# IIT Madras

## ONLINE DEGREE

(Refer Slide Time: 0:10)





Namaste. Welcome to the next video of the machine learning practice course. In this video, we will demonstrate linear regression with sklearn API. The objective of this collab is to demonstrate how to build a linear regression model with sklearn. We will be using the California housing dataset for this exercise. We will be implementing a linear regression API that uses the normal equation for training the model. Then we will also train the model with cross-validation.
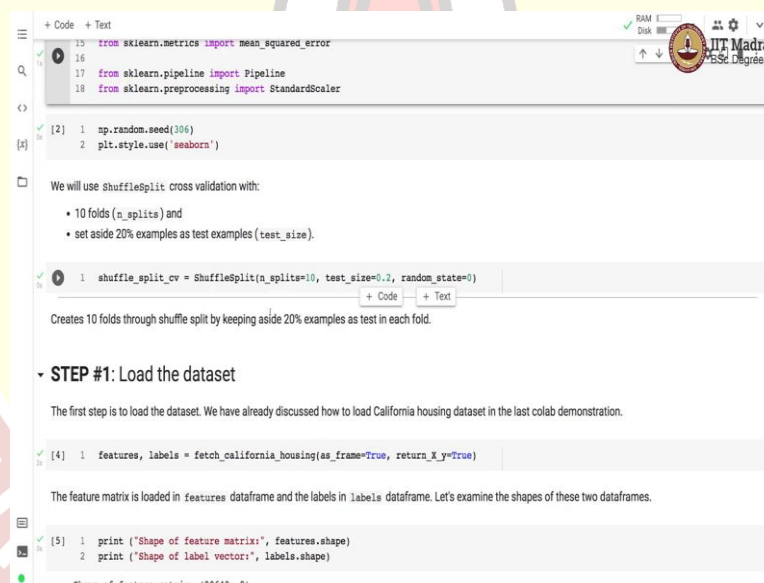
After the model is trained, we will evaluate the model with a score function that uses r2 score. We will also use cross val score with different scoring parameters in order to evaluate the

training model with different scoring metrics. We will also study the model diagnosis with learning curve and how to examine the learned model or weight vector. Let us start by importing some useful packages.

Here we are importing some basic Python packages for manipulation and visualization. Then we are, then we import fetch California Housing from sklearn. dataset module that helps us to fetch the dataset. The linear regression, the linear regression API that we are importing from sklearn.linear_model helps us to build the linear regression model.

Then we have a bunch of APIs imported from model selection module and that helps us to perform cross -validation, plot the learning curve and performing the shuffle split, shuffle split in cross -validation. We are going to use mean square error as an additional metrics to r2 and mean square error is implemented in sklearn.metrics module. We will be using pipelining for combining the preprocessing and the model in a pipeline manner.

(Refer Slide Time: 2:39)



It is always a good idea to set a random seed so that we get reproducible results across different runs. So we are going to set the seed to 306, you can set it to any number of your choice and we are going to use the plotting style that is seaborn. First, we will use shuffle split cross-validation with 10 folds and we will set aside 20 percent examples as test examples.

So, you can see that here we are instantiating a shuffle split object with n splits equal to 10, test size equal to 0.2 and we have set a random seed. This particular statement creates 10 folds through shuffles split by keeping the side 20 percent examples as test in each fold. We will be using the shuffle split strategy throughout this notebook.

(Refer Slide Time: 3:41)



The first step is to load the data. We have already discussed how to load the California housing dataset in the last colab demonstration. So, here we use Fetch California Housing API. And we set as the frame is equal to true and return_x_y parameter to true. So, because of this parameter, we only get features and labels. And since we are setting as frame equal to true, we get both these return values as a data frame.

So we have a feature data frame and we have a label data frame. So let us examine the shapes of the feature matrix and the label vector. We can see that we have 20,640 examples each with 8 features. And in the label vector, we have 20,640 entries. As a sanity check, we need to make sure that the number of rows in the feature matrix and labels match. And we can do that with assert statement by looking at the zeroth element of our shape in the feature matrix and label vector.

The second step is to perform data exploration which has been covered in detail in the last colab demonstration where we studied the California Housing dataset.

(Refer Slide Time: 5:21)



In the next part of Step 3, where we are going to do preprocessing and model building, we start by splitting the training data into training and test sets. We do not access this data till the end, or data exploration and tuning are performed on the training set. And by setting aside a small portion of training set as a dev or validation set. So we perform train test split by using a training test split class from sklearn.model selection module.

The train test split splits the given feature matrix and label vector into training and test feature matrices and train and test label vectors. Again, in the similar manner as feature metrics and label vector, we also examine the shapes of training and test sets. Here we have got 15,480 examples in training and 5160 example in test.

We also perform a sanity check to make sure that the training feature matrix has the same number of rows as training the label vector. And also, we perform the same check in the test set 2. So, we use again assert statements to make sure that these conditions are met.

(Refer Slide Time: 7:00)



Next, we build a linear model with default parameter settings of linear regression API. We will make use of pipeline API for combining data preprocessing and model building. We will use a standard scaler for scaling all the features to the same scale followed by the linear regression model.

So, this pipeline that we have built here has two parts. The first part is feature scaling, where we use standard scalar objects or standard scalar class. Here we instantiate the object of standard scalar plus, which helps us to bring all features on the same scale.

Then in the second step of the pipeline, we have linear regression object instantiated, which helps us to perform linear regression with normal equation. Once this pipeline is constructed, we can train the pipeline by calling the Fit function on the pipeline object. And we supply the training feature metrics and training label vector as parameters to the Fit function.

(Refer Slide Time: 8:10)



```
+ Code  + Text
After constructing the pipeline object, let's train it with training set.

[10]  1  # set up the linear regression model.
      2  lin_reg_pipeline = Pipeline([("feature_scaling", StandardScaler()),
      3                              ("lin_reg", LinearRegression())])
      4
      5  # train linear regression model with normal equation.
      6  lin_reg_pipeline.fit(train_features, train_labels)

    Pipeline(steps=[('feature_scaling', StandardScaler()),
                    ('lin_reg', LinearRegression())])

Now that we have trained the model, let's check the learnt/estimated weight vectors (intercept_ and coef_).

  1  print("intercept (w_0):", lin_reg_pipeline[-1].intercept_)
  2  print("weight vector (w_1, ..., w_m):", lin_reg_pipeline[-1].coef_)

  intercept (w_0): 2.0703489205426377
  weight vector (w_1, ..., w_m): [ 0.85210815  0.12065533 -0.30210555  0.34860575 -0.00164465 -0.04116356
   -0.89314697 -0.86784046]

A couple of things to notice:

 • We accessed the LinearRegression object as lin_reg_pipeline[-1] which is the last step in the pipeline.
 • The intercept can be obtained via intercept_ member variable and
 • The weight vector corresponding to features via coef_.
```
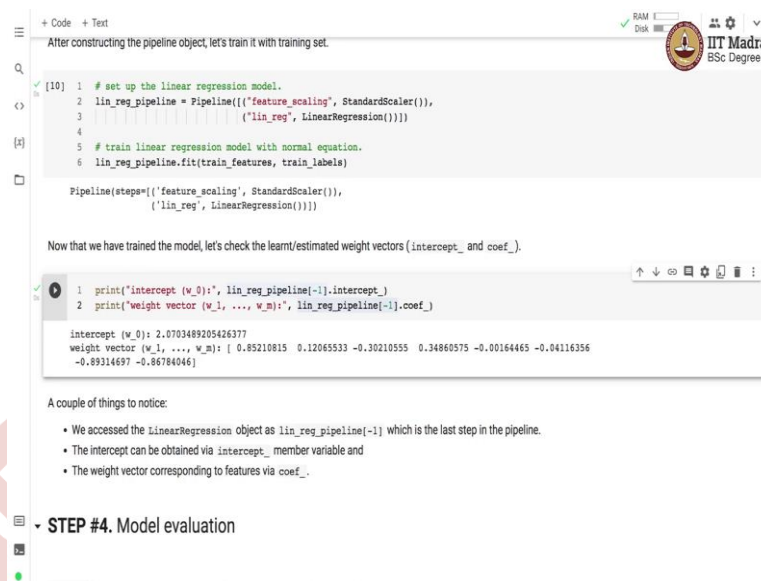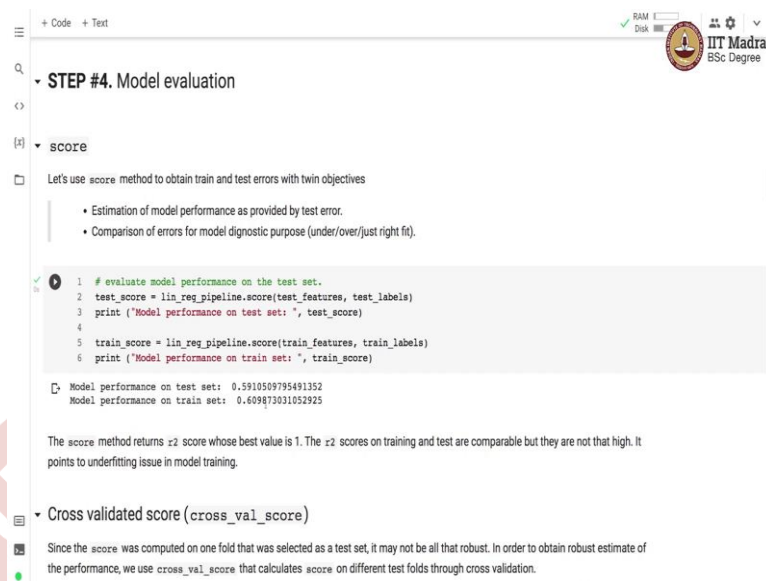
**STEP #4. Model evaluation**

Now that we have trained a model, let us check the estimated weight vectors. And these weight vectors can be accessed through the intercept_and coef_member variable of the linear regression object. So, here note that in the pipeline, we have two stages, first is feature scaling and the second is linear regression. So, we can access the linear regression object by using the index in the linear regression pipeline object of the pipeline class.

So, we can, we can access this particular estimator by lin_reg_pipeline of -1. And you are very well aware of how indexing works in Python. So, since we are trying to access -1 index, which is the last step, in this particular pipeline, which is linear regression. So, this particular expression gives us the linear regression object. And then on that object we call intercept, underscore, which is a member variable of this object that gives us the weight for the intercept, that is w0.

And when we access the coef_member variable of the same object, we get weights for the remaining parameters. So intercept, w0 is 2.07. And a weight vector that contains M entries w1, w2, all the way up to Wm, where Wm is a number of features. And since there are 8 features, you will see eight entries over here in the weight vector.
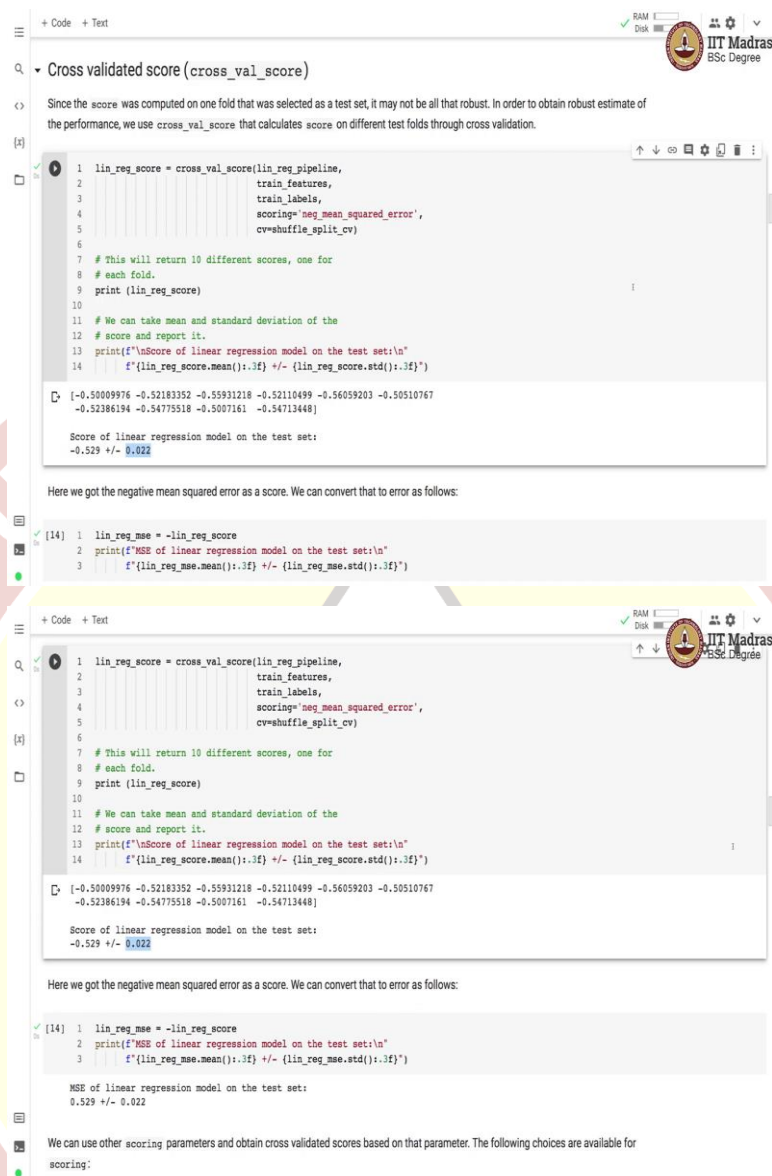
Now that model is trained, let us move on to the next step which is model evaluation. In evaluation, we will first use the scoring method to obtain train and test error with you between objectives. We will use the score obtained for estimation of model performance and at estimation is based on what we get in the test error. We also compare training and test error for model diagnostic purposes, where we want to figure out whether model is underfitting, overfitting or the model has got gotten the just tight fit.

So, what we do is, we call the score method on the pipeline object. So, what pipeline does it, it takes this feature matrix, first applies feature scaling and then called the predict and then it scores based on the predicted value and the actual value of the labels. So, we get this score and by default this score is r square or the coefficient of determination.

We perform the same thing with the training set. The score method returns r2 score. And since this is the score, the higher the better and you can see that we have slightly higher performance on the training set and the test set. And remember that the best value in r2 score is 1.

So, we are far from being the best, but we have a decent performance here. Another thing to note that the performance is comparable on test and training. But it is not that high and this could well point to the underfitting issue in the model.

(Refer Slide Time: 12:05)



In the score method, we have calculated the score on one fold that was selected as a test set. So that is why it may not be all that robust. In order to obtain a robust estimate of the performance, we use cross -validated score with cross_val_score API from sklearn.model_selection module. This helps us to calculate score on different test folds through cross -validation.

So, the way we call cross_val_score is by specifying the estimator, the training features, and labels, we can also specify what kind of scoring function we want to use, here we are specifying the negative mean squared error. So, this is negative of the mean squared error and this is scoring, so in scoring; the score higher the better and error if you are calculating error, the lower the error is a better error, is the better error.
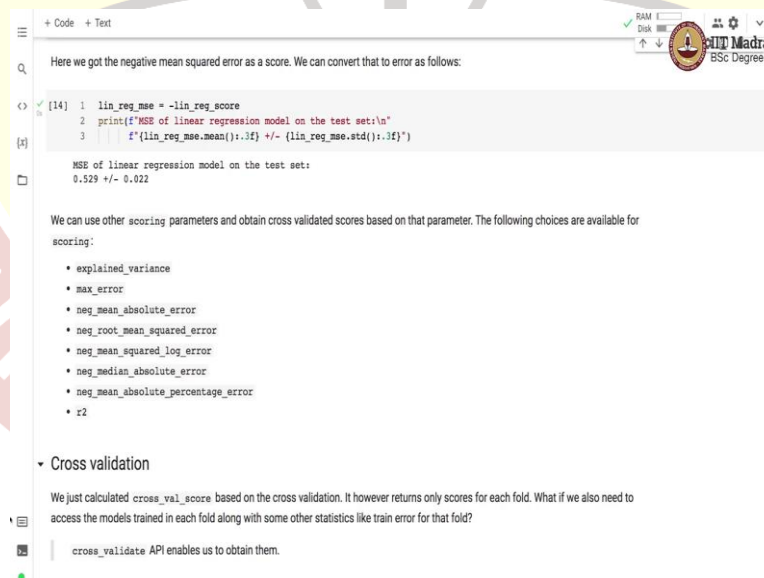
And we have seen this in our slides on this topic and we are going to use the shuffle split cv as a cross-validation strategy. Here we have initialized this cross shuffle split cv with 10 different folds. So, we will get 10 different scores one for each fold. For every fold, there are some examples that are set aside as test and we calculate this cross-validation score on that test set.

Since there are 10 different test sets, we get 10 different scores through this particular method. And these are 10 different scores that you can see here. And now since we have 10 different scores, we report the mean and standard deviation in this particular scores. So, we see that the negative mean squared error is the mean of that is -0.52. And the standard deviation is 0.022.

So, you can see that most of the scores are more or less comparable. And that is also that we can also verify through the standard deviation value which is fairly low. So, since we were calculating the score which is negative of the mean squared error, we got this particular number.

Now we want to convert let us say this score to the error we can convert it as follows. So we can take this particular score multiplied by -1 and we get the mean squared error. Now we can see that the mean square error is 0.52 and the standard deviation is 0.02. Here we use negative mean squared error as a scoring parameter.

(Refer Slide Time: 15:26)



There are some more choices that are available as scoring parameter, which is like explained variance, max error, negative mean absolute error, negative root mean squared error, negative mean squared log error, negative mean absolute error, negative mean absolute percentage error and r2. You have seen these, we have seen the choices in a bit more detail in the presentation on this topic.

(Refer Slide Time: 15:58)



So, the cross Val score just gave us score based on cross-validation. If we want to let us say access, the model is trained in each fold along with some of the statistics, like train error, we can use cross-validation API. And cross-validation API has pretty much the same parameters like cross val score, it takes the estimator, then the training, training feature metrics, training labels, the cross-validation strategy which is shuffled split cv, the scoring metric which is negative mean squared error.

And we also set these two parameters to true which are return train score and return estimator. So, return estimator when we set to true, we also get the estimators that were trained on the different fold. And when we said this particular parameter, we also get the score on the training folds.

And remember that these two by default are false since we need to explicitly set them to true. So as a result, we get a dictionary object that has got trained estimator time taken for fitting and scoring the models in cross-validation, we also have training score and the test score.

(Refer Slide Time: 17:21)



Let us print the content of the dictionary for us to examine. So, you can see that there is a key and a value. So there is this estimator, which has got 10 different estimators that were trained on 10 different folds, we have fit times which is again a list of 10 numbers. This was this is a time that was taken to fit the model in each fold, then the score time, in the same manner, the time it took to calculate the score in each fold.

Then we also get 10 different numbers here, which are test scores, as well as the train scores over here. And remember, these are the scores, and the higher the scores, the better the performance. So, you can see that our performance, in general, is better on the training set than the test set which is what you normally get.

(Refer Slide Time: 18:22)

We will compare the training and test error to assess the generalization performance of our model. However, we have training, training and test scores in cv results dictionary, we multiply these errors, we multiply the scores by -1 and convert them to errors. So, we multiply train score and test score by -1 and that way we have now the errors. So, the mean squared error for the training set is 0.51 and for the test set it is 0.52

So, in training set there is a smaller standard deviation than the test error. So, training has got a better error or the lower error, hence training performs better than the test. And these scores are comparable, but they are still high. And that actually points us to possibly the problem of underfitting. We will confirm that by plotting the learning curves.

(Refer Slide Time: 19:34)



The learning curves have got slightly different meanings in the context of sklearn. In the machine learning techniques course, we talked about learning curves as the iteration on the x-axis and the laws on the y axis. So, here this is a slightly different kind of learning curve where we have the training set size on x-axis and we have the error on Y-axis. So, what has really changed is the x-axis. Here, instead of iteration, we have training set size as the x-axis.

So for this, we will be using; for this, we will be using learning curve API that calculates cross-validation scores for the different number of samples as specified in the argument which is training_sizes. Let us check out the learning curve API. Meanwhile, here, there is a code snippet for plotting the learning curve and we are hiding it because this has got code for plotting the learning curves.

So, let us see how to call the learning curve API. So, the learning curve has got estimator training, and test feature metrics and labels. In the cross-validation strategy, the scoring metric is the number of jobs, which specifies whether we should be using parallelism and how much parallelism we should be using in this particular API. And we set return times to true and we specified training sizes. So, all these arguments up to this point, are exactly the same as cross val score or cross-validate.

So, we now have return times and training sizes and training sizes is the most important argument because here we have to specify what is the range of values that we should be used in training. So, we use values between 0.2 and 1 and there is 10 such kind of values that will be used.

We get some such kind of plot. After plotting the training examples versus mean squared error, we have training examples on the x-axis and mean squared error on the y-axis. So, you can see that there are 10 different points with varying numbers of training examples, and for each of these points, there is corresponding mean square error that you are plotting on the y axis. So, we can see that both curves have reached a plateau and are close and fairly high.

So, a few instances in the training means that model can fit those examples perfectly. And we can see that when the model is trained on a few instances, the error on the test set is quite high. And as we add more training examples, the test error gradually comes down. So, this learning curve is a type of underfitting model.

We also study how the model training scales as a function of number of examples. So, we have fit times we and we plot the training examples versus fit times. And you can see that as the number of training example grows, the time that it takes to fit the model also increases.

(Refer Slide Time: 23:11)





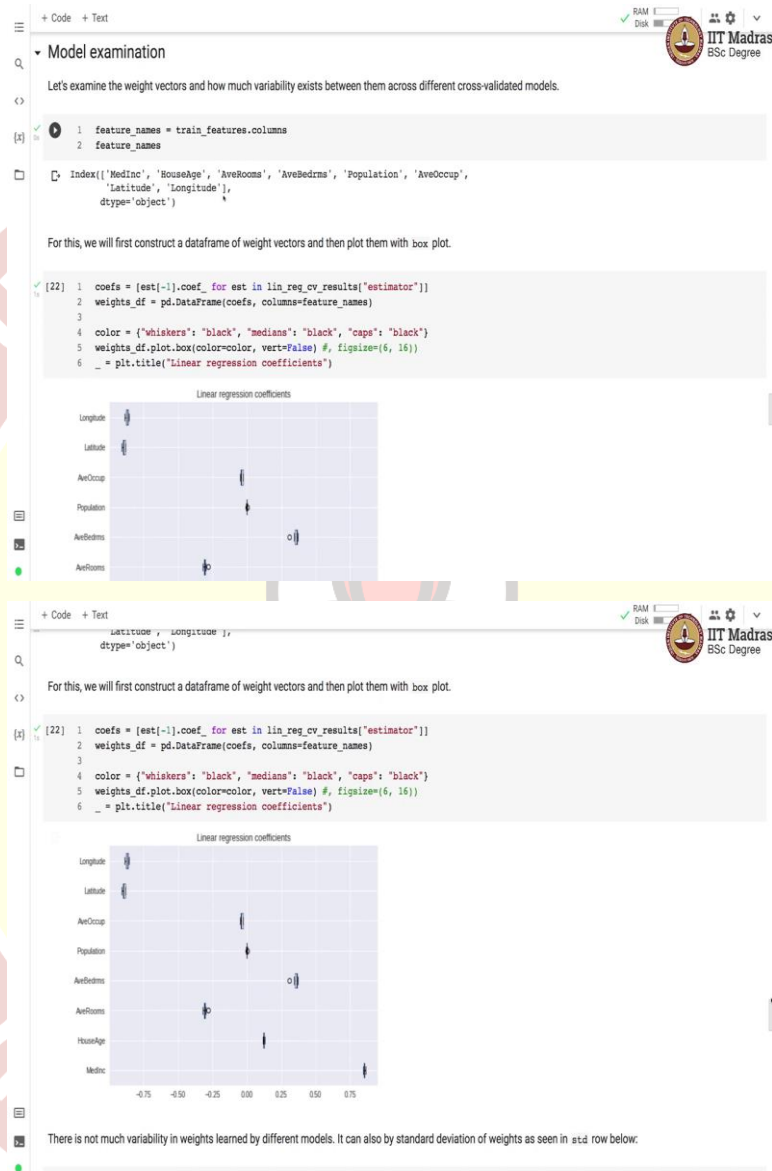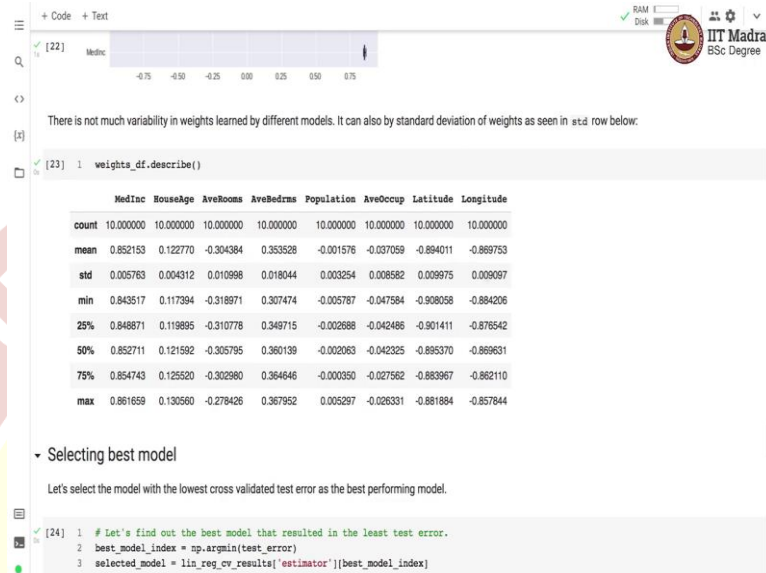Let us examine the weights and the variability in the weight vector across different folds in cross validation. So, here, we have a boxplot for linear regression coefficients that we obtain in in different folds. So, you can see that this is a boxplot for the weights that are assigned to longitude across different folds and the variability is quite small, because box plot is quite compact here.
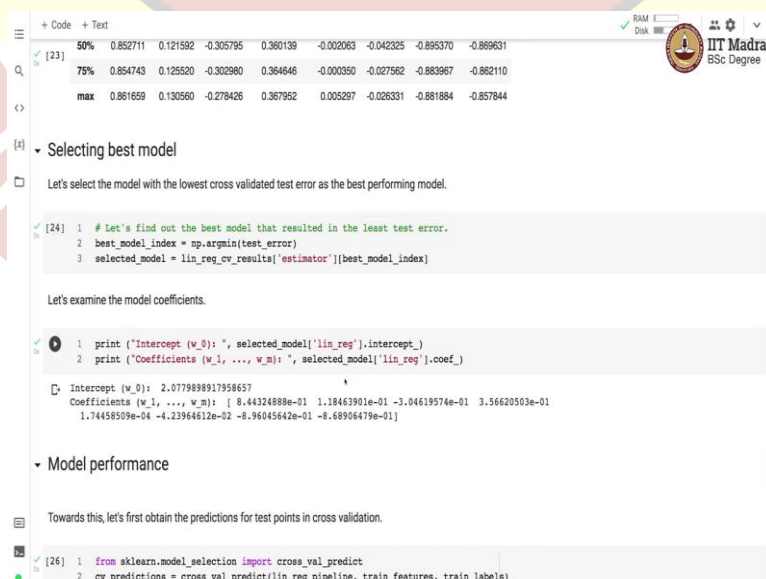
Whereas, in average bedroom, average rooms and population there are there seems to be some outlier weights. In most of the cases, the boxplot is quite compact, that means there is a lesser variability in the weights that we are obtained, that we obtained across different folds.

(Refer Slide Time: 24:21)



And we can check that by also calling the description on the weights data frame. And you can see that the standard deviation in most of the cases is fairly small, except for these two cases it is higher than the 10 other cases. And you can also see the presence of outliers in the case of its population where the seventy-fifth percentile and the max have some difference right.
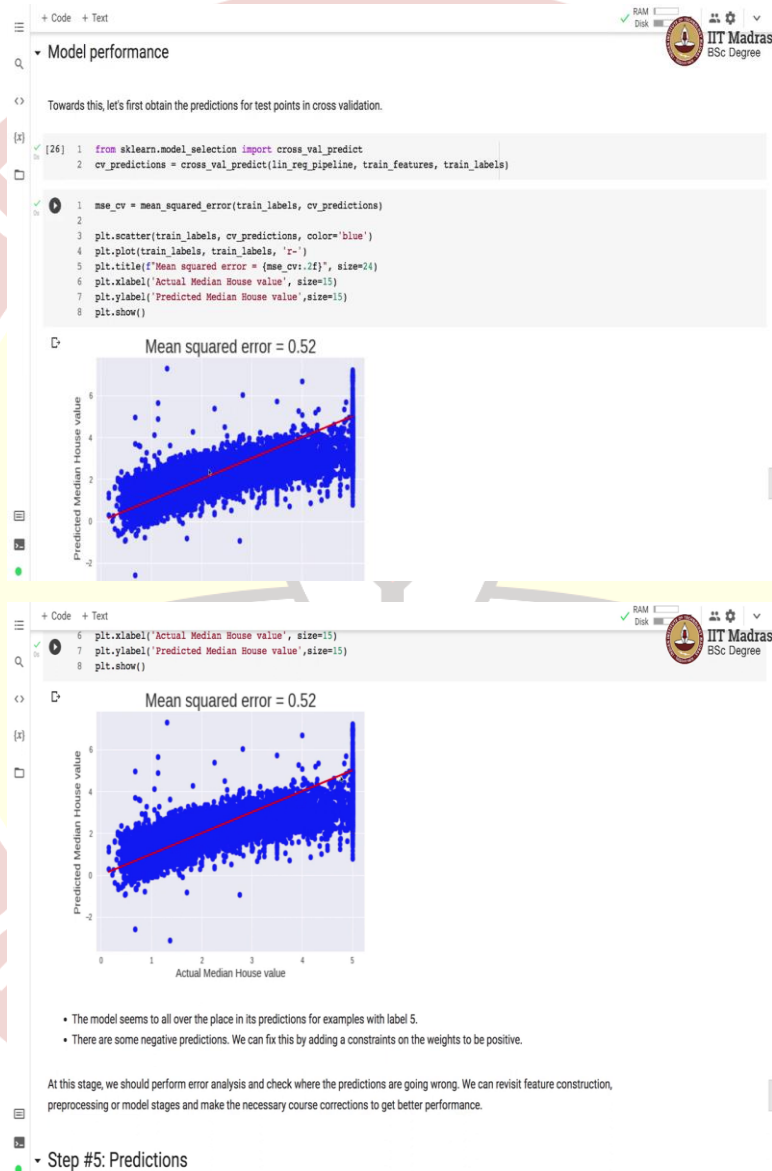
(Refer Slide Time: 25:09)



So, we have selected, we have actually obtained multiple models in cross validation, and we need to select one of the one model which is which gives us the best performance on the

validation set. So, we can we can find out the best model by calculating the one that results in the least error, we can found that by calculate, by calculating the np.argmin on the test error, we get the index of the best model and we can access the best model but from the from the, from the dictionary that was obtained in cross validate. We can examine the model coefficients of the selected model in the same way as we examined the coefficients earlier in the colab.

(Refer Slide Time: 26:06)



Finally, we also want to find out how the predictions on the test points look like compared to their actual values. So, we use cross val predict, we use cross val predict API in order to get the prediction on the test points. This particular API gives us the score for points when they were part of the test set in cross validation. We calculate the mean squared error based on the
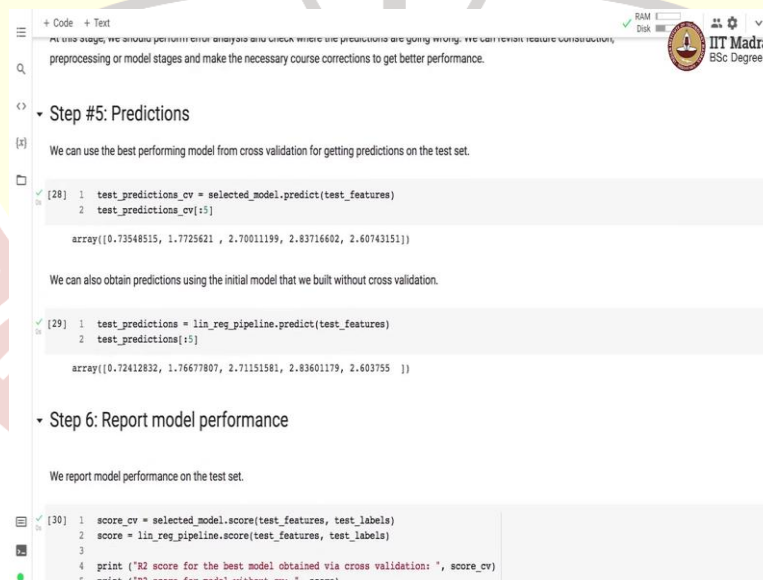
actual training labels and the labels that we have obtained in cross validation. And we simply plot the graph of the actual label and the predicted label.

Ideally, if we have perfect classification, all the labels should be along this 45-degree line, but we can see that most of the points are scattered around this particular line. And that shows the amount of error that we have in prediction.

So, if you carefully examine this graph, you will find that the model seems to be all over the place for this particular value which is label 5, for this value 5 we are predicting all possible values and there are some negative values that did not make much sense because house price cannot be negative. We can fix this by adding the constraint and the weights to be positive.

So, at this stage normally in the data in your machine learning project, we perform error analysis and check whether predictions are going wrong. Then we can revisit the model construction through feature by offering additional features you know, applying different preprocessing features, or changing the model. So, in that way, we make a necessary course correction to get better performance in the next run. So, you are where the iterative process of model building comes in.

(Refer Slide Time: 28:38)



And once you finally build a model, we obtain the predictions on the test set. And remember, this test set was set aside we did not use it throughout the training process. And now we are using the test process, now we are using the test set, now that the training is over in order to estimate how the model will perform on the future examples.

So, we call the predict function on the selected model by supplying the test feature matrix to it we get the test predictions. If the test labels are available to you, you can also perform, we can also plot such a plot. Normally in real life, the test labels may not be available when you are using the model in production. Hence, we have not added it over here.

(Refer Slide Time: 29:31)



Finally, we report the model performance on the test set. And for that we use the scoring method. We can also use any other metric of our interest like mean squared error or mean absolute error. And all the scoring metrics take the actual labels and predicted labels and provide us the metric

Now here we obtained two metric one with, the one with cv and another one without cv. You can observe that the cross validation base model has slightly lower mean squared error than the

model without cross validation. And hence this model performs we can conclude that the cross validator model performs better than the one without cross validation.

So, as an exercise, I would request you to change the scoring schemes and compare the results across different metrics. Another exercise is to repeat all the modeling step with SGDRegressor API that implements that that trains the linear regression model through iterative optimization. So, this was the demonstration of linear regression model with the normal equation.

We have tried to use different modeling steps in order to give you the flavor of different scale in APIs that are available at our disposal. In the next video, we will look at how to build a baseline model. Till then, Namaste and goodbye.