity, Windy and an attribute Play (class).

- How Naïve Bayes predicts the class for 4 examples as follows:

Outlook	Temp	Humidity	Windy	Play
Overcast	Cool	High	False	?
Rainy	Cool	High	False	?
Sunny	Hot	Normal	False	?
???	Hot	Normal	False	?

Each case of Probability:

Out	Outlook			Temperature			Humidity			Windy			Play		
	Yes	No		Yes	No		Yes	No		Yes	No	Yes	No		
Sunny	2	3	Hot	2	2	High	3	4	TRUE	3	3	9	5		
Overcast	4	0	Mild	4	2	Normal	6	1	FALSE	1	2	14			
Rainy	3	2	Cool	3	1										
Sunny	2/9	3/5	Hot	2/9	2/5	High	1/3	4/5	TRUE	1/3	3/5	9/14	5/14		
Overcast	4/9	0	Mild	4/9	2/5	Normal	2/3	1/5	FALSE	2/3	2/5				
Rainy	1/3	2/5	Cool	1/3	1/5										

Using likelihood for each case:

Outlook			Temp		Humidity		Windy			Play				
	Ye	N		Ye	No		Ye			Ye	No		Yes	No
	S	0		S			S	No		S				
Over-			Coo			High			FALS			Ye	0.021164	0.000000
cast	4/9	0	1	1/3	1/5		1/3	4/5	E	2/3	2/5	S	0	0
Rainy		0.	Coo			High			FALS			Ye	0.015873	0.009142
	1/3	4	1	1/3	1/5		1/3	4/5	E	2/3	2/5	S	0	9
Sunny		0.	Hot			Nor-			FALS			Ye	0.014109	0.006857
	2/9	6		2/9	2/5	mal	2/3	1/5	E	2/3	2/5	S	3	1
???			Hot			Nor-			FALS			Ye	0.063492	0.011428
	1	1		2/9	2/5	mal	2/3	1/5	E	2/3	2/5	S	1	6

Result

Outlook	Temp	Humidity	Windy	Play
Overcast	Cool	High	False	Yes
Rainy	Cool	High	False	Yes
Sunny	Hot	Normal	False	Yes
???	Hot	Normal	False	Yes

2. Naive Bayes Numerical features: Dataset Golf

-Naïve Bayes predicts the class for 4 examples as follows:

Outlook	Temp	Humidity	Windy	Play
Overcast	66	80	False	?
Rainy	73	90	False	?
Sunny	80	85	False	?
???	90	85	???	?

Each case of Probability: Firstly, I calculate each case of Probability:

i iistiy,	1 carce	iiaic c	acn case (JI I 100a	omity.									
Outlook			Temperature				Humidity			Windy			Play	
	Yes	No		Yes	No		Yes	No		Yes	No	Yes	No	
Sunny	2	3		83	85		86	85	FALSE	6	2	9	5	
Over- cast	4	0		70	80		96	90	TRUE	3	3	###		
Rainy	3	2		68	65		80	70						
				64	72		65	95						
				69	71		70	91						
				75			80							
				75			70							
				72			90							
				81			75							
Sunny	2/9	1/3	Mean	73	74.6 0	Mea n	79.11	86.2 0	FALSE	2/3	3/5	###	5/14	
Over-			Std.	6.164	7.89	Std.		9.73						
cast	4/9	0	dev.	4	3	dev.	10.216	14	TRUE	1/3	3/5			
Rainy	1/3	2/9												

Using likelihood for each case:

Like previous part, I using excel to calculate the likelihood of each case:

(Yes	No	Temp	Yes	No	Humidity	Yes	No	Windy	Yes	No	Play	Yes	No
st	4/9	0	66	0.03396	0.02792	80	0.03890	0.03347	FALSE	2/3	3/5	Yes	0.000251681	0.00
	1/3	2/9	73	0.06472	0.04952	90	0.02213	0.03799	FALSE	2/3	3/5	Yes	0.000204575	0.00
,	2/9	1/3	80	0.03396	0.04000	85	0.03307	0.04068	FALSE	2/3	3/5	No	0.000106981	0.00
	1	1	90	0.00144	0.00753	85	0.03307	0.04068	???	1	1	No	0.000030701	0.00

Result

Outlook	Temp	Humidity	Windy	Play
Overcast	66	80	False	Yes
Rainy	73	90	False	Yes
Sunny	80	85	False	No

222	0.0	0 =	222	
1,1,1,1	()()	05	,,,,,	N ₁
1 1 1	90	0.1	111	1 1 1 ()
• • •	70	0.0	• • •	110

3. Implement the program using **GaussianNB** in **scikit-learn** library.

The program requires 2 parameters:

- file name of trainset
- file name of testset

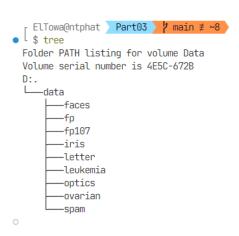
The program reports the classification results (accuracy, confusion matrix) for 6 datasets:

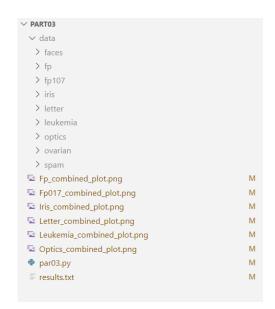
- Iris (.trn: trainset, .tst: testset)
- Optics (.trn: trainset, .tst: testset)
- Letter (.trn: trainset, .tst: testset)
- Leukemia (.trn: trainset, .tst: testset)
- Fp (.trn: trainset, .tst: testset)

In this report, I evaluated the performance of a Gaussian Naive Bayes classifier on five different datasets: Iris, Optics, Letter, Leukemia, and Fp. For each dataset, we trained the classifier for 10 epochs and analyzed its performance on the test set.

Source code:

The directory for Part 3:





```
import numpy as np
from sklearn.naive_bayes import GaussianNB
from sklearn.metrics import accuracy_score, confusion_matrix
import os
import matplotlib pyplot as plt
import seaborn as sns
def load_data(filename):
   try:
       data = np.loadtxt(filename, delimiter=",", dtype=float)
       data = np.loadtxt(filename, delimiter=" ", dtype=float)
   X = data[:, :-1]
   y = data[:, -1].astype(int)
   return X, y
def export_confusion_matrix(confusion, dataset_name):
   sns.heatmap(confusion, annot=True, fmt="d", cmap="Blues", cbar=False)
   plt.xlabel("Predicted Label")
   plt.ylabel("True Label")
   plt.title(f"Confusion Matrix - {dataset_name}")
   plt.tight_layout()
   plt.savefig(f"{dataset_name}_combined_plot.png")
   plt.close()
def print_confusion_matrix(confusion):
   for row in confusion:
       print(row)
def save_results_to_file(accuracy, confusion, dataset_name):
   with open("results.txt", "a") as f:
       f.write(f"Dataset: {dataset_name}\n")
       f.write(f"Accuracy: {accuracy}\n")
       f.write("Confusion Matrix:\n")
       np.savetxt(f, confusion, fmt="%d")
def test_model(clf, X_test, y_test):
   y_pred = clf.predict(X_test)
   # Calculate accuracy
   accuracy = accuracy_score(y_test, y_pred)
   return accuracy, y_pred
```

```
def main(trainset_filename, testset_filename, dataset_name=""):
   # Load train and test data
   X_train, y_train = load_data(trainset_filename)
   X_test, y_test = load_data(testset_filename)
   # Initialize Gaussian Naive Bayes classifier
   clf = GaussianNB()
   # Train classifier
   clf.fit(X_train, y_train)
   # Predict on test
   train_accuracy, _ = test_model(clf, X_train, y_train)
   accuracy, y_pred = test_model(clf, X_test, y_test)
   print("Train Accuracy:", train_accuracy)
   print("Test Accuracy:", accuracy)
   # Calculate confusion matrix
   confusion = confusion_matrix(y_test, y_pred)
   print("Confusion Matrix:")
   print_confusion_matrix(confusion)
   # Save into file
   export_confusion_matrix(confusion, dataset_name)
   save_results_to_file(accuracy, confusion, dataset_name)
if __name__ == "__main__":
   datasets = [
           "name": "Iris",
           "train_file": "data//iris//iris.trn",
           "test_file": "data//iris//iris.tst",
       },
           "name": "Optics",
           "train_file": "data//optics//optics.trn",
           "test_file": "data//optics//optics.tst",
       },
           "name": "Letter",
           "train_file": "data//letter//letter.trn",
           "test_file": "data//letter//letter.tst",
       },
           "name": "Leukemia".
           "train_file": "data//leukemia//leukemia.trn",
           "test_file": "data//leukemia//leukemia.tst",
```

```
{
    "name": "Fp",
    "train_file": "data//fp//fp.trn",
    "test_file": "data//fp//fp.tst",
},
{
    "name": "Fp017",
    "train_file": "data//fp107//fp107.trn",
    "test_file": "data//fp107//fp107.tst",
},
]

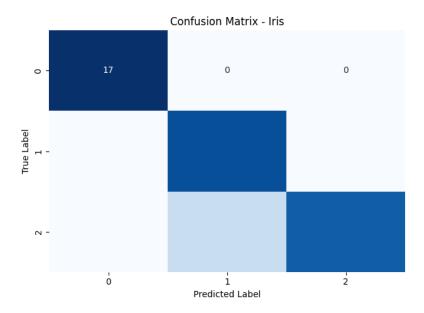
for dataset in datasets:
    print(f"Dataset: {dataset['name']}")
    trainset_path = os.path.join(dataset["train_file"])
    testset_path = os.path.join(dataset["test_file"])
    main(trainset_path, testset_path, dataset["name"])
    print("\n")
```

Dataset: Iris

• Training Accuracy: 98%

• Test Accuracy: 92%

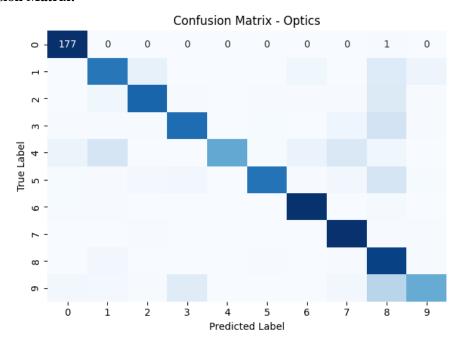
• Confusion Matrix:



Dataset: Optics

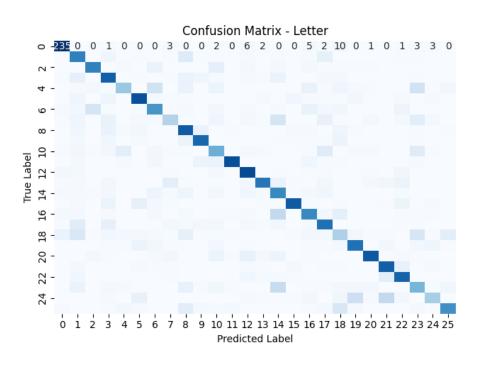
• Training Accuracy: 81.45%

- Test Accuracy: 78.63%
- Confusion Matrix:



Dataset: Letter

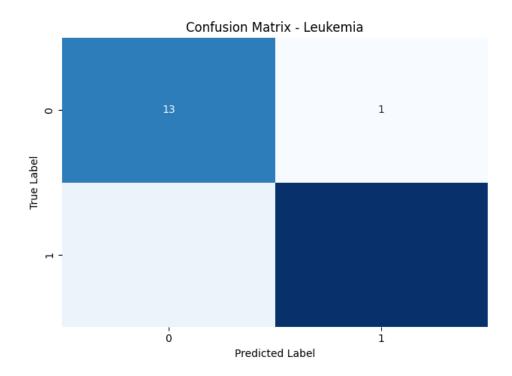
- Training Accuracy: 64.81%
- **Test Accuracy:** 63.16%
- Confusion Matrix:



Dataset: Leukemia

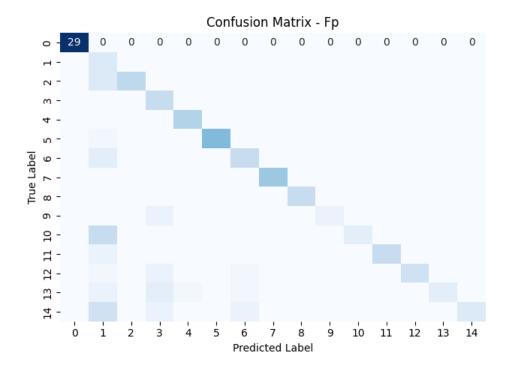
• Training Accuracy: 100%

- **Test Accuracy:** 91.18% **Confusion Matrix:**



Dataset: Fp

- Training Accuracy: 100% Test Accuracy: 75%
- **Confusion Matrix:**



Dataset: Fp107

Training Accuracy: 100%Test Accuracy: 96.816%

• Confusion Matrix:

