| Topic | Naive bayes |
|---|---|
| Class Description | Students learn the concepts of Naive bayes. They also learn about Bayes theorem. Students compare the Naive Bayes algorithm with Logistics regression and make conclusions. |
| Class | C120 |
| Class time | 45 mins |
| Goal | ● Explore the concept of Naive bayes algorithm.<br>● Create a prediction model using Naive bayes algorithm. |
| Resources Required | ● Teacher Resources<br> ○ Google Colab Notebook<br> ○ Laptop with internet connectivity<br> ○ Earphones with mic<br> ○ Notebook and pen<br><br>● Student Resources<br> ○ Google Colab Notebook<br> ○ Laptop with internet connectivity<br> ○ Earphones with mic<br> ○ Notebook and pen |

| Class structure | Warm Up<br>Teacher-led Activity<br>Student-led Activity<br>Wrap up | 5 mins<br>15 min<br>15 min<br>5 min |
|---|---|---|

<div align="center">

**CONTEXT**

● **Explore the concept of Naive Bayes algorithm and learn about bayes law.**

</div>

| Class Steps | Teacher Action | Student Action |
|---|---|---|
| Step 1:<br>Warm Up<br>(5 mins) | Hi <Student Name>!<br>How are you doing today? | **ESR:**<br>- We learned about the concept of Decision Tree.<br>- We wrote a supervised |

| | | |
|---|---|---|
| | Let's quickly revise what we did in last class? | learning algorithm called Decision Tree.<br>- We also drew a prediction flow chart of the data. |
| | Awesome. So in last class we saw a supervised learning algorithm which made predictions from the decision rules it learnt from prior data training. Today we are going to look at another algorithm which assumes that every feature in a dataset is independent and has its own contribution to the outcome.<br>Do you remember any other algorithm which we have studied that predicts the dependency of the variables in the dataset? | **ESR:**<br>Logistics regression. |
| | Yes, so as you see Naive bayes and Logistic regression seem the same so many times people get confused on which algorithm to use.<br>Before that let's learn more about the Naive bayes algorithm. | |
| <td colspan="3" style="text-align:center">**Teacher Initiates Screen Share**</td> |
| <td colspan="3">**CHALLENGE**<br>● **Create Naive bayes and logistics regression prediction models.**<br>● **See how each model performs on the basis of variable dependencies.**<br>● **Make a conclusion on the basis of the outcome.**</td> |
| **Step 2:**<br>**Teacher-led**<br>**Activity**<br>**(15 min)** | Naive Bayes algorithm is a supervised machine learning algorithm based on the Bayes Probability theorem. Naive Bayes | *Student listens and asks questions.* |

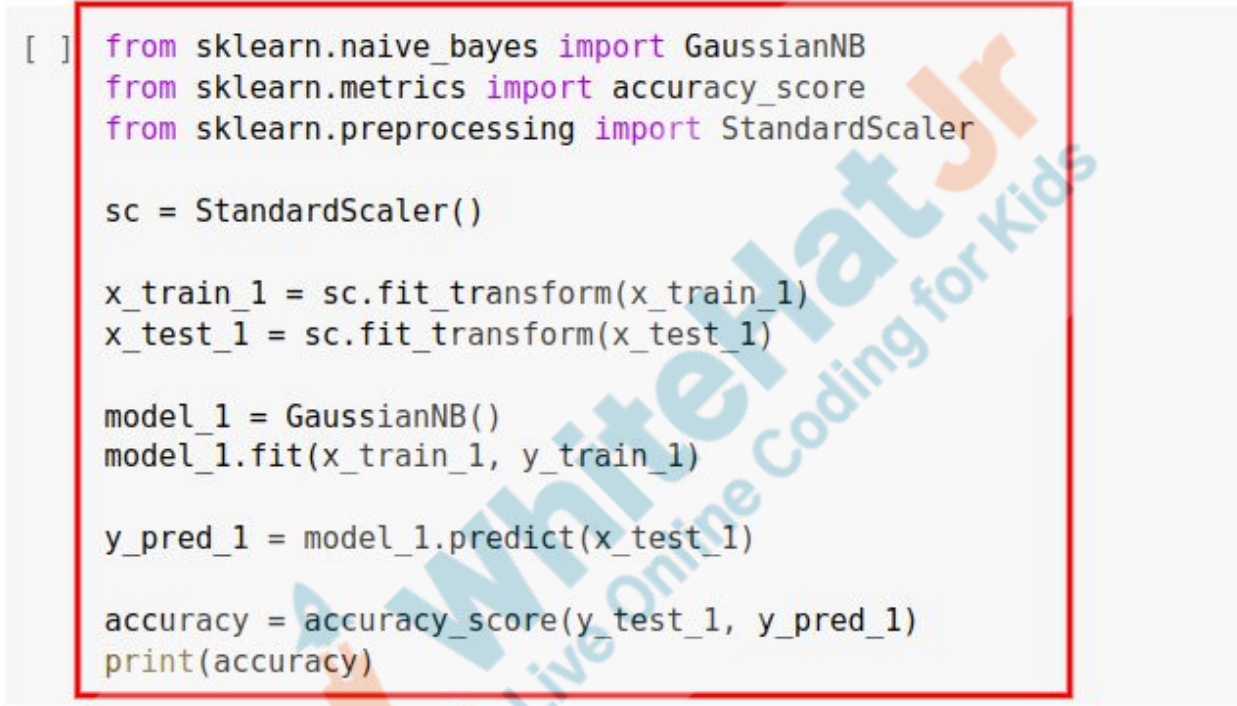| | | |
|---|---|---|
| | assumes that there is no correlation between the features in a dataset used to train the model.<br><br>Despite the oversimplified assumptions, Naive Bayes works very well in many real world complex problems. They require a relatively small number of training data samples to perform classification efficiently, compared to other algorithms like Logistic Regression and Decision trees, that we studied earlier. | |
| | Naive Bayes algorithm works on Bayes theorem.<br><br>Bayes theorem or Bayes' law or Bayes' rule describes the probability of an event, based on prior knowledge of conditions that might be related to the event.<br>What this means is that Bayes theorem describes the probability of a feature, based on prior knowledge of situations related to that feature.<br><br>For example, if the probability of someone having diabetes is related to his or her age, then by using the Bayes Theorem, the age can be used to more accurately predict the probability of having diabetes. | *Student listens and asks questions.* |
| | The word naive implies that every pair of features in the dataset is independent of each other. Naive Bayes works on the assumption that | |

| | | |
|---|---|---|
| | the value of a particular feature is independent of any other feature.<br><br>For example, how can you classify if a vegetable is a tomato?<br><br>Yes, but with Naive Bayes, each of these three features (shape, size and color) contributes independently to the probability that the vegetable is a tomato. Also, it assumes that there is no possible correlation between the shape, size and color. | **ESR:**<br>A vegetable may be classified as a tomato if it's round, about 4-5 cm in diameter, and red in color. |
| | Let's write some code to understand the differences between the two as we try to understand Naive Bayes a little more.<br><br>*<Teacher downloads the data from Teacher activity 1 and Opens the Google Colab Notebooks from Teacher activity 2>* | - |
| | *<Teacher uploads the data and reads it using pandas and prints it>*<br>Here we are using the data for the causes of diabetes.<br><br>Code:-<br>**#Uploading the csv**<br>**from google.colab import files**<br>**data_to_load = files.upload()**<br><br>**#Code to read the file.**<br>**import pandas as pd**<br><br>**df = pd.read_csv('diabetes.csv')** | *The student helps the teacher with the code to upload the data file, read it and print it's content.* |

| | **print(df.head())** | |
|---|---|---|



```
#Uploading the csv
from google.colab import files
data_to_load = files.upload()
```

Choose Files  No file chosen          Upload widget is only available whe
Saving diabetes.csv to diabetes.csv

```
import pandas as pd

df = pd.read_csv('diabetes.csv')

print(df.head())
```

```
   glucose  bloodpressure  diabetes
0       40             85         0
1       40             92         0
2       45             63         1
3       45             80         0
4       40             73         1
```

| | In the data that we have, we can see that we have glucose, bloodpressure and we know if the given person has diabetes or not.<br><br>Here, we will use the glucose and the bloodpressure to predict if the person has diabetes or not using Naive Bayes. | - |
|---|---|---|

| | Before that what is the first step we do with the data?<br><br>Perfect!.<br>*<Teacher codes to split the data for training and testing the model>*<br><br>Code:-<br>**from sklearn.model_selection import train_test_split**<br><br>**X = df[["glucose", "bloodpressure"]]**<br>**y = df["diabetes"]**<br><br>**x_train_1, x_test_1, y_train_1, y_test_1 = train_test_split(X, y, test_size=0.25, random_state=42)** | **ESR:**<br>We split the data into 2 parts to train and test the model.<br><br><br><br>*The student helps the teacher to code for splitting the data to train and test the model.* |

```
[ ]  from sklearn.model_selection import train_test_split

     X = df[["glucose", "bloodpressure"]]
     y = df["diabetes"]

     x_train_1, x_test_1, y_train_1, y_test_1 = train_test_split(X, y, test_size=0.25, random_state=42)
```

| | Now we'll code to train the model with Naive Bayes.<br><br>Code:-<br># first we are going to import GaussianNB module from sklearn naive_bayes module<br>- Gaussian Naive Bayes algorithm is a special type of NB algorithm. It's specifically used when the features have continuous values. It's also assumed that all the features are following a gaussian distribution i.e, | |

| | | |
|---|---|---|
| | normal distribution.<br><br>**from sklearn.naive_bayes import GaussianNB**<br><br>**from sklearn.metrics import accuracy_score**<br>**from sklearn.preprocessing import StandardScaler**<br><br>**sc = StandardScaler()**<br><br>**x_train_1 = sc.fit_transform(x_train_1)**<br>**x_test_1 = sc.fit_transform(x_test_1)**<br><br>**model_1 = GaussianNB()**<br>**model_1.fit(x_train_1, y_train_1)**<br><br>**y_pred_1 = model_1.predict(x_test_1)**<br><br>**accuracy = accuracy_score(y_test_1, y_pred_1)**<br>**print(accuracy)**<br><br>Can you tell me what accuracy_score and StandardScaler do? | **ESR:**<br>accuracy_score returns "accuracy classification score". What it does is the calculation of "How accurate the classification is".<br><br>StandardScaler standardizes a feature by subtracting the mean and then scaling to unit variance.<br>Unit variance means dividing all the values by the standard deviation. |

| | | |
|---|---|---|
| | Perfect. So what accuracy do we see here?<br><br>Now let's see if we can get this accuracy using the logistics regression. | **ESR:**<br>We can see an amazing accuracy of 94.4%. |

```
[ ]  from sklearn.naive_bayes import GaussianNB
     from sklearn.metrics import accuracy_score
     from sklearn.preprocessing import StandardScaler

     sc = StandardScaler()

     x_train_1 = sc.fit_transform(x_train_1)
     x_test_1 = sc.fit_transform(x_test_1)

     model_1 = GaussianNB()
     model_1.fit(x_train_1, y_train_1)

     y_pred_1 = model_1.predict(x_test_1)

     accuracy = accuracy_score(y_test_1, y_pred_1)
     print(accuracy)

  👤  0.9437751004016064
```

| | | |
|---|---|---|
| | So let's split the data to train and test our logistics regression model.<br><br>*<Teacher codes to split the data to train and test the model>*<br><br>Code:-<br><br>**from sklearn.model_selection import train_test_split** | *The student helps the teacher with the code.* |

| | | |
|---|---|---|
| | **X = df[["glucose", "bloodpressure"]]**<br>**y = df["diabetes"]**<br><br>**x_train_2, x_test_2, y_train_2, y_test_2 = train_test_split(X, y, test_size=0.25, random_state=42)** | |

```
[ ]  from sklearn.model_selection import train_test_split

     X = df[["glucose", "bloodpressure"]]
     y = df["diabetes"]

     x_train_2, x_test_2, y_train_2, y_test_2 = train_test_split(X, y, test_size=0.25, random_state=42)
```

| | | |
|---|---|---|
| | Now we have data ready, let's train our model on this data.<br><br>*<Teacher codes to train the logistics regression model>*<br><br>Code:-<br>**from sklearn.linear_model import LogisticRegression**<br>**from sklearn.metrics import accuracy_score**<br>**from sklearn.preprocessing import StandardScaler**<br><br>**sc = StandardScaler()**<br><br>**x_train_2 = sc.fit_transform(x_train_2)**<br>**x_test_2 = sc.fit_transform(x_test_2)**<br><br>**model_2 =** | *The student helps the teacher with the code.* |

| | | |
|---|---|---|
| | **LogisticRegression(random_state = 0)**<br>**model_2.fit(x_train_2, y_train_2)**<br><br>**y_pred_2 = model_2.predict(x_test_2)**<br><br>**accuracy = accuracy_score(y_test_2, y_pred_2)**<br>**print(accuracy)**<br><br>What accuracy do you see? | **ESR:**<br>I can see an accuracy of 91.6%. |

```
[ ]  from sklearn.linear_model import LogisticRegression
     from sklearn.metrics import accuracy_score
     from sklearn.preprocessing import StandardScaler

     sc = StandardScaler()

     x_train_2 = sc.fit_transform(x_train_2)
     x_test_2 = sc.fit_transform(x_test_2)

     model_2 = LogisticRegression(random_state = 0)
     model_2.fit(x_train_2, y_train_2)

     y_pred_2 = model_2.predict(x_test_2)

     accuracy = accuracy_score(y_test_2, y_pred_2)
     print(accuracy)

  ●  0.9156626506024096
```

| | | |
|---|---|---|
| | While the accuracy score for both the datasets was close, with Naive Bayes giving us an accuracy of 94.4% and logistic regression giving us an | **ESR:**<br>Varied |

| | | |
|---|---|---|
| | accuracy of 91.6%, Naive Bayes still performed better.<br><br>Can you guess why?<br><br>The reason for this is that if we look at our features again, we can see that the Glucose and the Blood Pressure had no correlation with each other. They both contribute individually to whether a person would have diabetes or not. This is exactly what Naive Bayes algorithm assumes, that all the features contribute individually to the outcome. | |
| | This was for the case of where Naive Bayes outperforms Logistic Regression, but let's see an example of the case where Logistic Regression outperforms Naive Bayes.<br>Can you try doing that? I'll guide you wherever you need help. | **ESR:**<br>Yes! |

| | | |
|---|---|---|
| | Now it's your turn. Please share your screen with me. | |

- **Ask Student to press ESC key to come back to panel**
- **Guide Student to start Screen Share**
- **Teacher gets into Fullscreen**

## ACTIVITY

- **Create Naive Bayes and logistics regression prediction models for different data.**
- **Make a conclusion on the basis of the outcome.**

| Step 3:<br>Student-Led<br>Activity<br>(15 min) | *Teacher helps the student to download the data and open the Colab notebook.* | *Student downloads the data from Student Activity 1 and Opens Colab Notebook from Student Activity 2.* |
|---|---|---|
| | Here we are using the income data of various people.<br>You can use pd.describe() to view some basic statistics details such as age, workclass etc.<br>*<Teacher helps the student with the code>*<br><br>Code:-<br>**#Uploading the csv**<br>**from google.colab import files**<br>**data_to_load = files.upload()**<br><br>**import pandas as pd**<br><br>**df = pd.read_csv('income.csv')**<br><br>**print(df.head())**<br>**print(df.describe())** | *Student uploads the data and prints the data.* |

```
[ ]  #Uploading the csv
     from google.colab import files
     data_to_load = files.upload()
```

Choose Files  No file chosen          Upload widget is only available when the cell has been executed in the current browser
Saving income.csv to income.csv

```
[ ]  import pandas as pd

     df = pd.read_csv('income.csv')

     print(df.head())
     print(df.describe())
```

```
    age          workclass  ...  native-country  income
0   39          State-gov  ...   United-States   <=50K
1   50   Self-emp-not-inc  ...   United-States   <=50K
2   38            Private  ...   United-States   <=50K
3   53            Private  ...   United-States   <=50K
4   28            Private  ...            Cuba   <=50K

[5 rows x 14 columns]
                age   education-num   capital-gain   capital-loss   hours-per-week
count   45222.000000    45222.000000   45222.000000   45222.000000     45222.000000
mean       38.547941       10.118460    1101.430344      88.595418        40.938017
std        13.217870        2.552881    7506.430084     404.956092        12.007508
min        17.000000        1.000000       0.000000       0.000000         1.000000
25%        28.000000        9.000000       0.000000       0.000000        40.000000
50%        37.000000       10.000000       0.000000       0.000000        40.000000
75%        47.000000       13.000000       0.000000       0.000000        45.000000
max        90.000000       16.000000   99999.000000    4356.000000        99.000000
```

| | | |
|---|---|---|
| | From the given data, we will consider the following fields to determine the salary of a person:<br><br>Age<br>Hours Per Week<br>Education Number<br>Capital Gain<br>Capital Loss<br><br>Now let's split the data to train and test the model.<br>*<Teacher helps student to write code to split the model>* | *Student codes to split the data to train and test the model.* |

| | Code:-<br><br>**from sklearn.model_selection import train_test_split**<br><br>**X = df[["age", "hours-per-week", "education-num", "capital-gain", "capital-loss"]]**<br>**y = df["income"]**<br><br>**x_train_1, x_test_1, y_train_1, y_test_1 = train_test_split(X, y, test_size=0.25, random_state=42)** | |

```
from sklearn.model_selection import train_test_split

X = df[["age", "hours-per-week", "education-num", "capital-gain", "capital-loss"]]
y = df["income"]

x_train_1, x_test_1, y_train_1, y_test_1 = train_test_split(X, y, test_size=0.25, random_state=42)
```

| | Now let's train the Naive Bayes model.<br><br>*<Teacher helps student to code for training the Naive Bayes model>*<br><br>Code:-<br>**from sklearn.naive_bayes import GaussianNB**<br>**from sklearn.metrics import accuracy_score**<br>**from sklearn.preprocessing import StandardScaler**<br><br>**sc = StandardScaler()**<br><br>**x_train_1 =** | *The student codes to train the Naive Bayes model.* |

| | | |
|---|---|---|
| | sc.fit_transform(x_train_1)<br>x_test_1 =<br>sc.fit_transform(x_test_1)<br><br>model_1 = GaussianNB()<br>model_1.fit(x_train_1, y_train_1)<br><br>y_pred_1 =<br>model_1.predict(x_test_1)<br><br>accuracy =<br>accuracy_score(y_test_1,<br>y_pred_1)<br>print(accuracy)<br><br><br>What accuracy can you see? | **ESR:**<br>We can see the accuracy of almost 79%. |

```
[ ] from sklearn.naive_bayes import GaussianNB
    from sklearn.metrics import accuracy_score
    from sklearn.preprocessing import StandardScaler

    sc = StandardScaler()

    x_train_1 = sc.fit_transform(x_train_1)
    x_test_1 = sc.fit_transform(x_test_1)

    model_1 = GaussianNB()
    model_1.fit(x_train_1, y_train_1)

    y_pred_1 = model_1.predict(x_test_1)

    accuracy = accuracy_score(y_test_1, y_pred_1)
    print(accuracy)

    0.7896692021935255
```

| | Alright now let's check the accuracy with logistics regression.<br><br>So let's split the data to train and test the logistics regression.<br><br>*<Teacher helps the student to split the data>*<br>Code:<br>**from sklearn.model_selection import train_test_split**<br><br>**X = df[["age", "hours-per-week", "education-num", "capital-gain", "capital-loss"]]**<br>**y = df["income"]**<br><br>**x_train_2, x_test_2, y_train_2, y_test_2 = train_test_split(X, y, test_size=0.25, random_state=42)** | *Student codes to split the data to train and test the model.* |
|---|---|---|

```
[ ]  from sklearn.model_selection import train_test_split

     X = df[["age", "hours-per-week", "education-num", "capital-gain", "capital-loss"]]
     y = df["income"]

     x_train_2, x_test_2, y_train_2, y_test_2 = train_test_split(X, y, test_size=0.25, random_state=42)
```

| | Now let's train the logistics regression model.<br><br>*<Teacher helps student to code to train the model and print the accuracy>*<br><br>Code:<br>**from sklearn.linear_model import** | *<Student codes to train the model and print the accuracy>* |
|---|---|---|

| | | |
|---|---|---|
| | ```
LogisticRegression
from sklearn.metrics import
accuracy_score
from sklearn.preprocessing import
StandardScaler

sc = StandardScaler()

x_train_2 =
sc.fit_transform(x_train_2)
x_test_2 =
sc.fit_transform(x_test_2)

model_2 =
LogisticRegression(random_state
= 0)
model_2.fit(x_train_2, y_train_2)

y_pred_2 =
model_2.predict(x_test_2)

accuracy =
accuracy_score(y_test_2,
y_pred_2)
print(accuracy)
``` | |

```
[ ] from sklearn.linear_model import LogisticRegression
    from sklearn.metrics import accuracy_score
    from sklearn.preprocessing import StandardScaler

    sc = StandardScaler()

    x_train_2 = sc.fit_transform(x_train_2)
    x_test_2 = sc.fit_transform(x_test_2)

    model_2 = LogisticRegression(random_state = 0)
    model_2.fit(x_train_2, y_train_2)

    y_pred_2 = model_2.predict(x_test_2)

    accuracy = accuracy_score(y_test_2, y_pred_2)
    print(accuracy)

    0.8116929064213692
```

|  | Now what accuracy do we see?<br><br>We can see that the logistics regression outperforms the Naive Bayes module. | **ESR:**<br>We see the accuracy of 81.1%. |
|  | In the first dataset, as we pointed out earlier, both the glucose and the blood pressure had little correlation, and both of them were contributing individually to whether a person has diabetes or not.<br><br>Conclusion: In these kinds of dataset, where all the features contribute |  |

individually to the outcome, Naive Bayes outperforms logistic regression and is highly efficient.

In the second dataset, Logistic Regression outperformed Naive Bayes. The reason is that in this dataset, not all features contribute individually to the outcome. For example, there have been people of all age groups earning both less than and more than 50K. There have also been people with all education numbers that have an income of both less and more than 50K. Here, the combination of all the features is a better predictor of whether a person is earning more than or less than 50K, instead of all features having their individual contribution.

| Teacher Guides Student to Stop Screen Share |
|---|

| **FEEDBACK** |
|---|
| ● **Appreciate the student for their efforts** |
| ● **Identify 2 strengths and 1 area of progress for the student** |

| **Step 4:**<br>**Wrap-Up**<br>**(5 min)** | So today we saw Naive Bayes algorithm for prediction.<br>Can you quickly revise what we did in today's class? | **ESR:**<br>- We learned about the Naive bayes algorithm and theorem.<br>- We did a comparison between Naive bayes and logistics regression and saw how they outperform each |

| | | |
|---|---|---|
| | | other in different circumstances. |
| | The next class is going to be special! We will use our understanding of image processing techniques to create an invisibility cloak. | - |
| **Project Overview** | **Naive Bayes**<br><br>**Goal of the Project:**<br><br>In this project you will apply what you learned in the class and create your own algorithm.<br><br><br>**Story:**<br><br>Suppose you are working as a product manager at a wine factory, where you have to classify the product in various categories, you have large data of wines. Naive Bayes is the most straightforward and fast classification algorithm, which is suitable for a large chunk of data.Apply this algorithm to the data you have and note what interesting insights you find.<br><br><br><br>I am very excited to see your project solution and I know you will do really well. | |

| | Bye Bye! | |
|---|---|---|

| | **Teacher Clicks** ✖ End Class | |
|---|---|---|

| **Additional Activities** | *Encourage the student to write reflection notes in their reflection journal using markdown.*<br><br>Use these as guiding questions:<br><br>● What happened today?<br>  - Describe what happened<br>  - Code I wrote<br>● How did I feel after the class?<br>● What have I learned about programming and developing games?<br>● What aspects of the class helped me? What did I find difficult? | *The student uses the markdown editor to write her/his reflection in a reflection journal.* |
|---|---|---|

| Activity | Activity Name | Links |
|---|---|---|
| Teacher Activity 1 | Data for diabetes | https://raw.githubusercontent.com/whitehatjr/datasets/master/C120/diabetes.csv |
| Teacher Activity 2 | Google Colab Notebook | https://colab.research.google.com/notebooks/intro.ipynb#recent=true |
| Teacher Activity 3 | Reference code | https://drive.google.com/file/d/1ZMw4bHrqhp69q-QtujVZVmHnA8xxlwuz/view?usp=sharing |

| Student Activity 1 | Data for income | https://github.com/whitehatjr/datasets/blob/master/C120/income.csv |
|---|---|---|
| Student Activity 2 | Google Colab Notebook | https://colab.research.google.com/notebooks/intro.ipynb#recent=true |