

## Question 1: Review the Titanic data(60)

In lecture 3, we use Titanic data as the example for data pre-processing. This dataset contains information about passengers on the Titanic, including features like age, gender, class, and whether they survived or not. Now we are going to fit this data to the three classification models we have discussed.

```
In [1]: 1 import pandas as pd
        2 import seaborn as sns
        3
        4 # Load the Titanic dataset
        5 titanic_data = sns.load_dataset('titanic')
        6
        7 print(titanic_data.head())
```

	survived	pclass	sex	age	sibsp	parch	fare	embarked	class	\
0	0	3	male	22.0	1	0	7.2500	S	Third	
1	1	1	female	38.0	1	0	71.2833	C	First	
2	1	3	female	26.0	0	0	7.9250	S	Third	
3	1	1	female	35.0	1	0	53.1000	S	First	
4	0	3	male	35.0	0	0	8.0500	S	Third	

	who	adult_male	deck	embark_town	alive	alone
0	man	True	NaN	Southampton	no	False
1	woman	False	C	Cherbourg	yes	False
2	woman	False	NaN	Southampton	yes	True
3	woman	False	C	Southampton	yes	False
4	man	True	NaN	Southampton	no	True

1. Perform the pre-processing steps we have done in the lecture 3, including cleaning the missing values, convert the target (survived) to a categorical variable and split the training and testing data. (10)

```
In [2]: 1 # Find the rows that have missing data
2 row_with_missing = titanic_data[titanic_data.isnull().any(axis =1)]
3 row_with_missing
```

Out[2]:

	survived	pclass	sex	age	sibsp	parch	fare	embarked	class	who	adult_male	deck
0	0	3	male	22.0	1	0	7.2500	S	Third	man	True	NaN
2	1	3	female	26.0	0	0	7.9250	S	Third	woman	False	NaN
4	0	3	male	35.0	0	0	8.0500	S	Third	man	True	NaN
5	0	3	male	NaN	0	0	8.4583	Q	Third	man	True	NaN
7	0	3	male	2.0	3	1	21.0750	S	Third	child	False	NaN
...	...	...	...	...	...	...	...	...	...	...	...	...
884	0	3	male	25.0	0	0	7.0500	S	Third	man	True	NaN
885	0	3	female	39.0	0	5	29.1250	Q	Third	woman	False	NaN
886	0	2	male	27.0	0	0	13.0000	S	Second	man	True	NaN
888	0	3	female	NaN	1	2	23.4500	S	Third	woman	False	NaN
890	0	3	male	32.0	0	0	7.7500	Q	Third	man	True	NaN

709 rows × 15 columns

```
In [3]: 1 # Convert the 'deck' column to a string data type.
2 titanic_data['deck'] = titanic_data['deck'].astype(str)
3
4 # Replace the string 'nan' in the 'deck' column with 'Unknown'.
5 titanic_data['deck'] = titanic_data['deck'].replace('nan', 'Unknown')
6
7 # Find rows with missing values (NaN) in any column and store them in 'row_
8 row_with_missing = titanic_data[titanic_data.isnull().any(axis=1)]
9
10 # Remove the 'alive' column from the 'titanic_data' DataFrame.
11 titanic_data = titanic_data.drop('alive', axis=1)
12
13 # Drop rows with any missing values (NaN) from the 'titanic_data' DataFrame
14 titanic_data = titanic_data.dropna()
```

```
In [4]: 1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.model_selection import train_test_split
4
5 X = titanic_data.drop('survived', axis=1)
6 y = titanic_data['survived']
7
8 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, r
```

2. Now only use the age and fare as the features, fit Naive bayes, LDA and QDA model. Report the classification table for each model. Which one performs the best? (20)

```
In [5]: 1 from sklearn.naive_bayes import GaussianNB
2 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
3 from sklearn.discriminant_analysis import QuadraticDiscriminantAnalysis
4
5 X = titanic_data[['age', 'fare']]
6
7
8 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, r
9
10 # Fit Naive Bayes model
11 nb_model = GaussianNB()
12 nb_model.fit(X_train, y_train)
13
14 # Fit LDA model
15 lda_model = LinearDiscriminantAnalysis()
16 lda_model.fit(X_train, y_train)
17
18 # Fit QDA model
19 qda_model = QuadraticDiscriminantAnalysis()
20 qda_model.fit(X_train, y_train)
21
22 # Compare accuracies
23 nb_accuracy = nb_model.score(X_test, y_test)
24 lda_accuracy = lda_model.score(X_test, y_test)
25 qda_accuracy = qda_model.score(X_test, y_test)
26
27 print("Naive Bayes accuracy: ", nb_accuracy)
28 print("LDA accuracy: ", lda_accuracy)
29 print("QDA accuracy: ", qda_accuracy)
```

Naive Bayes accuracy: 0.6293706293706294

LDA accuracy: 0.6223776223776224

QDA accuracy: 0.6153846153846154

```
In [6]: 1 from sklearn.metrics import classification_report
2
3 nb_predictions = nb_model.predict(X_test)
4 report = classification_report(y_test, nb_predictions,
5                               target_names=['Not Survived', 'Survived'])
6 print("Classification Report for Naive Bayes:\n", report)
7
8 lda_predictions = lda_model.predict(X_test)
9 report = classification_report(y_test, lda_predictions,
10                              target_names=['Not Survived', 'Survived'])
11 print("Classification Report for LDA:\n", report)
12
13 qda_predictions = qda_model.predict(X_test)
14 report = classification_report(y_test, qda_predictions,
15                              target_names=['Not Survived', 'Survived'])
16 print("Classification Report for QDA:\n", report)
```

Classification Report for Naive Bayes:

	precision	recall	f1-score	support
Not Survived	0.62	0.89	0.73	80
Survived	0.68	0.30	0.42	63
accuracy			0.63	143
macro avg	0.65	0.59	0.57	143
weighted avg	0.64	0.63	0.59	143

Classification Report for LDA:

	precision	recall	f1-score	support
Not Survived	0.61	0.89	0.72	80
Survived	0.67	0.29	0.40	63
accuracy			0.62	143
macro avg	0.64	0.59	0.56	143
weighted avg	0.64	0.62	0.58	143

Classification Report for QDA:

	precision	recall	f1-score	support
Not Survived	0.61	0.89	0.72	80
Survived	0.65	0.27	0.38	63
accuracy			0.62	143
macro avg	0.63	0.58	0.55	143
weighted avg	0.63	0.62	0.57	143

All three models performs very similarly. Naive Bayes performs slightly the best with a accuracy of 0.629 while QDA and LDA performs only very slightly worse with f1 scores of 0.615 and 0.622 respectively. Naive Bayes does have the highest accuracy. Overall, it is hard to say which model is exactly better, but Naive Bayes probably has a slight edge.

3. Make a data visualization to show the decision boundary for three models. (20)

```

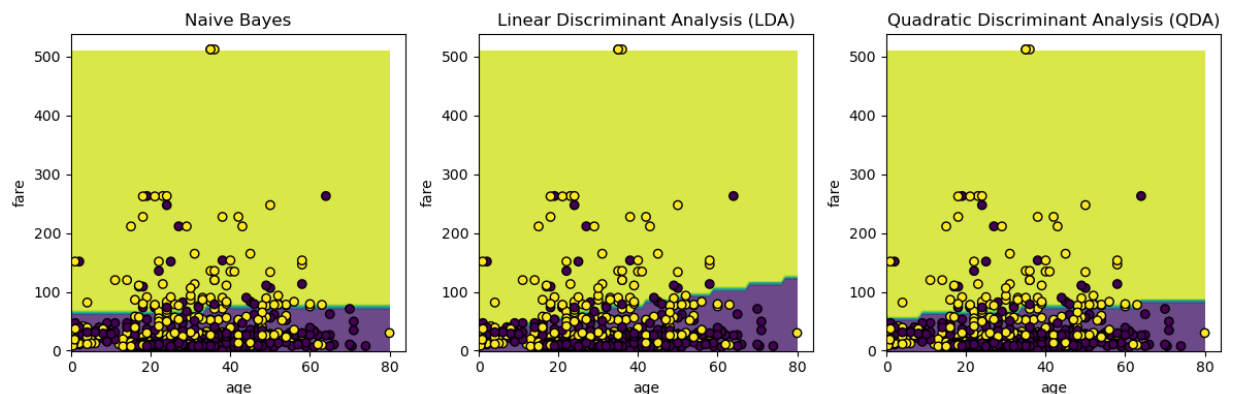
In [7]: 1 import warnings
2 import numpy as np
3 warnings.filterwarnings('ignore')
4
5 # Create a meshgrid of points to plot the decision boundary
6 x_min, x_max = X['age'].min() - 0.5, X['age'].max() + 0.5
7 y_min, y_max = X['fare'].min() - 0.5, X['fare'].max() + 0.5
8 xx, yy = np.meshgrid(np.arange(x_min, x_max, 1), np.arange(y_min, y_max, 1))
9
10 # Make predictions on the meshgrid points for each model
11 nb_pred = nb_model.predict(np.c_[xx.ravel(), yy.ravel()])
12 lda_pred = lda_model.predict(np.c_[xx.ravel(), yy.ravel()])
13 qda_pred = qda_model.predict(np.c_[xx.ravel(), yy.ravel()])
14
15 nb_pred = nb_pred.reshape(xx.shape)
16 lda_pred = lda_pred.reshape(xx.shape)
17 qda_pred = qda_pred.reshape(xx.shape)

```

```

In [8]: 1 # Plot the decision boundaries
2 plt.figure(figsize=(12, 4))
3
4 plt.subplot(1, 3, 1)
5 plt.contourf(xx, yy, nb_pred, alpha=0.8)
6 plt.scatter(X['age'], X['fare'], c=y, edgecolors='k')
7 plt.xlabel('age')
8 plt.ylabel('fare')
9 plt.title('Naive Bayes')
10
11 plt.subplot(1, 3, 2)
12 plt.contourf(xx, yy, lda_pred, alpha=0.8)
13 plt.scatter(X['age'], X['fare'], c=y, edgecolors='k')
14 plt.xlabel('age')
15 plt.ylabel('fare')
16 plt.title('Linear Discriminant Analysis (LDA)')
17
18 plt.subplot(1, 3, 3)
19 plt.contourf(xx, yy, qda_pred, alpha=0.8)
20 plt.scatter(X['age'], X['fare'], c=y, edgecolors='k')
21 plt.xlabel('age')
22 plt.ylabel('fare')
23 plt.title('Quadratic Discriminant Analysis (QDA)')
24
25 plt.tight_layout()
26 plt.show()

```



4. Now fit the models again with all variables. Make sure you have convert the categorical variables to factors. Report the classification table for each models. Which one performs the best? (20)

```
In [9]: 1 titanic_data = pd.get_dummies(titanic_data, columns=['pclass', 'sex', 'embark_town',
2                                     'class', 'who', 'adult_male',
3                                     'deck', 'embark_town'])
4
5 X = titanic_data.drop('survived', axis=1)
6 y = titanic_data['survived']
7
8 X_train, X_test, y_train, y_test = train_test_split(X, y,
9                                                     test_size=0.2, random_state=83)
```

```
In [10]: 1 # Fit Naive Bayes model
2 nb_model = GaussianNB()
3 nb_model.fit(X_train, y_train)
4
5 # Fit LDA model
6 lda_model = LinearDiscriminantAnalysis()
7 lda_model.fit(X_train, y_train)
8
9 # Fit QDA model
10 qda_model = QuadraticDiscriminantAnalysis()
11 qda_model.fit(X_train, y_train)
12
13 # Compare accuracies
14 nb_accuracy = nb_model.score(X_test, y_test)
15 lda_accuracy = lda_model.score(X_test, y_test)
16 qda_accuracy = qda_model.score(X_test, y_test)
17
18 print("Naive Bayes accuracy: ", nb_accuracy)
19 print("LDA accuracy: ", lda_accuracy)
20 print("QDA accuracy: ", qda_accuracy)
```

```
Naive Bayes accuracy: 0.7202797202797203
LDA accuracy: 0.7552447552447552
QDA accuracy: 0.6923076923076923
```

```
In [11]: 1 nb_predictions = nb_model.predict(X_test)
2 report = classification_report(y_test, nb_predictions,
3                               target_names=['Not Survived', 'Survived'])
4 print("Classification Report for Naive Bayes:\n", report)
5
6 lda_predictions = lda_model.predict(X_test)
7 report = classification_report(y_test, lda_predictions,
8                               target_names=['Not Survived', 'Survived'])
9 print("Classification Report for LDA:\n", report)
10
11 qda_predictions = qda_model.predict(X_test)
12 report = classification_report(y_test, qda_predictions,
13                               target_names=['Not Survived', 'Survived'])
14 print("Classification Report for QDA:\n", report)
```

Classification Report for Naive Bayes:

	precision	recall	f1-score	support
Not Survived	0.78	0.70	0.74	80
Survived	0.66	0.75	0.70	63
accuracy			0.72	143
macro avg	0.72	0.72	0.72	143
weighted avg	0.73	0.72	0.72	143

Classification Report for LDA:

	precision	recall	f1-score	support
Not Survived	0.75	0.84	0.79	80
Survived	0.76	0.65	0.70	63
accuracy			0.76	143
macro avg	0.76	0.74	0.75	143
weighted avg	0.76	0.76	0.75	143

Classification Report for QDA:

	precision	recall	f1-score	support
Not Survived	0.66	0.95	0.78	80
Survived	0.85	0.37	0.51	63
accuracy			0.69	143
macro avg	0.75	0.66	0.64	143
weighted avg	0.74	0.69	0.66	143

Naive and QDA performs the worst in this case in terms of accuracy. LDA performs much better than both of the other models in terms of accuracy. So, I would choose LDA over the other two models in this case because it is a more accurate.

## Question 2: Simulation study (10)

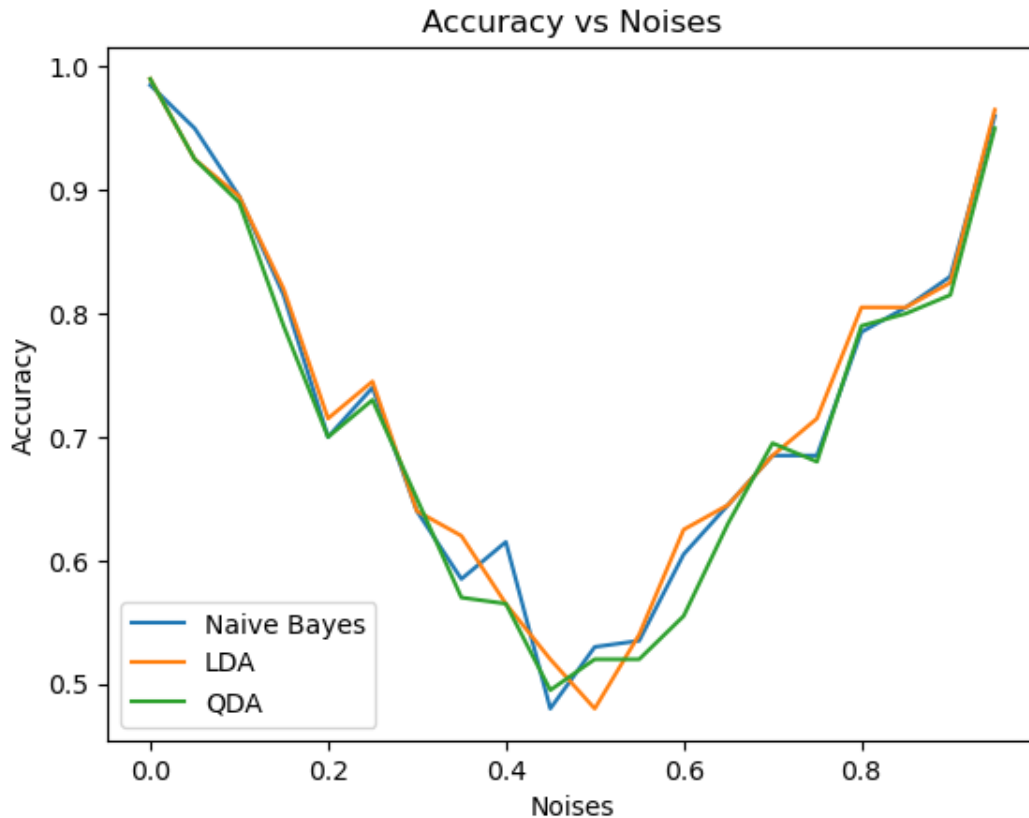
In the following simulation study, please write a sentence to discuss what this simulation code is doing and what you have seen in the figure.

```
In [12]: 1 import numpy as np
          2
          3 # Set the random seed for reproducibility
          4 np.random.seed(4400)
          5
          6 # Define the range of dataset sizes
          7 noises = np.arange(0,1,0.05)
          8
          9 # Initialize lists to store accuracy and time results
         10 accuracy_nb = []
         11 accuracy_lda = []
         12 accuracy_qda = []
```



```
In [13]: 1 import pandas as pd
2 from sklearn.naive_bayes import GaussianNB
3 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
4 from sklearn.discriminant_analysis import QuadraticDiscriminantAnalysis
5 from sklearn.model_selection import train_test_split
6
7 for noise in noises:
8
9     # Generate a random classification dataset
10    X1 = np.random.normal(5, 1, 1000)
11    X2 = np.random.normal(0, 1, 1000)
12    X3 = np.random.normal(2, 1, 1000)
13    X4 = np.random.normal(-3, 2, 1000)
14    X = pd.DataFrame({'X1': X1, 'X2': X2, 'X3': X3, 'X4': X4})
15    y = np.where(X1 > 5, 'group1', 'group2')
16
17    indices_to_change = np.random.choice(1000,
18                                         size=int(noise * 1000),
19                                         replace=False)
20
21    for index in indices_to_change:
22        if y[index] == 'group1':
23            y[index] = 'group2'
24        else:
25            y[index] = 'group1'
26
27    # Split the dataset into training and test sets
28    X_train, X_test, y_train, y_test = train_test_split(X, y,
29                                                         test_size=0.2,
30                                                         random_state=4400)
31
32    # Fit Naive Bayes model and calculate accuracy and time
33    nb_model = GaussianNB()
34    nb_model.fit(X_train, y_train)
35    accuracy_nb.append(nb_model.score(X_test, y_test))
36
37    # Fit LDA model and calculate accuracy and time
38    lda_model = LinearDiscriminantAnalysis()
39    lda_model.fit(X_train, y_train)
40    accuracy_lda.append(lda_model.score(X_test, y_test))
41
42    # Fit QDA model and calculate accuracy and time
43    qda_model = QuadraticDiscriminantAnalysis()
44    qda_model.fit(X_train, y_train)
45    accuracy_qda.append(qda_model.score(X_test, y_test))
```

```
In [14]: 1 import matplotlib.pyplot as plt
2
3 plt.plot(noises, accuracy_nb, label='Naive Bayes')
4 plt.plot(noises, accuracy_lda, label='LDA')
5 plt.plot(noises, accuracy_qda, label='QDA')
6 plt.xlabel('Noises')
7 plt.ylabel('Accuracy')
8 plt.title('Accuracy vs Noises')
9 plt.legend()
10 plt.show()
```



This simulation study appears to be conducting a study to assess the performance of three different classification models (Naive Bayes, Linear Discriminant Analysis, and Quadratic Discriminant Analysis) on a synthetic dataset with varying levels of noise. The noise is introduced by flipping the labels of a random subset of instances, causing misclassification in the dataset.

From the graph, we see that as the level of noise in the dataset increases, the classification accuracy of all three models (Naive Bayes, LDA, and QDA) decreases. This demonstrates that the introduced noise negatively impacts the models' performance, making it harder for them to correctly classify the data.

The simulation is evaluating the robustness of these classification models to noisy data, and the figure illustrates the impact of noise on their accuracy. From the figure, it seems like all three models are acting similarly. It does seem like LDA predicts a little better than the other two models as noise increases.

### Question 3: Compare the models (20)

Please summarize the similarity and difference between Naive Bayes, LDA and QDA. Write at least three similarities between the models and at least two difference for each model. Hint: think about the how the models are proposed, the assumptions and the decision boundary etc.

## Similarities:

- i) Naive Bayes, LDA, and QDA are all classification algorithms used to make predictions based on input features. All three are considered supervised learning techniques.
- ii) Probabilistic models: making use of probability distributions to estimate class membership probabilities.
- iii) Classification based on posterior probabilities. After estimating the probability distributions, these algorithms calculate the posterior probability of each class given the input features using Bayes' theorem. The class with the highest posterior probability is assigned as the predicted class label.

## Differences:

### Naive Bayes

1. Independence assumption: assumes that features are conditionally independent given the class label even though that may not be true in the real world, but it does simplify the model. The model is very simple, fast, and easy to implement.
2. Simple Model: Computationally efficient and requires less training data. No assumption for the decision boundary (but usually a simple one close to linear)

### Linear Discriminant Analysis

1. Assumption of Equal Covariance Matrices: LDA assumes that all classes share the same covariance matrix. This simplifies the model but may not hold if the classes have different covariance structures.
2. Linear Decision Boundary: LDA aims to find a linear decision boundary that maximizes the separability between classes. It's well-suited for problems with linearly separable classes.

### Quadratic Discriminant Analysis

1. Relaxation of Covariance Assumption and Higher Model Complexity: QDA relaxes the equal covariance matrix assumption by allowing each class to have its own covariance matrix. This makes it more flexible but requires more parameters. Potentially higher model complexity (overfitting).
2. Non-Linear Decision Boundary (Quadratic): QDA can capture non-linear decision boundaries which makes it more powerful in cases where classes are not linearly separable. It is a more flexible model however may require more data to estimate the covariance matrices accurately. Maximizes separation between classes and minimizes within-class variance.