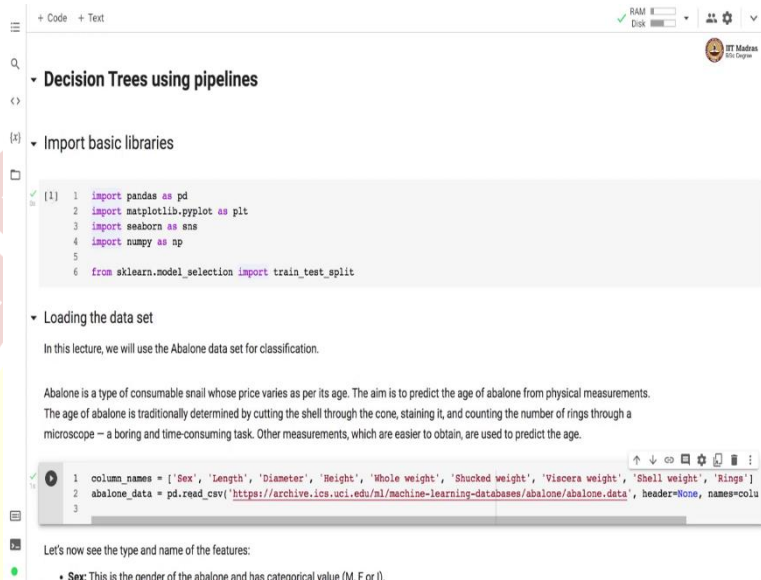


IIT Madras
ONLINE DEGREE

Machine Learning Practice
Professor Dr. Ashish Tendulkar
Indian Institute of Technology, Madras
Decision Trees for Classification – Abalone

(Refer Slide Time: 00:10)



```
[1] 1 import pandas as pd
    2 import matplotlib.pyplot as plt
    3 import seaborn as sns
    4 import numpy as np
    5
    6 from sklearn.model_selection import train_test_split
```

▼ Loading the data set

In this lecture, we will use the Abalone data set for classification.

Abalone is a type of consumable snail whose price varies as per its age. The aim is to predict the age of abalone from physical measurements. The age of abalone is traditionally determined by cutting the shell through the cone, staining it, and counting the number of rings through a microscope – a boring and time-consuming task. Other measurements, which are easier to obtain, are used to predict the age.

```
1 column_names = ['Sex', 'Length', 'Diameter', 'Height', 'Whole weight', 'Shucked weight', 'Viscera weight', 'Shell weight', 'Rings']
2 abalone_data = pd.read_csv('https://archive.ics.uci.edu/ml/machine-learning-databases/abalone/abalone.data', header=None, names=column_names)
3
```

Let's now see the type and name of the features:

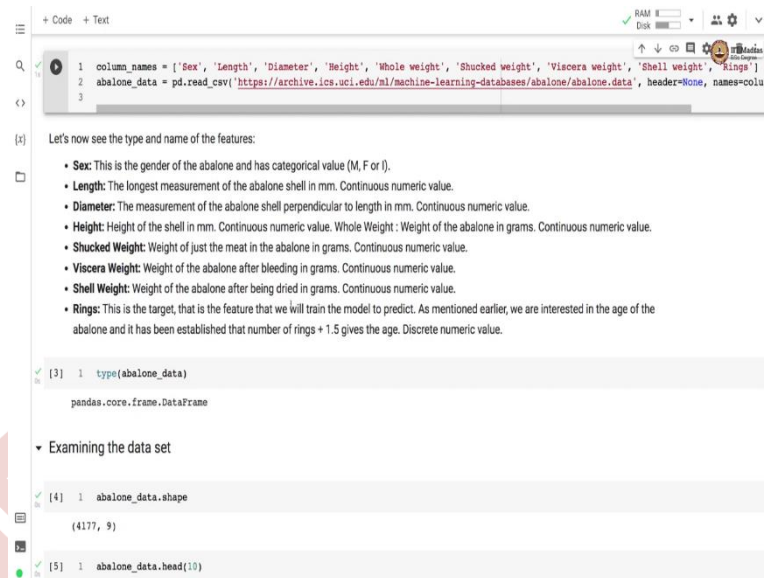
- Sex: This is the gender of the abalone and has categorical value (M, F or I).

Namaste! Welcome to the next video of Machine Learning Practice Course. In this video, we will implement decision trees using pipeline. We import basic libraries like pandas. Then for plotting, we import matplotlib.pyplot and Seaborn. We also import NumPy. And then for model selection, we import train_test_split. In this video, we will be using Abalone data for classification.

So, Abalone is a type of consumable snail, whose price varies as per its age. Our aim is to predict the age of Abalone from the physical measurements. The age of Abalone is traditionally determined by cutting the shell through the cone, staining it and counting the number of rings through a microscope, which is a boring and time consuming task. Other measurements which are easier to obtain, can be used to predict the age of the snail, which is an alternative method.

And that is what we are going to explore in this collab using decision trees. So, in this dataset, we have columns like the sex, the length, diameter, height, whole weight, shucked weight, viscera weight, shell weight, and rings. We load the Abalone data from UCI machine learning repository. And this data is available in CSV format so we use pandas.read_csv function to read the data.

(Refer Slide Time: 02:00)



```
1 column_names = ['Sex', 'Length', 'Diameter', 'Height', 'Whole weight', 'Shucked weight', 'Viscera weight', 'Shell weight', 'Rings']
2 abalone_data = pd.read_csv('https://archive.ics.uci.edu/ml/machine-learning-databases/abalone/abalone.data', header=None, names=column_names)
3
```

Let's now see the type and name of the features:

- **Sex:** This is the gender of the abalone and has categorical value (M, F or I).
- **Length:** The longest measurement of the abalone shell in mm. Continuous numeric value.
- **Diameter:** The measurement of the abalone shell perpendicular to length in mm. Continuous numeric value.
- **Height:** Height of the shell in mm. Continuous numeric value. Whole Weight : Weight of the abalone in grams. Continuous numeric value.
- **Shucked Weight:** Weight of just the meat in the abalone in grams. Continuous numeric value.
- **Viscera Weight:** Weight of the abalone after bleeding in grams. Continuous numeric value.
- **Shell Weight:** Weight of the abalone after being dried in grams. Continuous numeric value.
- **Rings:** This is the target, that is the feature that we will train the model to predict. As mentioned earlier, we are interested in the age of the abalone and it has been established that number of rings + 1.5 gives the age. Discrete numeric value.

```
[3]: type(abalone_data)
pandas.core.frame.DataFrame
```

Examining the data set

```
[4]: abalone_data.shape
(4177, 9)
```

```
[5]: abalone_data.head(10)
```

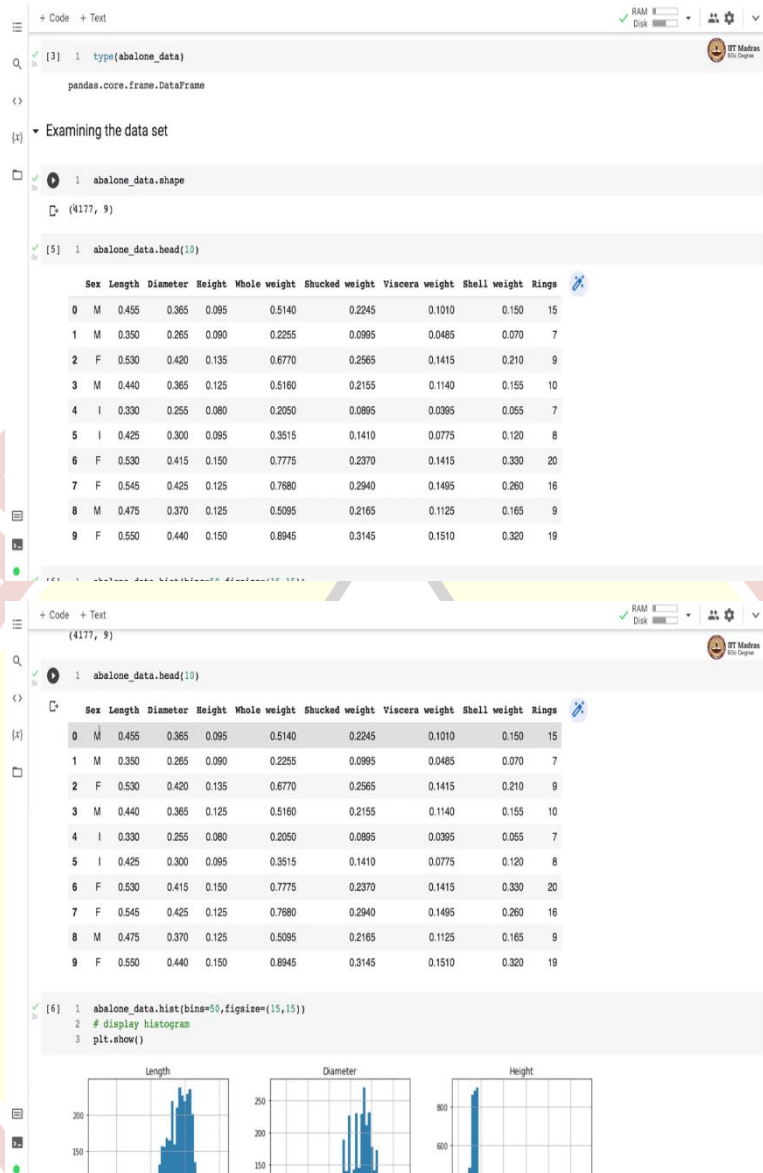
So, let us look at what each type of columns represent. So, we have sex that represent the gender of Abalone and has categorical values which are male, female, or infant. Then there is a length which is the longest measurement of the abalone shell in millimetre, it is continuous numerical value.

Then we have diameter which is the measurement of abalone shell perpendicular to the length in millimeter again a continuous numerical value than your height which is height of the shell in millimeter continuous numerical value. Whole weight is the weight of Abalone in gram which is again a continuous value.

Then there is a shucked weight which is weight of just meat in Abalone in grams. Viscera weight which is weight of abalone after bleeding in grams, then shell weight which is weight of Abalone after being dried in grams. And rings, this is the target, rings is a target variable for us. And we will train our model to predict the rings. And we are interested in age of Abalone.

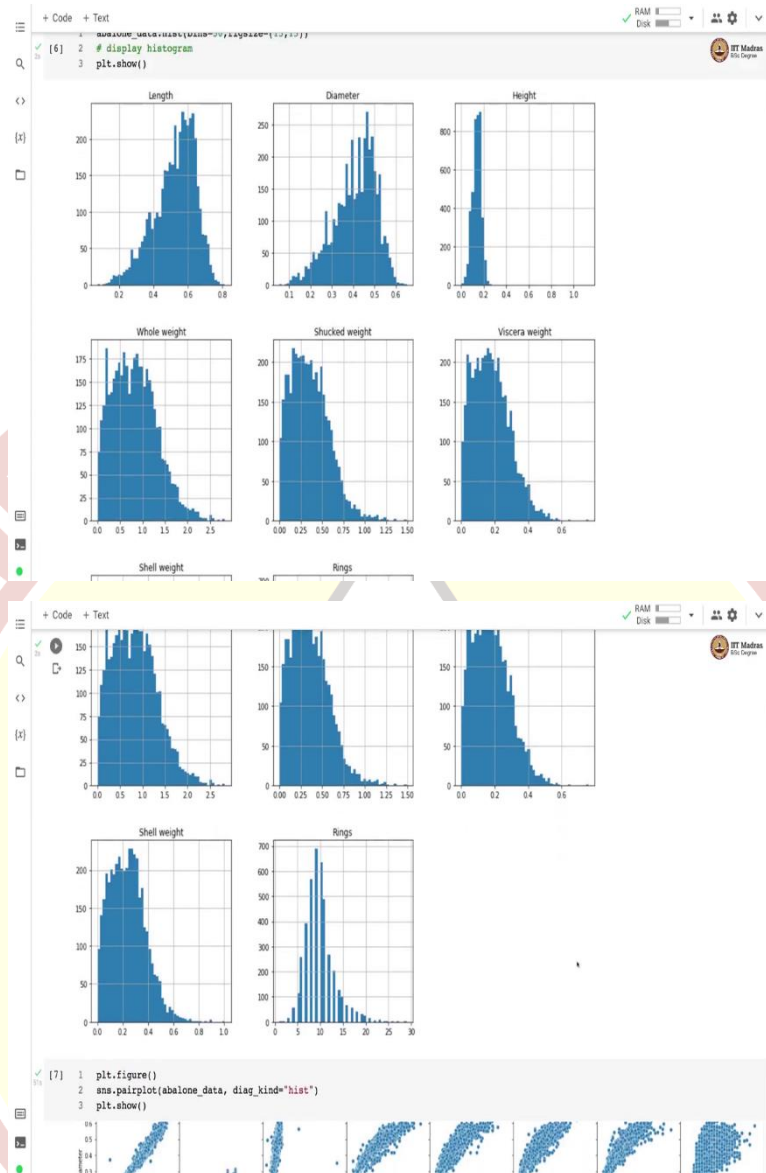
And it is established that the number of rings + 1.5 gives the age of Abalone and it's a discrete numeric value. So, we have bunch of features which are in continuous, which are continuous numerical values. And our target is a discrete numerical value. So, even though the target is discrete numerical value, here, what we are going to do is we are going to treat this problem as a classification problem. And we will try to predict these number of rings for Abalone.

(Refer Slide Time: 03:57)



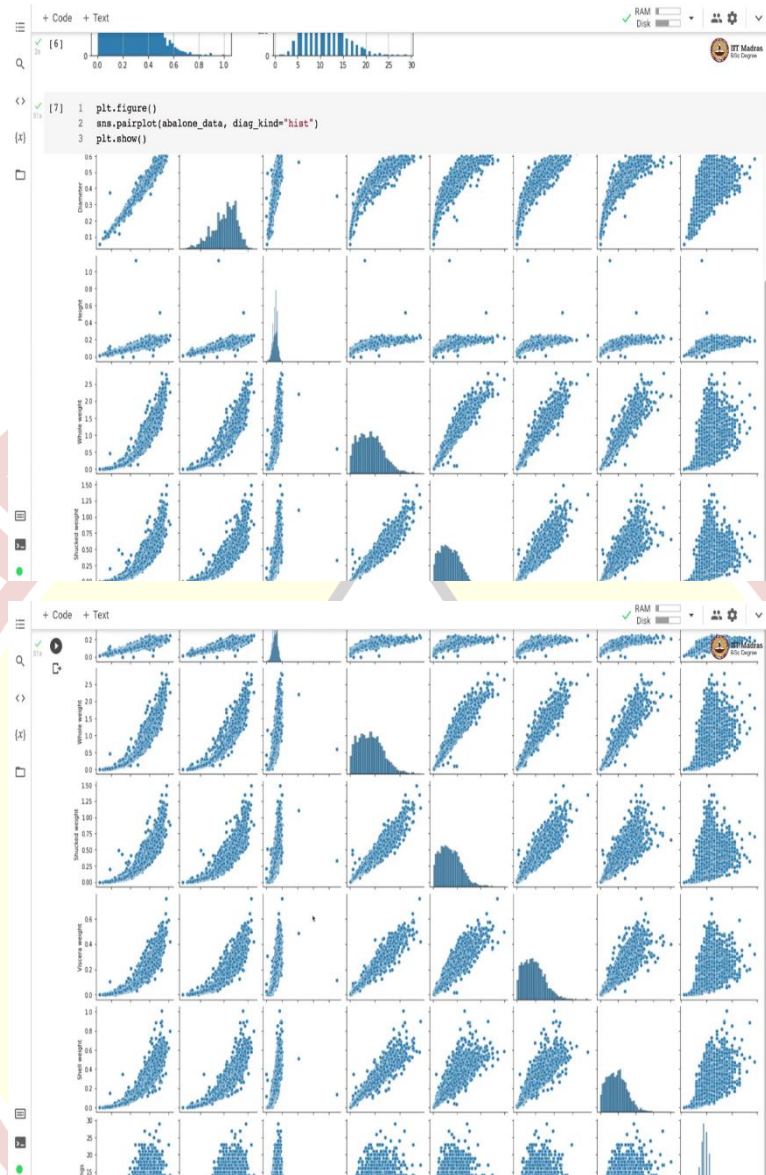
So, our data is in the format of the data frame. Let us examine the data. First thing to do is to look at the shape of the data. So, we have 4,177 examples in the dataset. And each example has got 9 features. Let us quickly examine the data you can see that the first 10 rows are printed on the screen. And you can check out the dataset or the first 10 sample rows on the screen.

(Refer Slide Time: 04:34)



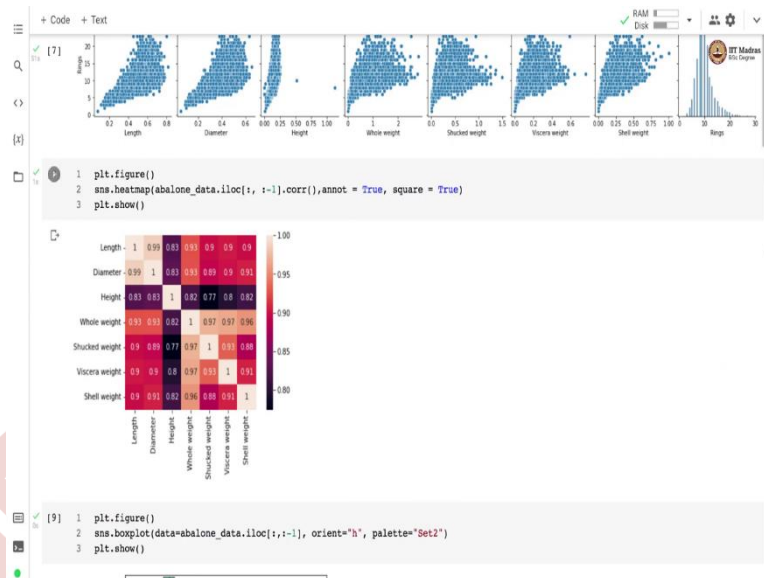
Next what we do is we plot the histogram of different features. And you can see that different features have got different distribution and different scales. And this is the distribution of rings which is our target variable. There are a lot of Abalone with the number of rings around 10. They seem to be quite abundant in this dataset, whereas, Abalone with rings > 15 are much smaller in this dataset

(Refer Slide Time: 05:21)



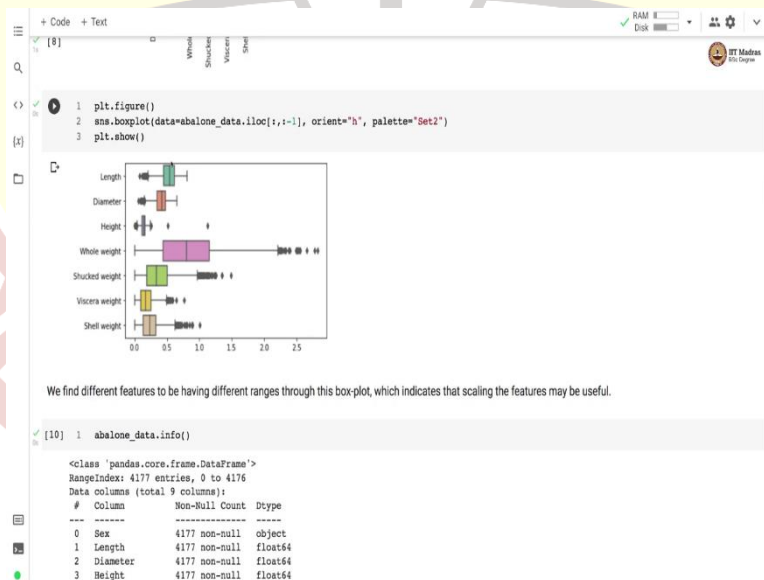
Then we perform pair plot between different features to see how these features are related to each other. And you can see that between certain features there seem to be very high correlation.

(Refer Slide Time: 05:38)



So, we can also establish correlation with heat map. And you can see that there is indeed a very strong correlation between certain features. For example, length and diameter there is a correlation of 0.99. Then between shucked weight and whole weight there is a correlation of 0.97. So, there is a lot of correlation between many features in this dataset.

(Refer Slide Time: 06:18)



Then we plot the distribution of different features in slightly different format, this time in a boxplot format. And we can see that, again different features have different ranges. And there are certain outliers as well. So, for example, in whole weight there are a bunch of outliers in shucked weight. So, you can see presence of outliers in each of these features.

(Refer Slide Time: 06:49)

```
+ Code + Text
We find different features to be having different ranges through this box-plot, which indicates that scaling the features may be useful.

1 abalone_data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4177 entries, 0 to 4176
Data columns (total 9 columns):
# Column      Non-Null Count  Dtype
---
0 Sex          4177 non-null   object
1 Length       4177 non-null   float64
2 Diameter     4177 non-null   float64
3 Height       4177 non-null   float64
4 Whole weight 4177 non-null   float64
5 Shucked weight 4177 non-null   float64
6 Viscera weight 4177 non-null   float64
7 Shell weight 4177 non-null   float64
8 Rings        4177 non-null   int64
dtypes: float64(7), int64(1), object(1)
memory usage: 293.8+ KB

From the information above, all features are continuous variables except for the Sex feature.

[11]: abalone_data.describe()

      Length  Diameter  Height  Whole weight  Shucked weight  Viscera weight  Shell weight  Rings
count  4177.000000  4177.000000  4177.000000  4177.000000  4177.000000  4177.000000  4177.000000  4177.000000
mean    0.523992    0.407881    0.139516    0.826742    0.359367    0.180594    0.238831    9.933684
std     0.120093    0.099240    0.041827    0.490389    0.221963    0.109614    0.139203    3.224169
min     0.075000    0.055000    0.000000    0.002000    0.001000    0.000500    0.001500    1.000000
```

So, we have already seen that in the dataset there are 4,177 entries, and there are 9 columns including the target variable which is rings for us. And all the, so most of the features are continuous numerical values. And 6 is one feature which is a categorical feature.

(Refer Slide Time: 07:17)

```
+ Code + Text

1 abalone_data.describe()

      Length  Diameter  Height  Whole weight  Shucked weight  Viscera weight  Shell weight  Rings
count  4177.000000  4177.000000  4177.000000  4177.000000  4177.000000  4177.000000  4177.000000  4177.000000
mean    0.523992    0.407881    0.139516    0.826742    0.359367    0.180594    0.238831    9.933684
std     0.120093    0.099240    0.041827    0.490389    0.221963    0.109614    0.139203    3.224169
min     0.075000    0.055000    0.000000    0.002000    0.001000    0.000500    0.001500    1.000000
25%     0.450000    0.350000    0.115000    0.441500    0.186000    0.063500    0.130000    8.000000
50%     0.545000    0.425000    0.140000    0.799500    0.336000    0.171000    0.234000    9.000000
75%     0.615000    0.480000    0.165000    1.153000    0.502000    0.253000    0.329000    11.000000
max     0.815000    0.650000    1.150000    2.825500    1.488000    0.760000    1.005000    29.000000

The count row shows that there are no missing values.
However, in the Height feature, the minimum value is zero. This possibility calls for a missing value in the data and we will process the missing value.

Next, take a look at the target in this case in the Rings column

[12]: abalone_data['Rings'].unique()

array([15, 7, 9, 10, 8, 20, 16, 19, 14, 11, 12, 18, 13, 5, 4, 6, 21,
       17, 22, 1, 3, 25, 23, 29, 2, 27, 25, 24])

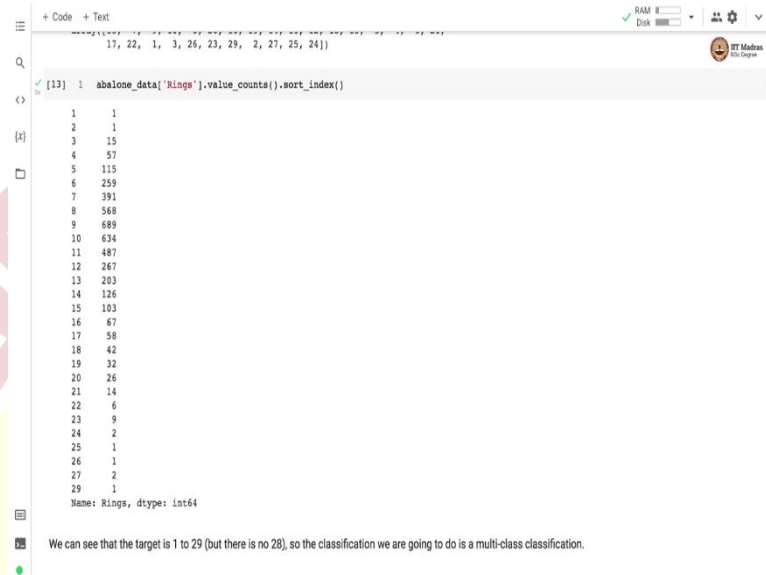
[13]: abalone_data['Rings'].value_counts().sort_index()

1      1
```

When we call describe on the data frame, and obtain different quartiles, we obtain mean, minimum, maximum and 25%, 50% and 75%. And we can examine this to check or reconfirm the presence of outliers. So, you can see here that the max over here is far higher than the 75th%, there are outliers in height, there are outliers in shucked weight, then also in shell weight.

And there are certain outliers probably on the other side where the values are much smaller than the 25th%. And you can see that in height the minimum value is 0. So, there is a possibility of a missing value in the data. And we need a special provision for processing these missing values.

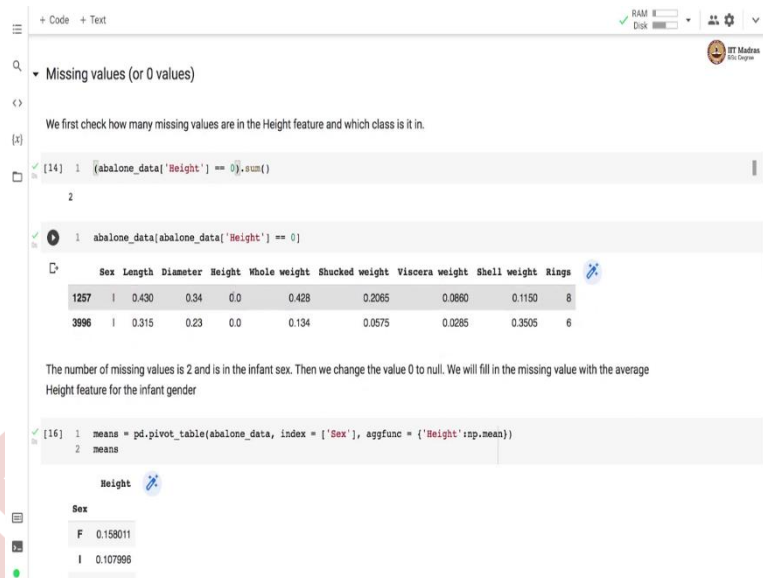
(Refer Slide Time: 08:16)



Then, let us look at a target column in this case, which is rings, and we count the number of rings. And we have seen this already in the plot that there are rings with count 10 or around 10, they form some kind of a majority in the class, but beyond 15 the number of rings are much lesser.

And here there is there are no rings present with number 28. So, there is another observation that you can make from this count. So, we are going to attempt a multi-class classification on this where there are 29 different classes and we will predict one of these 29 classes based on physical measurement of the Abalone.

(Refer Slide Time: 09:17)



Missing values (or 0 values)

We first check how many missing values are in the Height feature and which class is it in.

```
[14]: 1 (abalone_data['Height'] == 0).sum()
```

2

```
1 abalone_data[abalone_data['Height'] == 0]
```

| | Sex | Length | Diameter | Height | Whole weight | Shucked weight | Viscera weight | Shell weight | Rings |
|------|-----|--------|----------|--------|--------------|----------------|----------------|--------------|-------|
| 1257 | I | 0.430 | 0.34 | 0.0 | 0.428 | 0.2065 | 0.0660 | 0.1150 | 8 |
| 3996 | I | 0.315 | 0.23 | 0.0 | 0.134 | 0.0575 | 0.0285 | 0.3505 | 6 |

The number of missing values is 2 and is in the infant sex. Then we change the value 0 to null. We will fill in the missing value with the average Height feature for the infant gender

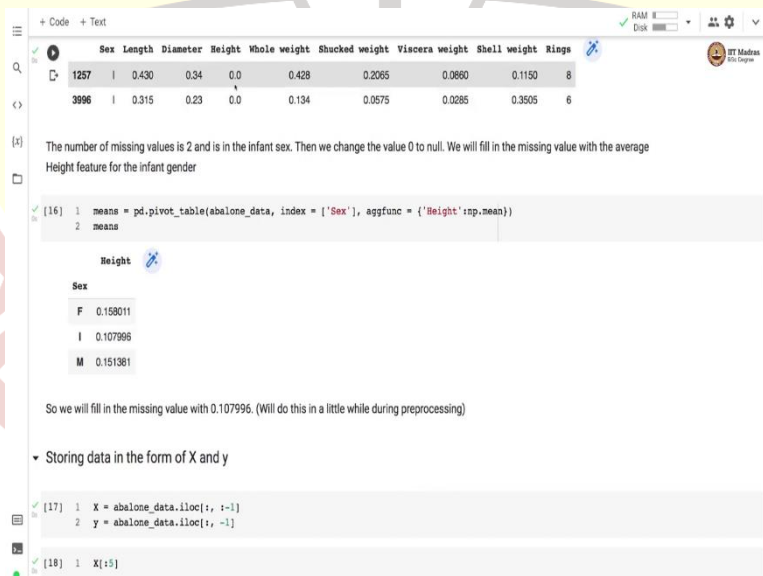
```
[16]: 1 means = pd.pivot_table(abalone_data, index = ['Sex'], aggfunc = ('Height':'np.mean'))
```

2 means

| | Height |
|-----|----------|
| Sex | |
| F | 0.158011 |
| I | 0.107996 |
| M | 0.151381 |

So, since there are missing values, we first handled the missing values, we found that in Abalone height column there are missing values and those missing values seem to be for 6 = infant and there are two missing values that you see here. So, we can fill up these missing values with average height feature of the infant gender.

(Refer Slide Time: 09:46)



| | Sex | Length | Diameter | Height | Whole weight | Shucked weight | Viscera weight | Shell weight | Rings |
|------|-----|--------|----------|--------|--------------|----------------|----------------|--------------|-------|
| 1257 | I | 0.430 | 0.34 | 0.0 | 0.428 | 0.2065 | 0.0660 | 0.1150 | 8 |
| 3996 | I | 0.315 | 0.23 | 0.0 | 0.134 | 0.0575 | 0.0285 | 0.3505 | 6 |

The number of missing values is 2 and is in the infant sex. Then we change the value 0 to null. We will fill in the missing value with the average Height feature for the infant gender

```
[16]: 1 means = pd.pivot_table(abalone_data, index = ['Sex'], aggfunc = ('Height':'np.mean'))
```

2 means

| | Height |
|-----|----------|
| Sex | |
| F | 0.158011 |
| I | 0.107996 |
| M | 0.151381 |

So we will fill in the missing value with 0.107996. (Will do this in a little while during preprocessing)

Storing data in the form of X and y

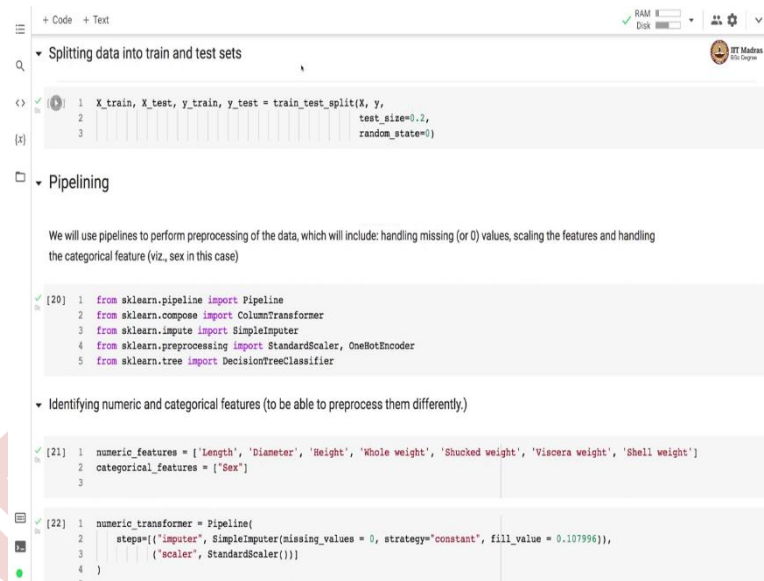
```
[17]: 1 X = abalone_data.iloc[:, :-1]
```

```
2 y = abalone_data.iloc[:, -1]
```

```
[18]: 1 X[:5]
```

So, here we calculate the mean height for different genders. And since there are missing values for gender I, we use its mean value to fill up these missing values.

(Refer Slide Time: 10:04)



The screenshot shows a Jupyter Notebook with the following content:

- Splitting data into train and test sets**
Code cell [1]:

```
1 X_train, X_test, y_train, y_test = train_test_split(X, y,  
2 test_size=0.2,  
3 random_state=0)
```
- Pipelining**
Text cell: "We will use pipelines to perform preprocessing of the data, which will include: handling missing (or 0) values, scaling the features and handling the categorical feature (viz, sex in this case)"
Code cell [20]:

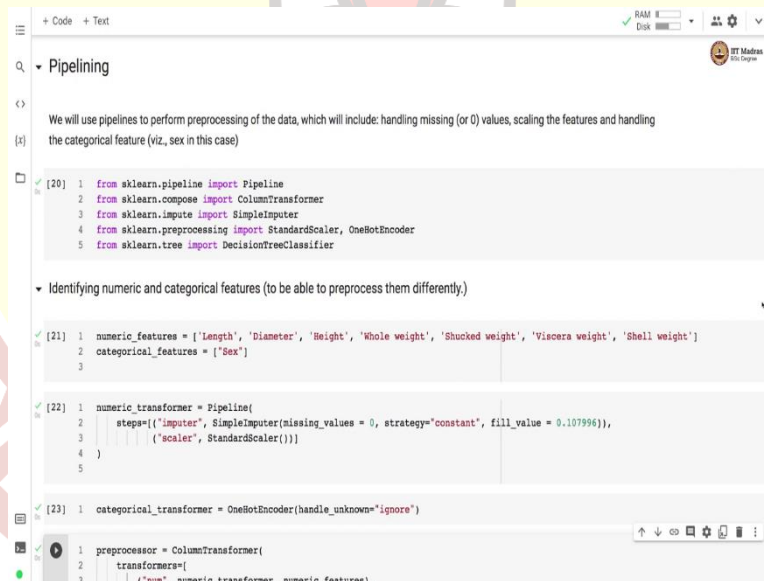
```
1 from sklearn.pipeline import Pipeline  
2 from sklearn.compose import ColumnTransformer  
3 from sklearn.impute import SimpleImputer  
4 from sklearn.preprocessing import StandardScaler, OneHotEncoder  
5 from sklearn.tree import DecisionTreeClassifier
```
- Identifying numeric and categorical features (to be able to preprocess them differently)**
Code cell [21]:

```
1 numeric_features = ['Length', 'Diameter', 'Height', 'Whole weight', 'Shucked weight', 'Viscera weight', 'Shell weight']  
2 categorical_features = ['Sex']  
3
```
- Code cell [22]:

```
1 numeric_transformer = Pipeline(  
2 steps=[('imputer', SimpleImputer(missing_values = 0, strategy='constant', fill_value = 0.107996)),  
3 ('scaler', StandardScaler())]  
4 )  
5
```

Next thing that we do is we split the data into training and test, we set aside 20% examples as test examples. And now we will define a model training pipeline.

(Refer Slide Time: 10:19)



The screenshot shows the continuation of the Jupyter Notebook with the following content:

- Pipelining**
Text cell: "We will use pipelines to perform preprocessing of the data, which will include: handling missing (or 0) values, scaling the features and handling the categorical feature (viz, sex in this case)"
- Code cell [20]:

```
1 from sklearn.pipeline import Pipeline  
2 from sklearn.compose import ColumnTransformer  
3 from sklearn.impute import SimpleImputer  
4 from sklearn.preprocessing import StandardScaler, OneHotEncoder  
5 from sklearn.tree import DecisionTreeClassifier
```
- Identifying numeric and categorical features (to be able to preprocess them differently)**
Code cell [21]:

```
1 numeric_features = ['Length', 'Diameter', 'Height', 'Whole weight', 'Shucked weight', 'Viscera weight', 'Shell weight']  
2 categorical_features = ['Sex']  
3
```
- Code cell [22]:

```
1 numeric_transformer = Pipeline(  
2 steps=[('imputer', SimpleImputer(missing_values = 0, strategy='constant', fill_value = 0.107996)),  
3 ('scaler', StandardScaler())]  
4 )  
5
```
- Code cell [23]:

```
1 categorical_transformer = OneHotEncoder(handle_unknown='ignore')  
2  
3
```
- Code cell [24]:

```
1 preprocessor = ColumnTransformer(  
2 transformers=[  
3 ('num', numeric_transformer, numeric_features),  
4 ('cat', categorical_transformer, categorical_features)]  
5 )
```

We will use pipeline construct to perform pre-processing of the data. This includes handling missing values, scaling the features, and handling categorical features. So, here we have mixed type of features that are majority of the features which are numerical features and we have a categorical feature as well. So, for that, we will be using ColumnTransformer. So, first identify numerical and categorical features and list them out. So, we have numerical features listed in this list and categorical features are listed in the other list.

(Refer Slide Time: 10:57)

```
+ Code + Text
1 from sklearn.pipeline import Pipeline
2 from sklearn.compose import ColumnTransformer
3 from sklearn.impute import SimpleImputer
4 from sklearn.preprocessing import StandardScaler, OneHotEncoder
5 from sklearn.tree import DecisionTreeClassifier

[20]

[x] Identifying numeric and categorical features (to be able to preprocess them differently.)

[21] 1 numeric_features = ['Length', 'Diameter', 'Height', 'Whole weight', 'Shucked weight', 'Viscera weight', 'Shell weight']
2     categorical_features = ['Sex']
3

[22] 1 numeric_transformer = Pipeline(
2     steps=[('imputer', SimpleImputer(missing_values = 0, strategy='constant', fill_value = 0.107996)),
3           ('scaler', StandardScaler())]
4 )
5

[23] 1 categorical_transformer = OneHotEncoder(handle_unknown='ignore')

1 preprocessor = ColumnTransformer(
2     transformers=[
3         ('num', numeric_transformer, numeric_features),
4         ('cat', categorical_transformer, categorical_features),
5     ]
6 )

[25] 1 # Append classifier to preprocessing pipeline.
2     # Now we have a full prediction pipeline.
3     clf = Pipeline(
4         steps=[('preprocessor', preprocessor),
5               ('classifier', DecisionTreeClassifier(
6                   max_depth=3, random_state=42))]
7     )
```

And we define numeric transformers and categorical transformers. So, numeric transformers are used to process the numerical features and categorical transformers will process the categorical features. So, here in numeric transformers, we have SimpleImputer and StandardScaler as two steps in the pre-processing transformation, the SimpleImputer uses the constant strategy and fill up the missing values with the values specified in fill_value parameter. In categorical transformer, we perform OneHotEncoder.

(Refer Slide Time: 11:41)

```
+ Code + Text

[22] 2 steps=[('imputer', SimpleImputer(missing_values = 0, strategy='constant', fill_value = 0.107996)),
3         ('scaler', StandardScaler())]
4 )
5

[23] 1 categorical_transformer = OneHotEncoder(handle_unknown='ignore')

1 preprocessor = ColumnTransformer(
2     transformers=[
3         ('num', numeric_transformer, numeric_features),
4         ('cat', categorical_transformer, categorical_features),
5     ]
6 )

[25] 1 # Append classifier to preprocessing pipeline.
2     # Now we have a full prediction pipeline.
3     clf = Pipeline(
4         steps=[('preprocessor', preprocessor),
5               ('classifier', DecisionTreeClassifier(
6                   max_depth=3, random_state=42))]
7     )

[26] 1 clf.fit(X_train, y_train)
2     print("model score: %.3f" % clf.score(X_test, y_test))

model score: 0.245

[27] 1 y_pred = clf.predict(X_test)

Let us compare the actual and predicted values of y.
```

Then we define a ColumnTransformer which has got two transformers one is the numerical transformer which is applied on numeric features and categorical transformer which is applied on categorical features.

(Refer Slide Time: 11:56)

```
+ Code + Text
1 preprocessor = ColumnTransformer(
2     transformers=[
3         ('num', numeric_transformer, numeric_features),
4         ('cat', categorical_transformer, categorical_features),
5     ]
6 )

[25] 1 # Append classifier to preprocessing pipeline.
2     # Now we have a full prediction pipeline.
3     clf = Pipeline(
4         steps=[('preprocessor', preprocessor),
5               ('classifier', DecisionTreeClassifier(
6                   max_depth=3, random_state=42))]
7     )

[26] 1 clf.fit(X_train, y_train)
2     print("model score: %.3f" % clf.score(X_test, y_test))

model score: 0.245

[27] 1 y_pred = clf.predict(X_test)

Let us compare the actual and predicted values of y.

[28] 1 comparison = np.concatenate((y_pred.reshape(len(y_pred),1), y_test.values.reshape(len(y_test),1)),1)
2
3     for each in comparison:
4         print(each)

[ 7 10]
[ 8  9]
[11 17]
[10 14]
```

Now, we will define our model along with the pre-processor in the pipeline construct. So, we have a pipeline with two steps one is a pre-processor and second is a DecisionTreeClassifier with max_depth set to 3. So, we will be building a tree with a maximum depth of 3. So, we train the tree by calling the fit function on the pipeline object with feature metrics and training label vector.

We can see that on the training set, we obtained the score of 0.245. We perform the prediction by calling the predict function on the pipeline object by supplying the test feature matrix and we obtain the predictions for each of the example in the test set.

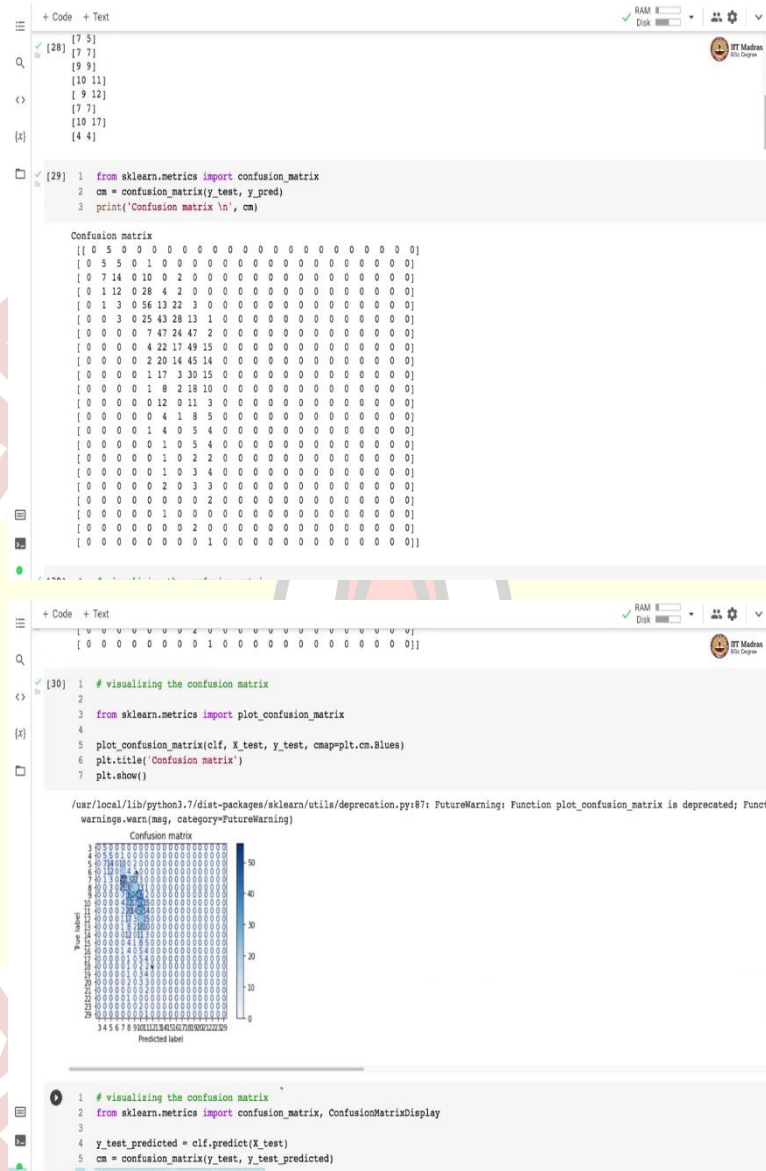
(Refer Slide Time: 12:53)

```
+ Code + Text
1 comparison = np.concatenate((y_pred.reshape(len(y_pred),1), y_test.values.reshape(len(y_test),1)),1)
2
3     for each in comparison:
4         print(each)

[ 7 10]
[ 8  9]
[11 17]
[10 14]
[ 7  7]
[ 7 16]
[ 9  9]
[ 8  8]
[ 7 13]
[10 10]
[ 8 12]
[ 8  8]
[11 11]
[11 12]
[10  9]
[10 14]
[ 9  9]
[10 11]
[10  9]
[ 4  3]
[ 8 14]
[ 8 10]
[ 7  7]
[ 8  7]
[ 5  6]
[ 7  7]
[ 9  7]
[ 8  9]
[10  9]
[ 7  7]
[ 8  6]
```


So, here, we have printed out, here we have printed out the actual label and the predicted label.

(Refer Slide Time: 13:03)



The image displays two screenshots of a Jupyter Notebook interface. The top screenshot shows the calculation of a confusion matrix using the `sklearn.metrics.confusion_matrix` function. The code is as follows:

```
[29] 1 from sklearn.metrics import confusion_matrix
      2 cm = confusion_matrix(y_test, y_pred)
      3 print('Confusion matrix \n', cm)
```

The output is a 26x26 confusion matrix, where each row represents an actual class and each column represents a predicted class. The matrix is mostly zeros, with some non-zero values along the diagonal and a few off-diagonal elements, indicating some misclassification.

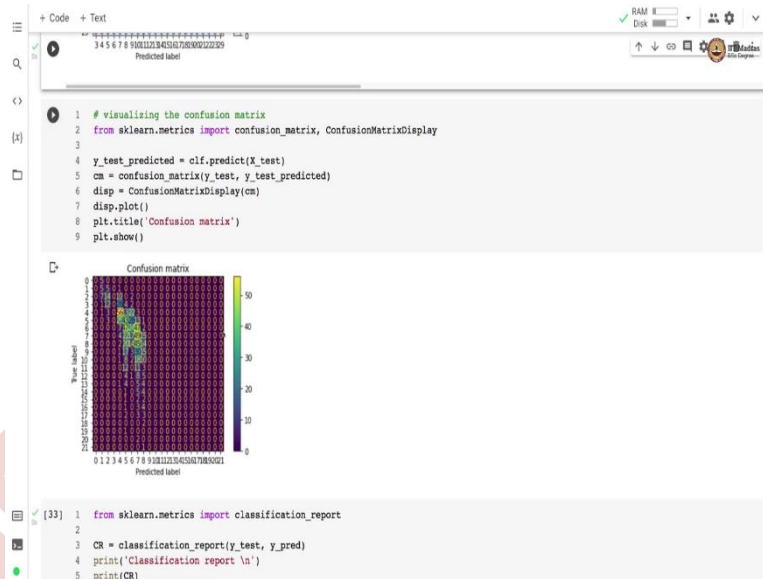
The bottom screenshot shows the visualization of the confusion matrix using the `sklearn.metrics.plot_confusion_matrix` function. The code is as follows:

```
[30] 1 #visualizing the confusion matrix
      2
      3 from sklearn.metrics import plot_confusion_matrix
      4
      5 plot_confusion_matrix(clf, X_test, y_test, cmap=plt.cm.Blues)
      6 plt.title('Confusion matrix')
      7 plt.show()
```

The output is a heatmap visualization of the confusion matrix. The x-axis is labeled 'Predicted label' and the y-axis is labeled 'True label'. The color scale ranges from 0 (dark blue) to 50 (light blue). The diagonal elements are the most prominent, showing high counts for correct classifications. There are also some off-diagonal elements, indicating misclassifications.

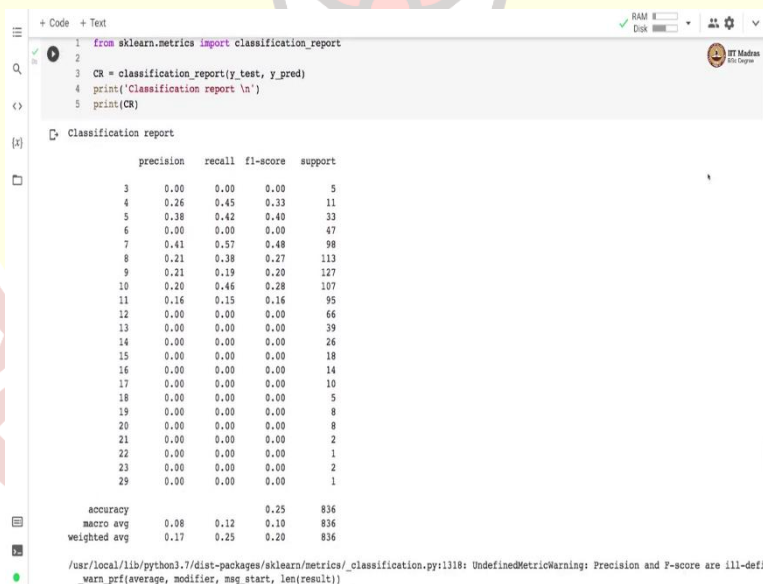
And we will actually calculate a confusion matrix and visualize the confusion matrix. So, you can see that there seem to be some kind of confusion in this particular part between different types of rings.

(Refer Slide Time: 13:25)



So, let us visualize the confusion matrix. And this visual confusion matrix shows the amount of confusion and you can see that in this particular region there is confusion that a lot of zeros because test set does not contain any examples from these ring types.

(Refer Slide Time: 13:50)



So, let us get a classification report and get precision recall and f1-score by the classes. You can see that bunch of classes have got 0 precision recall. The highest f1-score is obtained for number of rings equal to 7. And overall accuracy is 25%.

(Refer Slide Time: 14:24)

```
+ Code + Text
weighted avg 0.17 0.25 0.20 836

/usr/local/lib/python3.7/dist-packages/sklearn/metrics/_classification.py:1318: UndefinedMetricWarning: Precision and F-score are ill-defined for weights not in range [0.0, 1.0]
warn_prf(average, modifier, msg_start, len(result))
/usr/local/lib/python3.7/dist-packages/sklearn/metrics/_classification.py:1318: UndefinedMetricWarning: Precision and F-score are ill-defined for weights not in range [0.0, 1.0]
warn_prf(average, modifier, msg_start, len(result))
/usr/local/lib/python3.7/dist-packages/sklearn/metrics/_classification.py:1318: UndefinedMetricWarning: Precision and F-score are ill-defined for weights not in range [0.0, 1.0]
warn_prf(average, modifier, msg_start, len(result))

Cross-validation

[34] 1 from sklearn.model_selection import cross_val_score
2 acc = cross_val_score(estimator = clf, X = X_train, y = y_train, cv = 10)
3 print(type(acc))
4 print('Accuracy of each fold ', list(acc*100))
5 print('Accuracy: {:.2f} %'.format(acc.mean()*100))

/usr/local/lib/python3.7/dist-packages/sklearn/model_selection/_split.py:680: UserWarning: The least populated class in y has only 1 member
UserWarning,
<class 'numpy.ndarray'>
Accuracy of each fold [27.46268656716418, 22.45508982035928, 23.952095808383234, 24.550898203592812, 24.251497005988025, 23.952095808383234]
Accuracy: 25.53 %

Visualizing the decision tree

[35] 1 from sklearn import tree
2
3 #plt the figure, setting a black background
4 plt.figure(figsize=(48,8), facecolor='w')
5
6 #create the tree plot
```

Let us also perform cross-validation and opt in the cross-validation score. So, the accuracy in cross-validation is 25.53%. Now, next what we will do is we can visualize the decision tree that was learned.

(Refer Slide Time: 14:41)

```
+ Code + Text

/usr/local/lib/python3.7/dist-packages/sklearn/model_selection/_split.py:680: UserWarning: The least populated class in y has only 1 member
UserWarning,
<class 'numpy.ndarray'>
Accuracy of each fold [27.46268656716418, 22.45508982035928, 23.952095808383234, 24.550898203592812, 24.251497005988025, 23.952095808383234]
Accuracy: 25.53 %

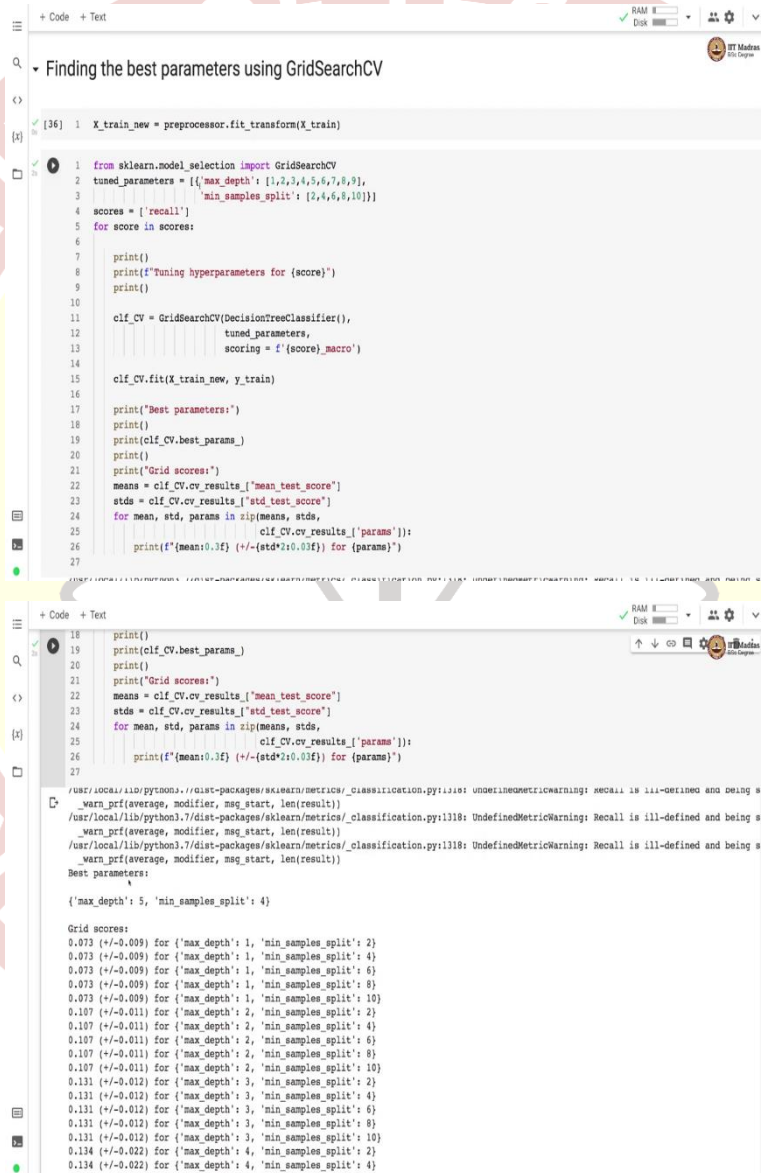
Visualizing the decision tree

1 from sklearn import tree
2
3 #plt the figure, setting a black background
4 plt.figure(figsize=(48,8), facecolor='w')
5
6 #create the tree plot
7 a = tree.plot_tree(clf['classifier'],
8 #use the feature names stored
9 feature_names = column_names,
10 rounded = True,
11 filled = True,
12 fontsize=12)
13 #show the plot
14 plt.show()
```

So, for that we use a tree API from sklearn and called plot _tree on the tree class. So, the learn tree is on the screen, so you can see that a top-level node is viscera weight and we have printed the gini index as well as the number of samples and the values over here. So, as you go down the sub trees, the gini index has improved and you can see that the number of samples.

So, here we got a split where 930 samples were pushed on the left and the remaining 2,411 samples are pushed to the right. So, in the second step, we use length as a feature where we want to perform the split, here we use viscera weight, again as a feature to perform the split and so on. So, you can see that the leaf node, the first leaf node has got 63 samples, second one has got 100 samples and so on, and you can see the values of the target in each of the leaf node which is also printed over here.

(Refer Slide Time: 16:14)



```
[36] 1 X_train_new = preprocessor.fit_transform(X_train)

2 from sklearn.model_selection import GridSearchCV
3 tuned_parameters = [{'max_depth': [1,2,3,4,5,6,7,8,9],
4   'min_samples_split': [2,4,6,8,10]}]
5 scores = ['recall']
6 for score in scores:
7     print()
8     print(f'Tuning hyperparameters for {score}')
9     print()
10
11     clf_cv = GridSearchCV(DecisionTreeClassifier(),
12                           tuned_parameters,
13                           scoring = f'{score}_macro')
14
15     clf_cv.fit(X_train_new, y_train)
16
17     print("Best parameters:")
18     print()
19     print(clf_cv.best_params_)
20     print()
21     print("Grid scores:")
22     means = clf_cv.cv_results_["mean_test_score"]
23     stds = clf_cv.cv_results_["std_test_score"]
24     for mean, std, params in zip(means, stds,
25                                 clf_cv.cv_results_["params"]):
26         print(f'{mean:0.3f} +/- (std*2):{0.03f} for {params}')
27
```

```
18 print()
19 print(clf_cv.best_params_)
20 print()
21 print("Grid scores:")
22 means = clf_cv.cv_results_["mean_test_score"]
23 stds = clf_cv.cv_results_["std_test_score"]
24 for mean, std, params in zip(means, stds,
25                             clf_cv.cv_results_["params"]):
26     print(f'{mean:0.3f} +/- (std*2):{0.03f} for {params}')
27
```

```
Best parameters:
{'max_depth': 5, 'min_samples_split': 4}

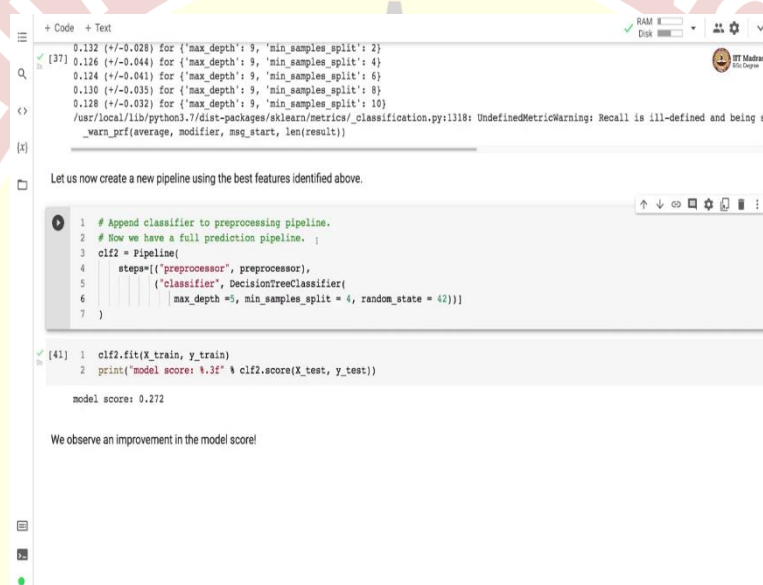
Grid scores:
0.073 +/-0.009 for {'max_depth': 1, 'min_samples_split': 2}
0.073 +/-0.009 for {'max_depth': 1, 'min_samples_split': 4}
0.073 +/-0.009 for {'max_depth': 1, 'min_samples_split': 6}
0.073 +/-0.009 for {'max_depth': 1, 'min_samples_split': 8}
0.073 +/-0.009 for {'max_depth': 1, 'min_samples_split': 10}
0.107 +/-0.011 for {'max_depth': 2, 'min_samples_split': 2}
0.107 +/-0.011 for {'max_depth': 2, 'min_samples_split': 4}
0.107 +/-0.011 for {'max_depth': 2, 'min_samples_split': 6}
0.107 +/-0.011 for {'max_depth': 2, 'min_samples_split': 8}
0.107 +/-0.011 for {'max_depth': 2, 'min_samples_split': 10}
0.131 +/-0.012 for {'max_depth': 3, 'min_samples_split': 2}
0.131 +/-0.012 for {'max_depth': 3, 'min_samples_split': 4}
0.131 +/-0.012 for {'max_depth': 3, 'min_samples_split': 6}
0.131 +/-0.012 for {'max_depth': 3, 'min_samples_split': 8}
0.131 +/-0.012 for {'max_depth': 3, 'min_samples_split': 10}
0.134 +/-0.022 for {'max_depth': 4, 'min_samples_split': 2}
0.134 +/-0.022 for {'max_depth': 4, 'min_samples_split': 4}
0.134 +/-0.022 for {'max_depth': 4, 'min_samples_split': 6}
0.134 +/-0.022 for {'max_depth': 4, 'min_samples_split': 8}
0.134 +/-0.022 for {'max_depth': 4, 'min_samples_split': 10}
```

Now, as usual, there are certain configurable parameters in decision tree, like maximum depth and then minimum number of samples to split. So, the depth we want to expand with the depth between 1 to 9 and the min _samples _split between 2 to 10. We are going to use recall as a major.

So, here we perform a grid search by defining GridSearchCV object with DecisionTreeClassifier as the estimator, then the grid, which is specified in tuned _parameters. And using the recall _macro, as a scoring measure. We call the fit function on the GridSearchCV object with the training feature matrix and label vector as input.

The base parameters are printed with by accessing best _params _ member variable of the GridSearchCV object. You also print the mean _test _score under standard deviation of the test score. So, you can see that the best parameters are obtained with max_depth equal=5 and min _samples _split = 4.

(Refer Slide Time: 17:45)



The screenshot shows a Jupyter Notebook interface. The top cell displays the results of a GridSearchCV search for a DecisionTreeClassifier. The output lists several parameter combinations with their corresponding mean test scores and standard deviations. The best parameters identified are max_depth=5 and min_samples_split=4, achieving a mean score of approximately 0.132. A warning message is also visible: 'UndefinedMetricWarning: Recall is ill-defined and being a _warn_prf(average, modifier, msg_start, len(result))'.

The second cell contains a text prompt: 'Let us now create a new pipeline using the best features identified above.' Below this, a new pipeline is created and fitted. The code defines a Pipeline with a DecisionTreeClassifier as the estimator, using the best parameters found in the previous search. The pipeline is fitted on the training data, and the model score is printed, resulting in 0.272.

```
[37] 0.132 (+/-0.028) for {'max_depth': 5, 'min_samples_split': 2}
      0.126 (+/-0.044) for {'max_depth': 5, 'min_samples_split': 4}
      0.124 (+/-0.041) for {'max_depth': 5, 'min_samples_split': 6}
      0.130 (+/-0.035) for {'max_depth': 5, 'min_samples_split': 8}
      0.128 (+/-0.032) for {'max_depth': 5, 'min_samples_split': 10}
/usr/local/lib/python3.7/dist-packages/sklearn/metrics/_classification.py:1318: UndefinedMetricWarning: Recall is ill-defined and being a
_warn_prf(average, modifier, msg_start, len(result))

Let us now create a new pipeline using the best features identified above.

1 # Append classifier to preprocessing pipeline.
2 # Now we have a full prediction pipeline.
3 clf2 = Pipeline(
4     steps=[('preprocessor', preprocessor),
5           ('classifier', DecisionTreeClassifier(
6               max_depth=5, min_samples_split=4, random_state=42))]
7 )

[41] 1 clf2.fit(X_train, y_train)
     2 print('model score: %.3f' % clf2.score(X_test, y_test))

model score: 0.272

We observe an improvement in the model score!
```

Now, we create a new pipeline with the best features identified above. So, we set the max _depth = 5 and min _samples _split = 4. And we define the pipeline object and we fit the pipeline object with the training data. Here we obtained the model score of 0.27 which is improvement from what we obtained earlier which was 25%. So, here we have got the score of 0.27 as compared to earlier score of 0.25.

So, we observe an input improvement in the model score. So, as an exercise you should perform the evaluation of this classifier on the test set. Other interesting exercise to do is to also attempt this Abalone problem as a regression problem and see if you get a better model. In general, the accuracy is quite low on this dataset. So, this was a demonstration of decision trees with continuous and discrete features. We use decision tree to predict the number of rings in Abalone.