

## Hands-on Practice Lab: Model Evaluation and Refinement

Estimated time needed: 45 minutes

In this lab, you will use the skills acquired throughout the module, and try to refine your model's performance in predicting the price of a laptop, given the attribute values.

### Objectives

After completing this lab you will be able to:

- Use training, testing and cross validation to improve the performance of the dataset.
- Identify the point of overfitting of a model
- Use Ridge Regression to identify the change in performance of a model based on its hyperparameters
- Use Grid Search to identify the best performing model using different hyperparameters

### Setup

For this lab, we will be using the following libraries:

- `skillsnetwork` for downloading the dataset
- `pandas` for managing the data.
- `numpy` for mathematical operations.
- `sklearn` for machine learning and machine-learning-pipeline related functions.
- `seaborn` for visualizing the data.
- `matplotlib` for additional plotting tools.

### Installing Required Libraries

The following required libraries are pre-installed in the Skills Network Labs environment. However, if you run this notebook commands in a different Jupyter environment (e.g. Watson Studio or Ananconda), you will need to install these libraries by removing the `#` sign before `%pip` in the code cell below.

The following required libraries are **not** pre-installed in the Skills Network Labs environment. **You will need to run the following cell** to install them:

```
In [1]: import piplite
await piplite.install('seaborn')
```

### Importing Required Libraries

We recommend you import all required libraries in one place (here):

```
In [ ]: from tqdm import tqdm
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
from sklearn.model_selection import train_test_split, cross_val_score, GridSearchCV
from sklearn.linear_model import LinearRegression, Ridge
from sklearn.preprocessing import PolynomialFeatures
```

### Importing the Dataset

Run the cell below to download the dataset into the console.

```
In [4]: from pyodide.http import pyfetch

async def download(url, filename):
    response = await pyfetch(url)
    if response.status == 200:
        with open(filename, "wb") as f:
            f.write(await response.bytes())

In [5]: filepath = 'https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMDeveloperSkillsNetwork-DA0101EN-Coursera/laptop_pricing_dataset_mod2.csv'

In [6]: await download(filepath, "laptops.csv")
file_name="laptops.csv"

In [7]: df = pd.read_csv(file_name, header=0)
```

Note: This version of the lab is working on JupyterLite, which requires the dataset to be downloaded to the interface. While working on the downloaded version of this notebook on their local machines (Jupyter Anaconda), the learners can simply **skip the steps above**, and simply use the URL directly in the `pandas.read_csv()` function. You can uncomment and run the statements in the cell below.

Import the data set into a data frame.

```
In [8]: #filepath = 'https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMDeveloperSkillsNetwork-DA0101EN-Coursera/Laptop_pricing_dataset_mod2.csv'
#df = pd.read_csv(filepath, header=None)
```

Print the value of `df.head()`.

```
In [9]: df.head()
```

```
Out[9]:
```

	Unnamed: 0.1	Unnamed: 0	Manufacturer	Category	GPU	OS	CPU_core	Screen_Size_inch	CPU_frequency	RAM_GB	Storage_GB_SSD	Weight_pounds	Price	Price-binned	Screen-Full_HD	IPS
0	0	0	Acer	4	2	1	5	14.0	0.551724	8	256	3.52800	978	Low	0	
1	1	1	Dell	3	1	1	3	15.6	0.689655	4	256	4.85100	634	Low	1	
2	2	2	Dell	3	1	1	7	15.6	0.931034	8	256	4.85100	946	Low	1	
3	3	3	Dell	4	2	1	5	13.3	0.551724	8	128	2.69010	1244	Low	0	
4	4	4	HP	4	2	1	7	15.6	0.620690	8	256	4.21155	837	Low	1	

Drop the two unnecessary columns that have been added into the file, 'Unnamed: 0' and 'Unnamed: 0.1'. Use `drop` to delete these columns.

```
In [10]: df.drop(['Unnamed: 0', 'Unnamed: 0.1'], axis=1, inplace=True)
```

## Task 1 : Using Cross validation to improve the model

Divide the dataset into `x_data` and `y_data` parameters. Here `y_data` is the "Price" attribute, and `x_data` has all other attributes in the data set.

```
In [11]: # Write your code below and press Shift+Enter to execute
y_data = df['Price']
x_data = df.drop('Price', axis=1)
```

[Click here for the solution](#)

Split the data set into training and testing subests such that you reserve 10% of the data set for testing purposes.

```
In [12]: # Write your code below and press Shift+Enter to execute
x_train, x_test, y_train, y_test = train_test_split(x_data, y_data, test_size=0.10, random_state=1)
print("number of test samples :", x_test.shape[0])
print("number of training samples:", x_train.shape[0])
```

```
number of test samples : 24
number of training samples: 214
```

[Click here for the solution](#)

Create a single variable linear regression model using "CPU\_frequency" parameter. Print the  $R^2$  value of this model for the training and testing subsets.

```
In [13]: # Write your code below and press Shift+Enter to execute
lr = LinearRegression()
lr.fit(x_train[['CPU_frequency']], y_train)
print(lr.score(x_test[['CPU_frequency']], y_test))
print(lr.score(x_train[['CPU_frequency']], y_train))
```

```
-0.06599437350393766
0.14829792099817962
```

[Click here for the solution](#)

Run a 4-fold cross validation on the model and print the mean value of  $R^2$  score along with its standard deviation.

```
In [16]: # Write your code below and press Shift+Enter to execute
lre = LinearRegression()
scores = cross_val_score(lre, x_data[['CPU_frequency']], y_data, cv=4)
print(scores.mean())
```

```
-0.1610923238859522
```

[Click here for the solution](#)

## Task 2: Overfitting

Split the data set into training and testing components again, this time reserving 50% of the data set for testing.

```
In [17]: # Write your code below and press Shift+Enter to execute
x_train, x_test, y_train, y_test = train_test_split(x_data, y_data, test_size=0.5, random_state=0)
print("number of test samples :", x_test.shape[0])
print("number of training samples:", x_train.shape[0])
```

```
number of test samples : 119
number of training samples: 119
```

[Click here for the solution](#)

To identify the point of overfitting the model on the parameter "CPU\_frequency", you'll need to create polynomial features using the single attribute. You need to evaluate the  $R^2$  scores of the model created using different degrees of polynomial features, ranging from 1 to 5. Save this set of values of  $R^2$  score as a list.

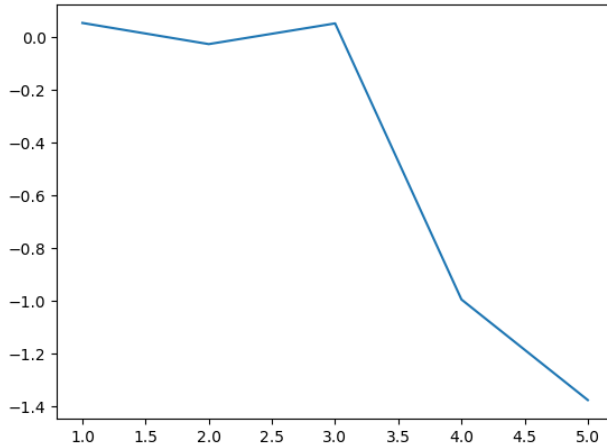
```
In [19]: # Write your code below and press Shift+Enter to execute
poly_degree = np.arange(1,6,1)
lrm = LinearRegression()
rs_test = []
for n in poly_degree:
    pf = PolynomialFeatures(degree=n)
    x_train_pr = pf.fit_transform(x_train[['CPU_frequency']])
    x_test_pr = pf.fit_transform(x_test[['CPU_frequency']])
    lrm.fit(x_train_pr,y_train)
    rs_test.append(lrm.score(x_test_pr,y_test))
```

► [Click here for the solution](#)

Plot the values of  $R^2$  scores against the order. Note the point where the score drops.

```
In [20]: # Write your code below and press Shift+Enter to execute
plt.plot(poly_degree,rs_test)
```

```
Out[20]: [ <matplotlib.lines.Line2D at 0xb6cc398>]
```



► [Click here for the solution](#)

## Task 3 : Ridge Regression

Now consider that you have multiple features, i.e. 'CPU\_frequency', 'RAM\_GB', 'Storage\_GB\_SSD', 'CPU\_core','OS','GPU' and 'Category'. Create a polynomial feature model that uses all these parameters with degree=2. Also create the training and testing attribute sets.

```
In [21]: # Write your code below and press Shift+Enter to execute
features = [ 'CPU_frequency', 'RAM_GB', 'Storage_GB_SSD', 'CPU_core', 'OS', 'GPU' , 'Category' ]
pf = PolynomialFeatures(degree=2)
x_train_pr = pf.fit_transform(x_train[features])
x_test_pr = pf.fit_transform(x_test[features])
```

► [Click here for the solution](#)

Create a Ridge Regression model and evaluate it using values of the hyperparameter alpha ranging from 0.001 to 1 with increments of 0.001. Create a list of all Ridge Regression  $R^2$  scores for training and testing data.

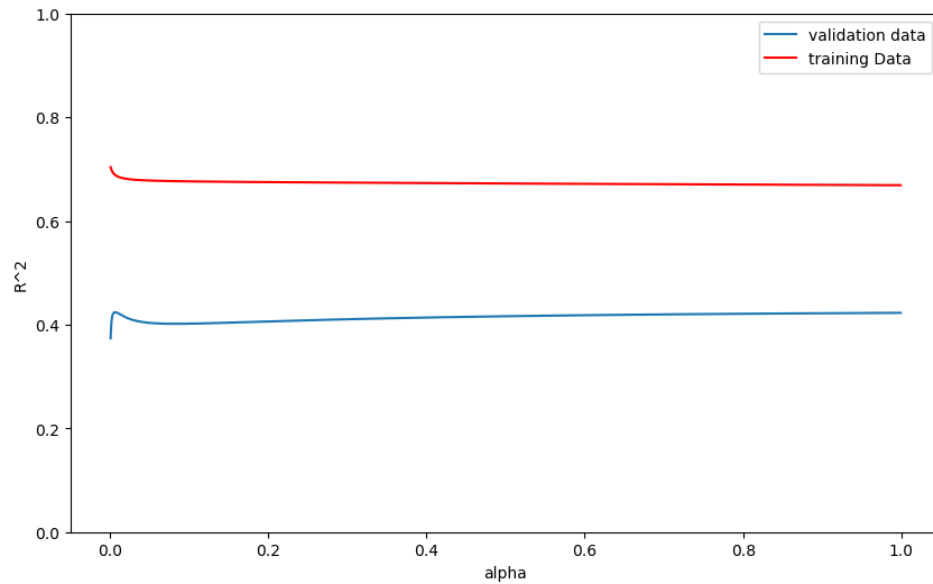
```
In [24]: # Write your code below and press Shift+Enter to execute
Alpha = np.arange(0.001,1,0.001)
rs_train = []
rs_test = []
for alpha in Alpha:
    rm = Ridge(alpha=alpha_)
    rm.fit(x_train_pr,y_train)
    rs_train.append(rm.score(x_train_pr,y_train))
    rs_test.append(rm.score(x_test_pr,y_test))
```

► [Click here for the solution](#)

Plot the  $R^2$  values for training and testing sets with respect to the value of alpha

```
In [25]: # Write your code below and press Shift+Enter to execute
plt.figure(figsize=(10, 6))
plt.plot(Alpha, rs_test, label='validation data')
plt.plot(Alpha, rs_train, 'r', label='training Data')
plt.xlabel('alpha')
plt.ylabel('R^2')
plt.ylim(0, 1)
plt.legend()
```

```
Out[25]: <matplotlib.legend.Legend at 0xb6b3f70>
```



► [Click here for the solution](#)

## Task 4: Grid Search

Using the raw data and the same set of features as used above, use GridSearchCV to identify the value of alpha for which the model performs best. Assume the set of alpha values to be used as

```
math
{0.0001, 0.001, 0.01, 0.1, 1, 10}
```

In [26]: *# Write your code below and press Shift+Enter to execute*  
`parameters = [{"alpha": [0.0001, 0.001, 0.01, 0.1, 1, 10]}`

► [Click here for the solution](#)

Create a Ridge instance and run Grid Search using a 4 fold cross validation.

In [27]: *# Write your code below and press Shift+Enter to execute*  
`rr = Ridge()  
grid = GridSearchCV(rr, parameters, cv=4)`

► [Click here for the solution](#)

Fit the Grid Search to the training data.

In [28]: `grid.fit(x_train[features], y_train)`

Out[28]:

```
GridSearchCV
  estimator: Ridge
    Ridge
```

► [Click here for the solution](#)

Print the R^2 score for the test data using the estimator that uses the derived optimum value of alpha.

In [29]: *# Write your code below and press Shift+Enter to execute*  
`bestrr = grid.best_estimator_  
print("Best R^2 score : ", bestrr.score(x_test[features], y_test))`

Best R^2 score : 0.3009905048691819

► [Click here for the solution](#)

## Congratulations! You have completed the lab

### Authors

[Abhishek Gagneja](#)

[Vicky Kuo](#)

Copyright © 2023 IBM Corporation. All rights reserved.