

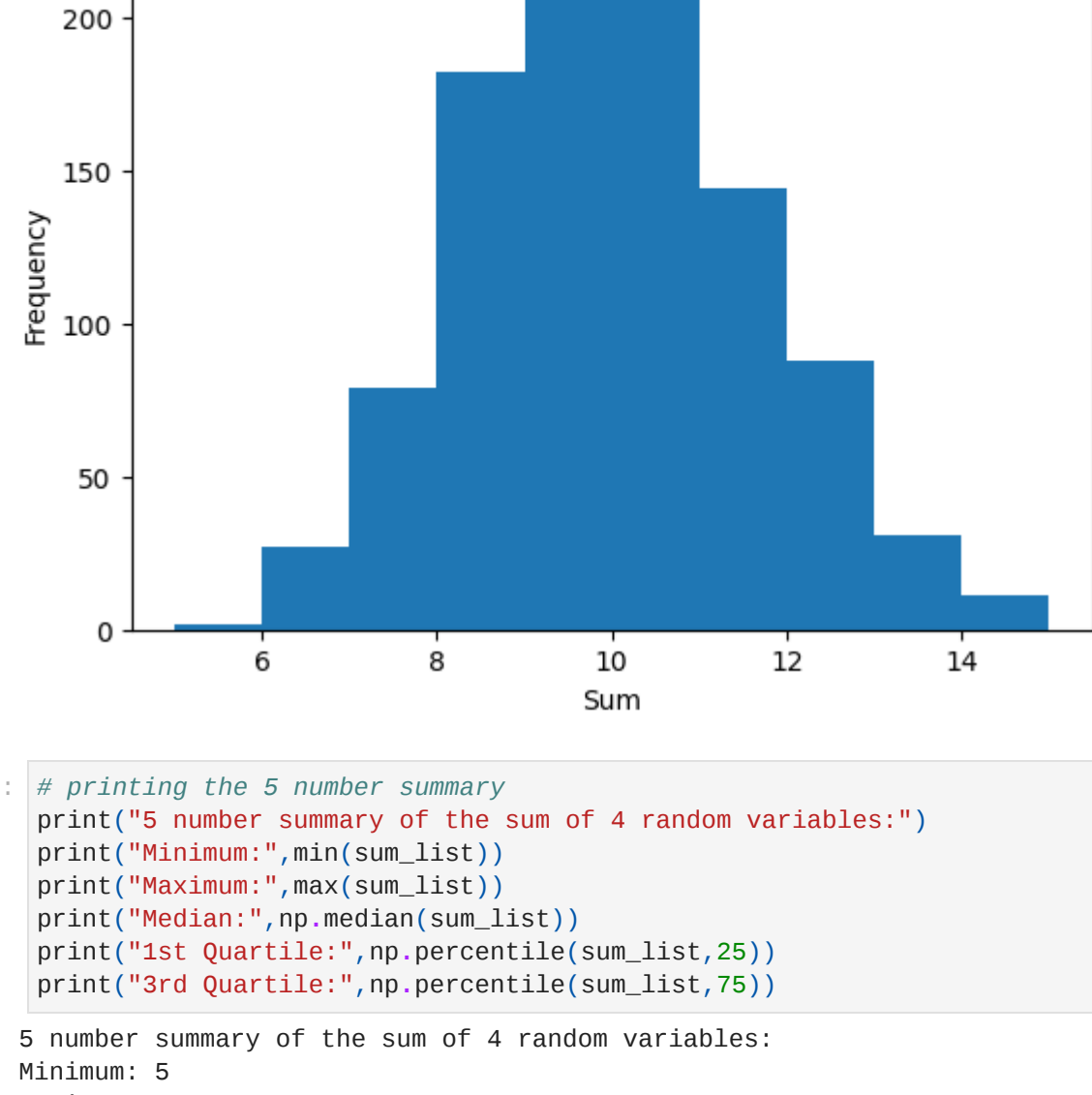
```
In [ ]: import random
import matplotlib.pyplot as plt
import numpy as np
from tabulate import tabulate
```

## Naive Bayes

### Part A

```
In [ ]: p=[0.125,0.5,0.25,0.125]
sum_list=[]
for i in range(9,1000):
    x=random.choices([1,2,3,4],weights=p,k=4)
    sum_list.append(sum(x))
```

```
In [ ]: # plotting the frequency distribution histogram
plt.hist(sum_list,label='Sum of 4 random variables')
plt.xlabel('Sum')
plt.ylabel('Frequency')
plt.show()
```



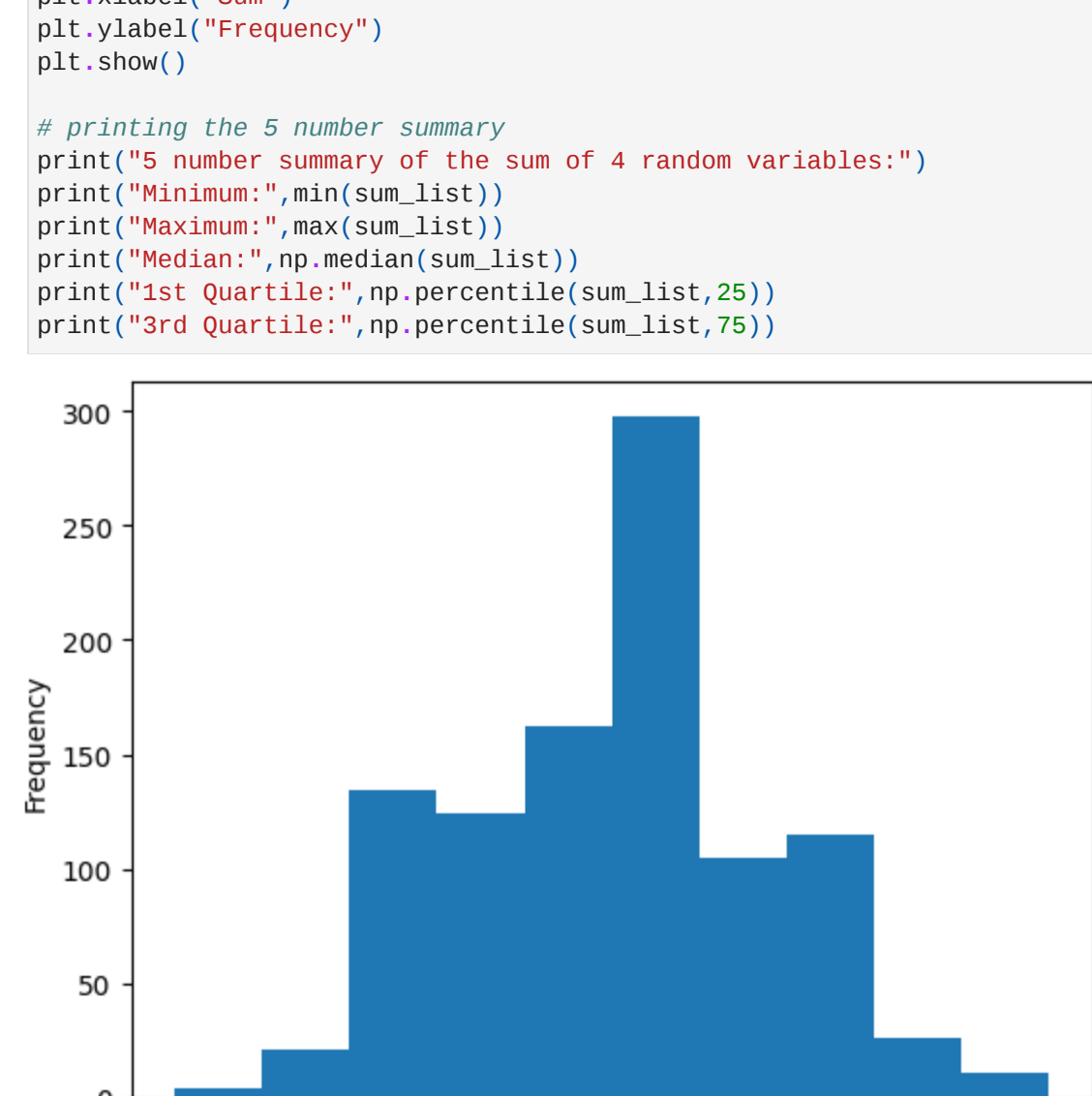
```
In [ ]: # printing the 5 number summary
print("% number summary of the sum of 4 random variables:")
print("Minimum:",min(sum_list))
print("Maximum:",max(sum_list))
print("Median:",np.median(sum_list))
print("1st Quartile:",np.percentile(sum_list,25))
print("3rd Quartile:",np.percentile(sum_list,75))
5 number summary of the sum of 4 random variables:
Minimum: 5
Maximum: 16
Median: 9.0
1st Quartile: 8.0
3rd Quartile: 11.0
Expected Value of dice roll: E[X] = 1 1/2 + 2 1/2 + 3 1/2 + 4 1/2 = 2.375
```

The theoretical Expected Sum ,  $E[X] = E[X_1 + X_2 + X_3 + X_4] = E[X_1] + E[X_2] + E[X_3] + E[X_4] = 4 * E[X_1] = 9.5$

```
In [ ]: Mean = np.mean(sum_list)
print("Mean/Expected sum from the python simulation :",Mean)
Mean/Expected sum from the python simulation : 9.517
So Theoretical Expected sum is close to the Expected sum from the simulation
```

Next we try for k = 4 and randomly roll the die 8 times and calculate the sum of the upward face value.

```
In [ ]: p=[0.125,0.5,0.25,0.125]
sum_list=[]
for i in range(9,1000):
    x=random.choices([1,2,3,4],weights=p,k=8)
    sum_list.append(sum(x))
# plotting the histogram between 4 and 16
plt.hist(sum_list,label='Sum of 4 random variables')
plt.xlabel('Sum')
plt.ylabel('Frequency')
plt.show()
# printing the 5 number summary
print("% number summary of the sum of 4 random variables:")
print("Minimum:",min(sum_list))
print("Maximum:",max(sum_list))
print("Median:",np.median(sum_list))
print("1st Quartile:",np.percentile(sum_list,25))
print("3rd Quartile:",np.percentile(sum_list,75))
```



5 number summary of the sum of 4 random variables:

Minimum: 12  
Maximum: 26  
Median: 18.0  
1st Quartile: 17.0  
3rd Quartile: 21.0

Expected Value of dice roll:  $E[X] = 1 1/2 + 2 1/2 + 3 1/2 + 4 1/2 = 2.375$

The theoretical Expected Sum ,  $E[X] = E[X_1 + X_2 + X_3 + X_4] = E[X_1] + E[X_2] + E[X_3] + E[X_4] = 4 * E[X_1] = 19$

```
In [ ]: Mean = np.mean(sum_list)
print("Mean/Expected sum from the python simulation :",Mean)
Mean/Expected sum from the python simulation : 18.974
So Theoretical Expected sum is close to the Expected sum from the simulation
```

Next we do the same for k=16

```
In [ ]: p=[1/(2**(i-1)) for i in range(2,17)]
p.insert(0,1/(2**(15)))
sum_list=[]
for i in range(9,1000):
    x=random.choices([1 for i in range(1,17)],weights=p,k=4)
    sum_list.append(sum(x))
# plotting the histogram between 4 and 16
plt.hist(sum_list,label='Sum of 4 random variables')
plt.xlabel('Sum')
plt.ylabel('Frequency')
plt.show()
# printing the 5 number summary
print("% number summary of the sum of 4 random variables:")
print("Minimum:",min(sum_list))
print("Maximum:",max(sum_list))
print("Median:",np.median(sum_list))
print("1st Quartile:",np.percentile(sum_list,25))
print("3rd Quartile:",np.percentile(sum_list,75))
```



5 number summary of the sum of 4 random variables:

Minimum: 8  
Maximum: 25  
Median: 12.0  
1st Quartile: 10.0  
3rd Quartile: 14.0

```
In [ ]: print("Mean/Expected sum from the python simulation :",np.mean(sum_list))
print("Theoretical expected sum:",np.round(np.sum([1*p[i-1] for i in range(1,17)])**4,5))
Mean/Expected sum from the python simulation : 12.177
Theoretical expected sum: 11.99792
```

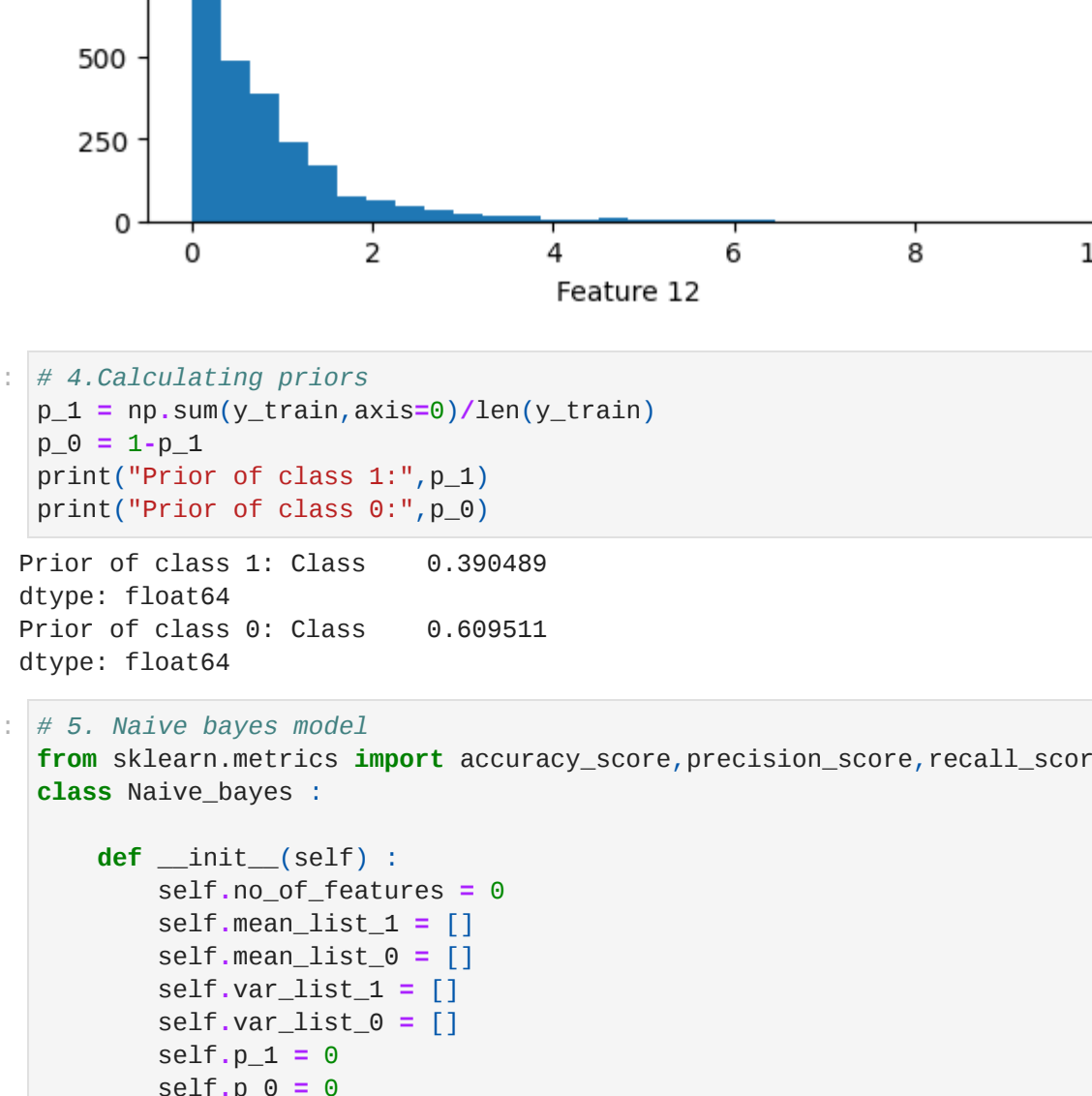
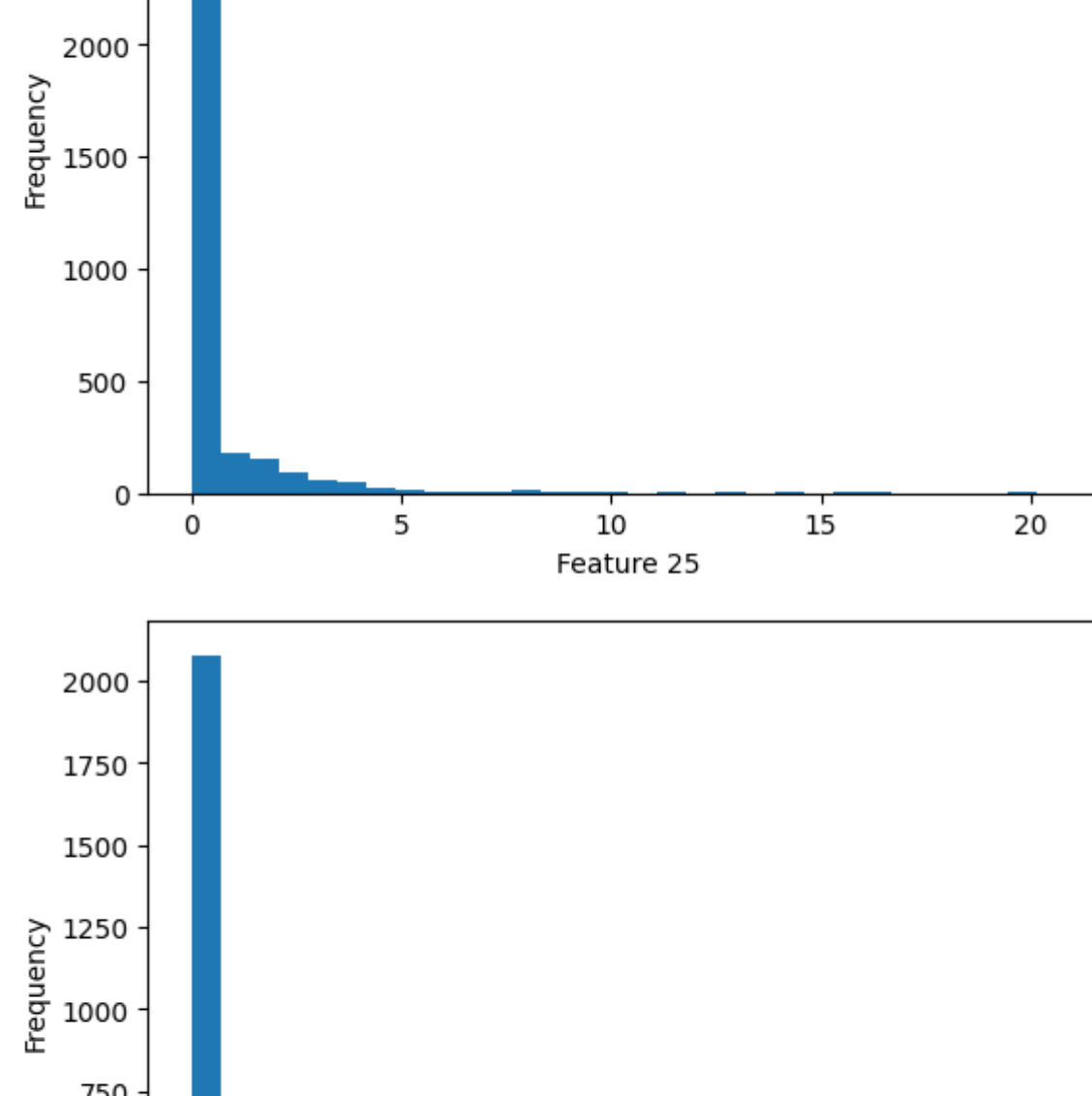
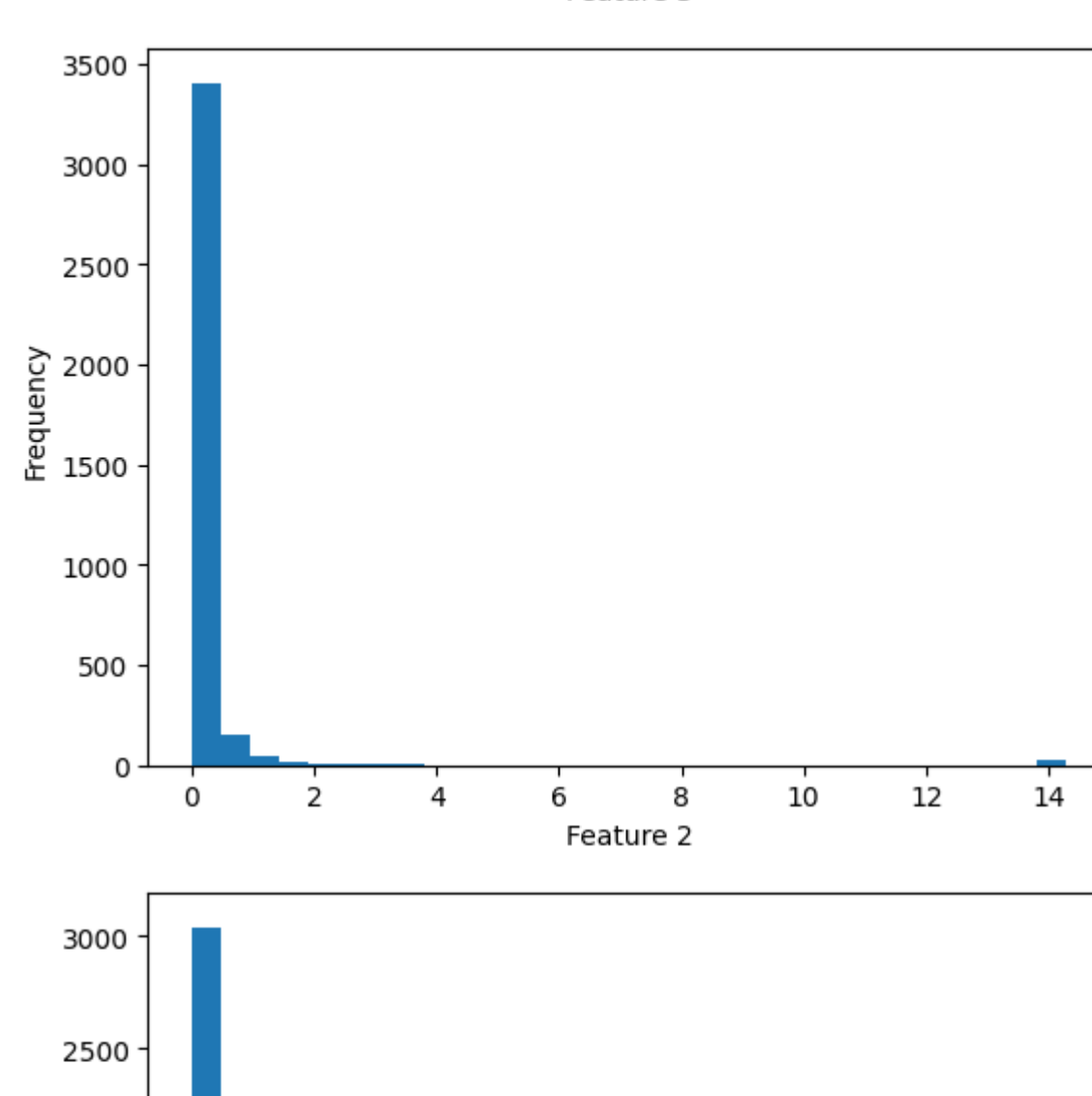
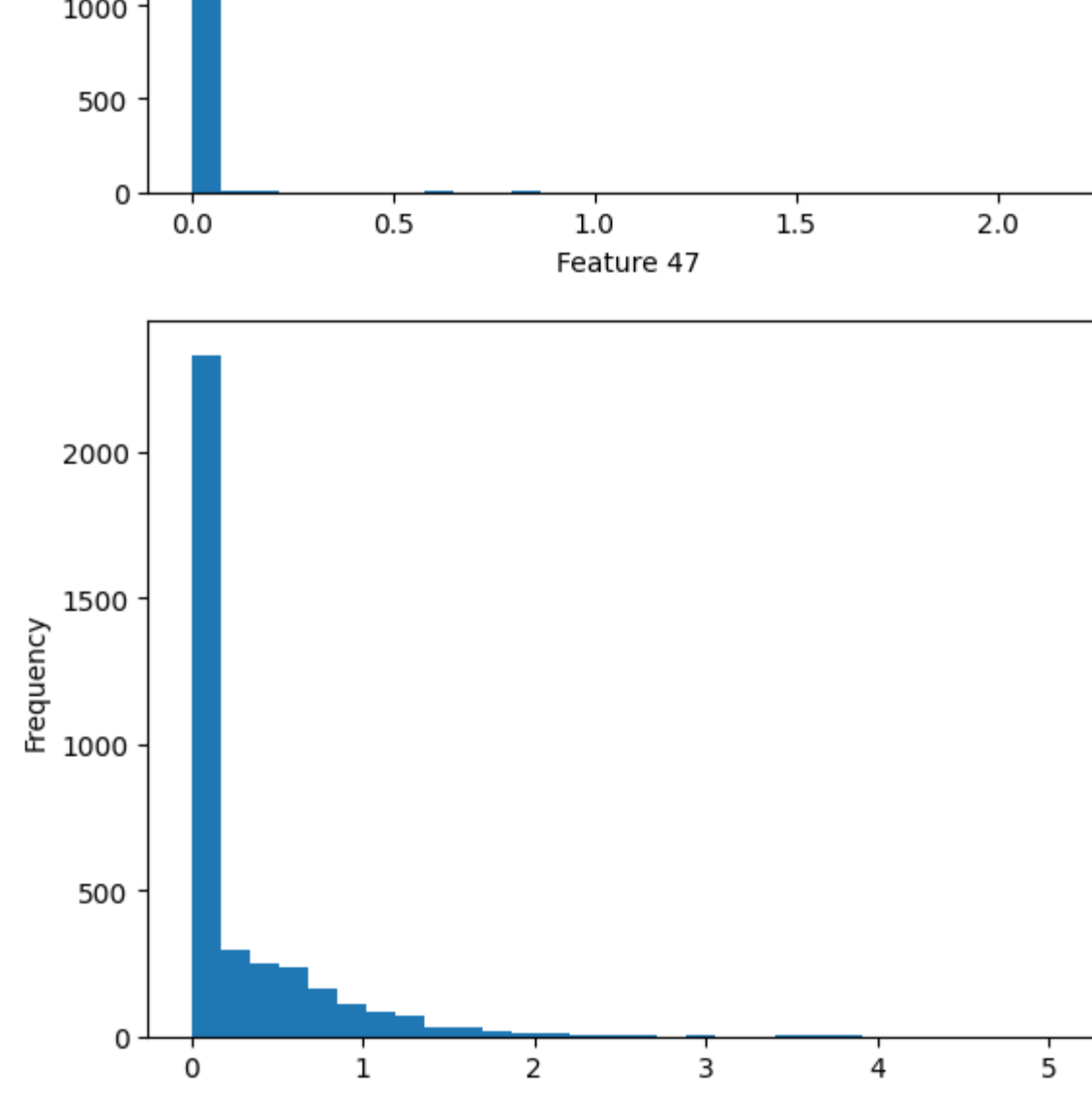
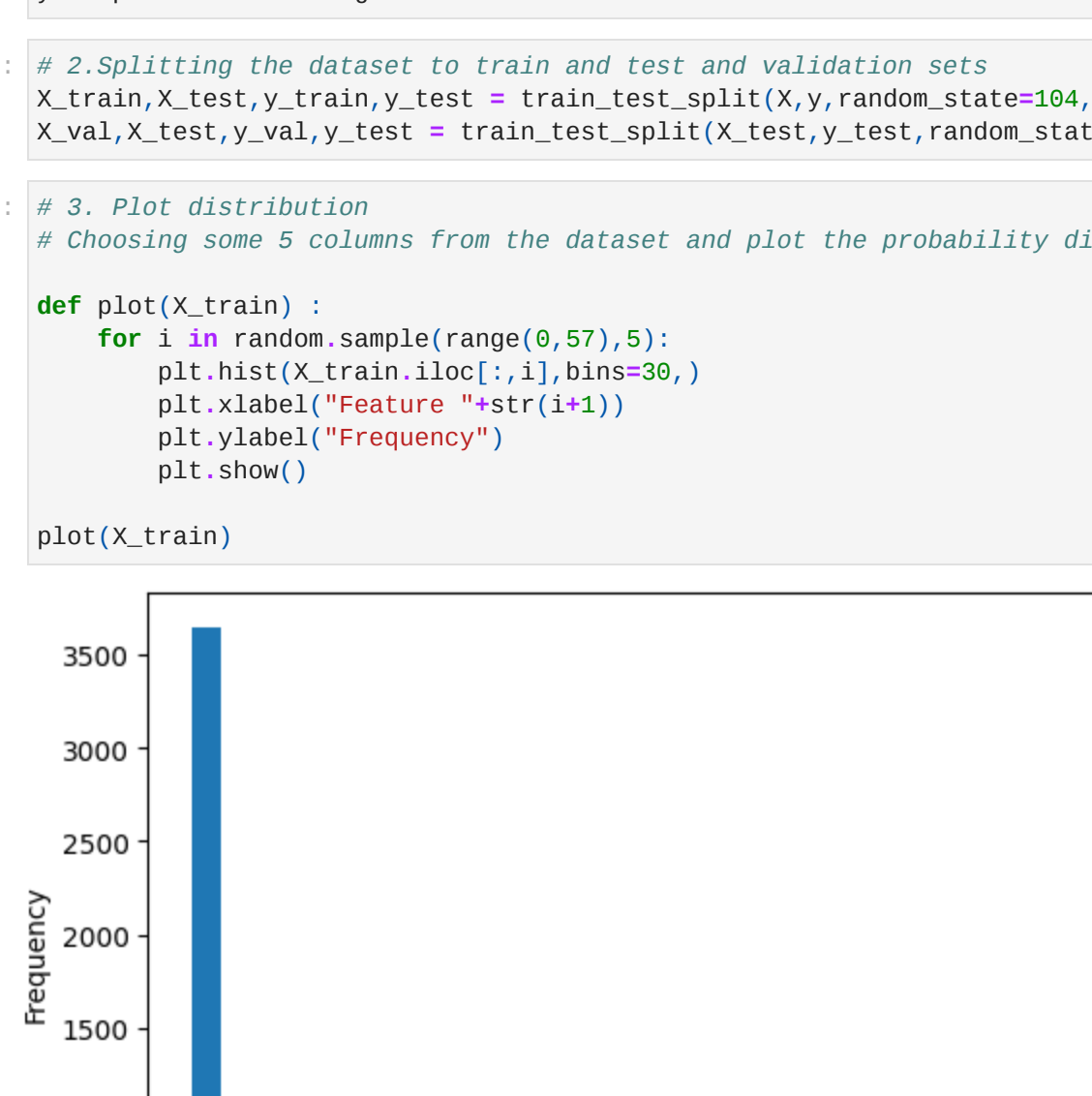
### Part B : Implementation of Naive Bayes

```
In [ ]: import pandas as pd
import sklearn
from ucmlrepo import fetch_ucmlrepo
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
```

```
In [ ]: # 1. fetch dataset
spambase = fetch_ucmlrepo(id=94)
# data (as pandas dataframes)
X = spambase.data.features
y = spambase.data.targets
```

```
In [ ]: # 2. Splitting the dataset to train and test and validation sets
X_train,X_test,y_train,y_test = train_test_split(X,y,random_state=94,train_size=0.8,shuffle=True)
X_val,X_test,y_val,y_test = train_test_split(X_test,y_test,random_state=184,train_size=0.5,shuffle=True)
```

```
In [ ]: # 3. Plot distribution
# Choosing some 5 columns from the dataset and plot the probability distribution.
def plot(X,train):
    for i in range(sample(range(0,87),5)):
        plt.hist(X_train.iloc[:,i],bins=20)
        plt.xlabel('Feature '+str(i+1))
        plt.ylabel('Frequency')
        plt.show()
plot(X,train)
```



```
In [ ]: # 4. Calculating priors
p_1 = np.sum(y_train,axis=0)/len(y_train)
p_0 = 1-p_1
print("Prior of class 1:",p_1)
print("Prior of class 0:",p_0)
```

Prior of class 1: Class 0.396489  
dtype: float64  
Prior of class 0: Class 0.603511  
dtype: float64

```
In [ ]: # 5. Naive Bayes model
from sklearn.metrics import accuracy_score,precision_score,recall_score,f1_score
class Naive_bayes:
```

```
    def __init__(self):
        self.no_of_features = 0
        self.mean_list_1 = []
        self.mean_list_0 = []
        self.var_list_1 = []
        self.var_list_0 = []
        self.p_1 = 0
        self.p_0 = 0
        self.y_pred = []
        self.no_of_parameters = 0
```

```
    def fit(self,X_train,y_train):
        self.no_of_features = X_train.shape[1]
        sum1 = np.sum(X_train,axis=0)
        self.p_1 = sum1/y_train.shape[0]
        self.p_0 = 1-self.p_1
        for i in range(self.no_of_features):
            self.mean_list_1.append(np.mean(X_train.loc[(y_train['Class']==1),X_train.columns[i]]))
            self.mean_list_0.append(np.mean(X_train.loc[(y_train['Class']==0),X_train.columns[i]]))
            self.var_list_1.append(np.var(X_train.loc[(y_train['Class']==1),X_train.columns[i]])
```

```
            self.var_list_0.append(np.var(X_train.loc[(y_train['Class']==0),X_train.columns[i]])
        self.mean_list_1 = np.array(self.mean_list_1)
        self.mean_list_0 = np.array(self.mean_list_0)
        self.var_list_1 = np.array(self.var_list_1)
        self.var_list_0 = np.array(self.var_list_0)
        self.no_of_parameters = 2*self.no_of_features + 2
```

```
    def predict(self,X_test):
        for i in range(X_test.shape[0]):
            p1 = np.exp(-1/2*np.sum((self.var_list_1)**-1*(X_test.iloc[i,:]-self.mean_list_1)**2)/(2*self.var_list_1))
            p0 = np.exp(-1/2*np.sum((self.var_list_0)**-1*(X_test.iloc[i,:]-self.mean_list_0)**2)/(2*self.var_list_0))
            if (self.p_1 > p0*self.p_0).all():
                self.y_pred.append(1)
            else:
                self.y_pred.append(0)
        # return np array
        return self.y_pred
```

```
    def accuracy(self,y_test):
        return accuracy_score(y_test,self.y_pred)
    def precision(self,y_test):
        return precision_score(y_test,self.y_pred)
    def recall(self,y_test):
        return recall_score(y_test,self.y_pred)
    def f1_score(self,y_test):
        return f1_score(y_test,self.y_pred)
```

```
In [ ]: # naive bayes model without log transformation
nb = Naive_bayes()
nb.fit(X_train,y_train)
y_pred = nb.predict(X_test)
print("For the naive bayes model without log transformation:")
print("Accuracy:",nb.accuracy(y_test))
print("Precision:",nb.precision(y_test))
print("Recall:",nb.recall(y_test))
print("F1 score:",nb.f1_score(y_test))
```

For the naive bayes model without log transformation:  
Accuracy: 0.83731892277657  
Precision: 0.76135932388664  
Recall: 0.92158927458969  
F1 score: 0.837328824833782

Accuracy: 0.83731892277657  
Precision: 0.76135932388664  
Recall: 0.92158927458969  
F1 score: 0.837328824833782

```
In [ ]: # naive bayes model with log transformation
X_train_log = np.log(X_train+0.1)
X_test_log = np.log(X_test+0.1)
nb_log = Naive_bayes()
nb_log.fit(X_train_log,y_train)
nb_log.predict(X_test_log)
```

For the naive bayes model with log transformation:  
Accuracy: 0.843817787418655  
Precision: 0.77731892277657  
Recall: 0.96869274589692  
F1 score: 0.83734672381195

### 8. Discussion

The changes that we notice in the scores from the two datasets (with and without log transformation) are very similar . However we notice a slight increase in the accuracy after applying the log transformation . On the other hand , there is a slight dip in recall value .

So though the difference is small , we can say that model with log transformation would generate better , however the higher recall of the model without log transformation would imply that there is smaller chance that this model will classify a non-spam mail as spam when compared to model with log transformation

### Part C : sklearn implementation of Naive Bayes

```
In [ ]: # 2. Training the model with and without log transformation
from sklearn.naive_bayes import GaussianNB
from sklearn.metrics import accuracy_score
nb = GaussianNB()
nb.fit(X_train,y_train.values.ravel())
y_pred = nb.predict(X_test)
```

# accuracy  
print("Accuracy of the model:",accuracy\_score(y\_test,y\_pred))  
Accuracy of the model: 0.8394793926247288

```
In [ ]: # applying log transformation on the data
X_train_log = np.log(X_train+0.1)
X_test_log = np.log(X_test+0.1)
nb_log = GaussianNB()
nb_log.fit(X_train_log,y_train.values.ravel())
y_pred = nb_log.predict(X_test_log)
```

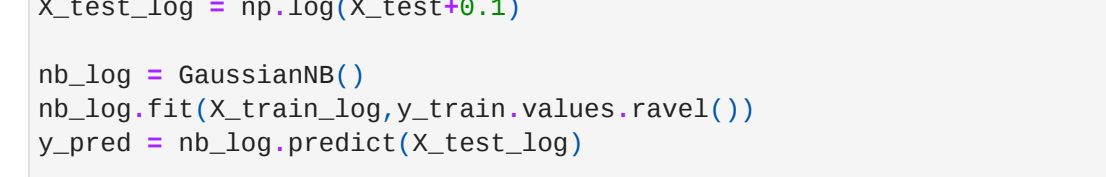
# accuracy  
print("Accuracy of the model after log transformation:",accuracy\_score(y\_test,y\_pred))  
Accuracy of the model after log transformation: 0.843817787418655

```
In [ ]: # plot ROC curve for both the models
from sklearn.metrics import roc_curve
from sklearn.metrics import roc_auc_score
from matplotlib import pyplot
```

# predict probabilities  
probs1 = nb.predict\_proba(X\_test)  
# keep probabilities for the positive outcome only  
probs1 = probs1[:, 1]  
# calculate AUC  
auc1 = roc\_auc\_score(y\_test, probs1)  
print("AUC for the model without log transformation: %.3f" % auc1)  
# calculate roc curve  
fpr, tpr, thresholds = roc\_curve(y\_test, probs1)  
# plot the roc curve for the model  
pyplot.plot(fpr, tpr, marker='.',label="ROC curve without log transformation")  
# show the plot  
pyplot.xlabel("False Positive Rate")  
pyplot.ylabel("True Positive Rate")  
pyplot.legend()

# predict probabilities  
probs = nb\_log.predict\_proba(X\_test\_log)  
# keep probabilities for the positive outcome only  
probs = probs[:, 1]  
# calculate AUC  
auc = roc\_auc\_score(y\_test, probs)  
print("AUC for the model with log transformation: %.3f" % auc)  
# calculate roc curve  
fpr, tpr, thresholds = roc\_curve(y\_test, probs)  
# plot the roc curve for the model  
pyplot.plot(fpr, tpr, marker='.',label="ROC curve with log transformation")  
# show the plot  
pyplot.legend()

AUC for the model without log transformation: 0.941  
AUC for the model with log transformation: 0.948



We look at the two curves and prefer the one with log transformation because it has the higher auc score

```
In [ ]: from sklearn.svm import SVC
# normalise X
scaler = StandardScaler()
X = scaler.fit_transform(X)
```

# 2. Splitting the dataset to train and test and validation sets  
X\_train,X\_test,y\_train,y\_test = train\_test\_split(X,y,random\_state=94,train\_size=0.8,shuffle=True)  
X\_val,X\_test,y\_val,y\_test = train\_test\_split(X\_test,y\_test,random\_state=184,train\_size=0.5,shuffle=True)

# train the model  
svm = SVC(kernel='linear')  
svm.fit(X\_train,y\_train.values.ravel())  
# predict  
y\_pred = svm.predict(X\_test)

```
In [ ]: # tabulate the accuracy , precision , recall and f1 score using tabulate
print("The SVM model : ")
print(tabulate(["Accuracy",accuracy_score(y_test,y_pred)],["Precision",precision_score(y_test,y_pred)],["Recall",recall_score(y_test,y_pred)],["F1 score",f1_score(y_test,y_pred)]))
```

# tabulate the accuracy , precision , recall and f1 for naive bayes model with log transformation using tabulate  
print("\n\nThe Naive Bayes model with log transformation : ")  
print("\n\n")

The SVM model :  
Metrics Value  
Accuracy 0.91974  
Precision 0.94824  
Recall 0.98747  
F1 score 0.96571

The Naive Bayes model with log transformation :  
Metric With log transformation  
Accuracy 0.843818  
Precision 0.777319  
Recall 0.968693  
F1 score 0.873104

We notice that the SVM model is able to better fit the data , and is more likely to correctly classify a mail as spam or not spam . The higher accuracy of the SVM model over the Naive Bayes helps us to deduce this . However we must also look at the recall values of the two models as well , the naive bayes has a better recall than the SVM model . Infact this seems like a more important measure than accuracy , as we wouldnt want non spam emails to be classified as spam , even though some spam emails be classified as non-spam

