

# Victor Hugo Avila Felipe - A01794425



## Data Analysis with Python

**IBM: DA0101EN**

### LEARNING OBJECTIVES

In this course you will learn about:

- Data Acquisition
- How to Obtain Basic Insight From a Dataset
- Data Wrangling
- Exploratory Data Analysis
- Model Development
- Model Evaluation

## 1 Introduction to Data Analysis with Python

- Problem requiring data analysis
- dataset to analyze in python
- overview of packages
- import and export data
- Basic insights

Can we estimate the price of used cars?

### The problem

why data analysis? data everywhere, helps discovery of information.

Tom wants to sell his car, but wants the best price. what affects the price?

### Understand the data

There are documentation and the .csv file.

The first attribute, "symboling", corresponds to the insurance risk level of a car. Cars are initially assigned a risk factor symbol associated with their price. Then, if an automobile is more risky, this symbol is adjusted by moving it up the scale. A value of +3 indicates that the auto is risky, -3 that it is probably pretty safe. The second attribute "normalized-losses" is the relative average loss payment per insured vehicle year. This value is normalized for all autos within a particular size classification (two-door small, station wagons, sports/specialty, etc...), and represents the average loss per car per year. The values range from 65 to 256. The other attributes are easy to understand.

*Thus, the goal of this project is to predict "price" in terms of other car features.*

No.	Attribute name	attribute range	No.	Attribute name	attribute range
1	symboling	-3, -2, -1, 0, 1, 2, 3.	14	curb-weight	continuous from 1488 to 4066.
2	normalized-losses	continuous from 65 to 256.	15	engine-type	dohc, dohcvt, l, ohc, ohcf, ohcv, rotor.
3	make	audi, bmw, etc.	16	num-of-cylinders	eight, five, four, six, three, twelve, two.
4	fuel-type	diesel, gas.	17	engine-size	continuous from 61 to 326.
5	aspiration	std, turbo.	18	fuel-system	1bbl, 2bbl, 4bbl, idl, mfi, mpfi, spdi, spfi.
6	num-of-doors	four, two.	19	bore	continuous from 2.54 to 3.94.
7	body-style	hardtop, wagon, etc.	20	stroke	continuous from 2.07 to 4.17.
8	drive-wheels	4wd, fwd, rwd.	21	compression-ratio	continuous from 7 to 23.
9	engine-location	front, rear.	22	horsepower	continuous from 48 to 288.
10	wheel-base	continuous from 86.6 120.9.	23	peak-rpm	continuous from 4150 to 6600.
11	length	continuous from 141.1 to 208.1.	24	city-mpg	continuous from 13 to 49.
12	width	continuous from 60.3 to 72.3.	25	highway-mpg	continuous from 16 to 54.
13	height	continuous from 47.8 to 59.8.	26	price	continuous from 5118 to 45400.

Target (Label)

## Python packages

A Python library is a collection of functions and methods that allow you to perform lots of actions without writing any code. The libraries usually contain built-in modules providing different functionalities, which you can use directly. And there are extensive libraries, offering a broad range of facilities.

## Scientific Computing Libraries in Python

### 1. Scientific Computing Libraries



**Pandas**  
(Data structures & tools)



**NumPy**  
(Arrays & matrices)



**SciPy**  
(Integrals, solving differential equations, optimization)

## Visualization Libraries in Python



## Algorithmic Libraries in Python



### importing and exporting data

Format and file path. `read_csv()` `df` `df.head(n)` `df.tail(n)`

Add headers `df.columns = headers` `headers = ["a", "b", "c"]`

export df to csv: `df.to_csv(path)`

csv, json, excel, sql

### Analyzing data

Check data type data distribution

object, float, int y datetime

- potential info and type mismatch
- compatibility with python methods

dataframe.dtypes

df.describe(include="all") -> count, mean, std deviation, min, 25%, 50%, 75%, max all ->  
UNIQUE, top, freq

## Data lab

- Data source: <https://archive.ics.uci.edu/ml/machine-learning-databases/autos/imports-85.data> ([https://archive.ics.uci.edu/ml/machine-learning-databases/autos/imports-85.data?utm\\_medium=Exinfluencer&utm\\_source=Exinfluencer&utm\\_content=000026UJ&utm\\_term=10006555&utm\\_id=NA-SkillsNetwork-Channel-SkillsNetworkCoursesIBMDveloperSkillsNetworkDA0101ENSkillNetwork20235326-2021-0](https://archive.ics.uci.edu/ml/machine-learning-databases/autos/imports-85.data?utm_medium=Exinfluencer&utm_source=Exinfluencer&utm_content=000026UJ&utm_term=10006555&utm_id=NA-SkillsNetwork-Channel-SkillsNetworkCoursesIBMDveloperSkillsNetworkDA0101ENSkillNetwork20235326-2021-0))
- Data type: csv

```
In [10]: #install specific version of libraries used in Lab
import sys
!{sys.executable} -m pip install pandas
!{sys.executable} -m pip install numpy
!{sys.executable} -m pip install matplotlib
!{sys.executable} -m pip install scipy
!{sys.executable} -m pip install seaborn
!{sys.executable} -m pip install ipywidgets
Requirement already satisfied: pandas in /home/flynn/anaconda3/lib/python3.9/site-packages (1.4.2)
Requirement already satisfied: python-dateutil>=2.8.1 in /home/flynn/anaconda3/lib/python3.9/site-packages (from pandas) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in /home/flynn/anaconda3/lib/python3.9/site-packages (from pandas) (2021.3)
Requirement already satisfied: numpy>=1.18.5 in /home/flynn/anaconda3/lib/python3.9/site-packages (from pandas) (1.21.5)
Requirement already satisfied: six>=1.5 in /home/flynn/anaconda3/lib/python3.9/site-packages (from python-dateutil>=2.8.1->pandas) (1.16.0)
Requirement already satisfied: numpy in /home/flynn/anaconda3/lib/python3.9/site-packages (1.21.5)
Requirement already satisfied: matplotlib in /home/flynn/anaconda3/lib/python3.9/site-packages (3.5.1)
Requirement already satisfied: kiwisolver>=1.0.1 in /home/flynn/anaconda3/lib/python3.9/site-packages (from matplotlib) (1.3.2)
Requirement already satisfied: numpy>=1.17 in /home/flynn/anaconda3/lib/python3.9/site-packages (from matplotlib) (1.21.5)
Requirement already satisfied: fonttools>=4.22.0 in /home/flynn/anaconda3/lib/python3.9/site-packages (from matplotlib) (4.22.0)
```

## Read Data

We use `pandas.read_csv()` function to read the csv file. In the brackets, we put the file path

along with a quotation mark so that pandas will read the file into a dataframe from that address. The file path can be either an URL or your local file address.

Because the data does not include headers, we can add an argument `headers = None` inside the `read_csv()` method so that pandas will not automatically set the first row as a header.

You can also assign the dataset to any variable you create.

This dataset was hosted on IBM Cloud object. Click [HERE](https://cocl.us/DA101EN_object_storage?utm_medium=Exinfluencer&utm_source=Exinfluencer&utm_content=000026UJ&utm_term=10006555&utm_id=NA-SkillsNetwork-Channel-SkillsNetworkCoursesIBMDveloperSkillsNetworkDA0101ENSkillsNetwork20235326-2021-01-01) ([https://cocl.us/DA101EN\\_object\\_storage?utm\\_medium=Exinfluencer&utm\\_source=Exinfluencer&utm\\_content=000026UJ&utm\\_term=10006555&utm\\_id=NA-SkillsNetwork-Channel-SkillsNetworkCoursesIBMDveloperSkillsNetworkDA0101ENSkillsNetwork20235326-2021-01-01](https://cocl.us/DA101EN_object_storage?utm_medium=Exinfluencer&utm_source=Exinfluencer&utm_content=000026UJ&utm_term=10006555&utm_id=NA-SkillsNetwork-Channel-SkillsNetworkCoursesIBMDveloperSkillsNetworkDA0101ENSkillsNetwork20235326-2021-01-01)) for free storage.

## Import pandas library

```
import pandas as pd
```

## Read the online file by the URL provides above, and assign it to variable "df"

```
df = pd.read_csv(path, header=None)
```

After reading the dataset, we can use the `dataframe.head(n)` method to check the top n rows of the dataframe, where n is an integer. Contrary to `dataframe.head(n)`, `dataframe.tail(n)` will show you the bottom n rows of the dataframe.

## Save Dataset

Correspondingly, Pandas enables us to save the dataset to csv. By using the `dataframe.to_csv()` method, you can add the file path and name along with quotation marks in the brackets.

For example, if you would save the dataframe `df` as **automobile.csv** to your local machine, you may use the syntax below, where `index = False` means the row names will not be written.

## Read/Save Other Data Formats

Data Formate	Read	Save
csv	<code>pd.read_csv()</code>	<code>df.to_csv()</code>
json	<code>pd.read_json()</code>	<code>df.to_json()</code>
excel	<code>pd.read_excel()</code>	<code>df.to_excel()</code>

Data Formate	Read	Save
	<code>hdf</code> <code>pd.read_hdf()</code>	<code>df.to_hdf()</code>
	<code>sql</code> <code>pd.read_sql()</code>	<code>df.to_sql()</code>

### Question 1

1/1 point (graded)

What does CSV stand for?

Comma-separated values

Car sold values

Car state values

None of the above

✓ Save

**Submit** You have used 1 of 2 attempts

### Question 2

0/1 point (graded)

In the data set, which of the following represents an attribute or feature?

Row

Column

Each element in the dataset

✗

**Submit** You have used 2 of 2 attempts

### Question 3

1/1 point (graded)

What is the name of what we want to predict?

Target

Feature

Dataframe

✓ Save

**Submit** You have used 1 of 2 attempts

### Question 4

1/1 point (graded)

What is the command to display the first five rows of a dataframe `df`?

`df.head()`

`df.tail()`

✓

Submit You have used 1 of 1 attempt

---

**Question 5**

1/1 point (graded)

What command do you use to get the data type of each row of the dataframe `df`?

df.dtypes

df.head()

df.tail()

✓

Save

Submit You have used 1 of 2 attempts

---

**Question 6**

1/1 point (graded)

How do you get a statistical summary of a dataframe `df`?

df.describe()

df.head()

df.tail()

✓

Save

Submit You have used 1 of 2 attempts

---

**Question 7**

1/1 point (graded)

If you use the method `describe()` without changing any of the arguments, you will get a statistical summary of all the columns of type "object".

False

True

✓

Submit You have used 1 of 1 attempt

## Module 2 - Cleaning and Preparing the Data

Pre-Processing Data in Python.- convert raw to for another data.

- identify and handling missing value
- data formatting
- Data normalization
- Data Binning

- Turning categorical values to numerical values

## Dealing with a missing value

? "N/A", 0

- check with the data collection source
- Drop the missing values: drop variable or data entry
- replace the missing values: with average, by frequency, based on other functions.

dataframes.dropna() axis 0 the entire row, 1 drops the entire column, inplace true writes the result back

### Replace missing values

The diagram illustrates the use of the `df.replace(missing_value, new_value)` method. It shows two data frames. The first data frame has a 'normalized-losses' column containing values 164, 164, NaN, 158, and ... . The second data frame shows the result after replacement, where the NaN value has been replaced by the mean (162).

```

How to replace missing values in Python

Use dataframe.replace(missing_value, new_value):



| normalized-losses | make |
|-------------------|------|
| ...               | ...  |
| 164               | audi |
| 164               | audi |
| NaN               | audi |
| 158               | audi |
| ...               | ...  |


→


| normalized-losses | make |
|-------------------|------|
| ...               | ...  |
| 164               | audi |
| 164               | audi |
| 162               | audi |
| 158               | audi |
| ...               | ...  |



mean = df["normalized-losses"].mean()
df["normalized-losses"].replace(np.nan, mean)

COGNITIVE CLASS.ai

```

## How to deal with missing data?

### Check with the data collection source

### Drop the missing values

- drop the variable
- drop the data entry

### Replace the missing values

- replace it with an average (of similar datapoints)
- replace it by frequency
- replace it based on other functions

### Leave it as missing data



## Data Formatting

- Data are usually collected from different places and stored in different formats.
- Bringing data into a common standard of expression allows users to make meaningful comparison.



## Applying calculations to an entire column

- Convert "mpg" to "L/100km" in Car dataset.

The diagram shows a transformation from one column to another. On the left, a table is labeled 'city-mpg' and contains the following data: 21, 21, 19, ... An arrow points to the right, leading to a table labeled 'city-L/100km' which contains: 11.2, 11.2, 12.4, ...

```
df["city-mpg"] = 235/df["city-mpg"]
```

```
df.rename(columns={"city_mpg": "city-L/100km"}, inplace=True)
```

## Incorrect data types

- Sometimes the wrong data type is assigned to a feature.

```
df["price"].tail(5)
```

200	16845
201	19045
202	21485
203	22470
204	22625

Name: price, dtype: object

## Correcting data types

To identify data types:

- Use `dataframe.dtypes()` to identify data type.

To convert data types:

- Use `dataframe.astype()` to convert data type.

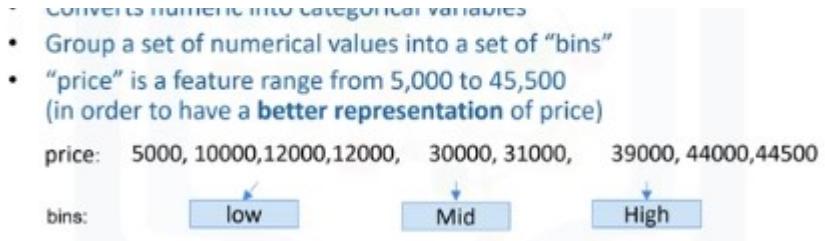
Example: convert data type to integer in column "price"

```
df["price"] = df["price"].astype("int")
```

Binning

Ver más ta... Compartir

- Binning: Grouping of values into "bins"
- Converts numeric into categorical variables



## Binning in Python pandas

A diagram showing a transformation process. On the left, there is a table with a single column labeled 'price' containing several numerical values. An arrow points from this table to another table on the right. The right table has two columns: 'price' and 'price-binned'. The 'price-binned' column uses color coding to represent the bins: 'Low' (orange), 'Medium' (light orange), and 'High' (yellow). The original price values are mapped into these categories based on the defined bins.

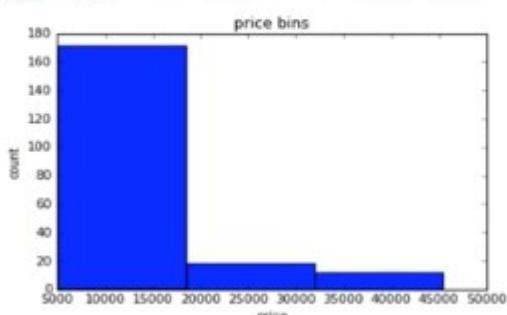
price
13495
16500
18920
41315
5151
6295
...

price	price-binned
13495	Low
16500	Low
18920	Medium
41315	High
5151	Low
6295	Low
...	...

```
bins = np.linspace(min(df["price"]), max(df["price"]), 4)
group_names = ["Low", "Medium", "High"]
df["price-binned"] = pd.cut(df["price"], bins, labels=group_names, include_lowest=True )
```

## Visualizing binned data

- E.g., Histograms



## Categorical → Numeric

### Solution:

- Add dummy variables for each unique category
- Assign 0 or 1 in each category

A diagram showing a transformation of categorical data into binary features. On the left, there is a table with columns 'Car', 'Fuel', and 'gas'. The 'Fuel' column has values 'gas' and 'diesel'. The 'gas' column has values '1' and '0'. An arrow points from this table to another table on the right. This second table has columns 'Car', 'Fuel\_gas', 'Fuel\_diesel', 'gas\_1', and 'gas\_0'. The 'Fuel\_gas' and 'Fuel\_diesel' columns are binary indicators. The 'gas\_1' and 'gas\_0' columns are binary indicators. The data shows that for Car A, Fuel is gas and gas\_1 is 1; for Car B, Fuel is diesel and gas\_0 is 1; for Car C, Fuel is gas and gas\_1 is 1; and for Car D, Fuel is gas and gas\_0 is 1.

Car	Fuel	...	gas	diesel
A	gas	...	1	0
B	diesel	...	0	1
C	gas	...	1	0
D	gas	...	1	0

"One-hot encoding"

## Dummy variables in Python pandas

- Use `pandas.get_dummies()` method.
- Convert categorical variables to dummy variables (0 or 1)

fuel
gas
diesel
gas
gas

gas	diesel
1	0
0	1
1	0
1	0

```
pd.get_dummies(df['fuel'])
```

## LAB

### Question #1:

Based on the example above, replace NaN in "stroke" column with the mean value.

```
avg_stroke = df["stroke"].astype("float").mean(axis = 0) print("Average of stroke:", avg_stroke)  
df["stroke"].replace(np.nan, avg_stroke, inplace = True)
```

Calculate the mean value for the "horsepower" column

Question #2:

According to the example above, transform mpg to L/100km in the column of "highway-mpg" and change the name of column to "highway-L/100km".

```
df["highway-mpg"] = 235/df["highway-mpg"]  
df.rename(columns={"highway-mpg":'highway-L/100km'}, inplace=True)  
df.head()
```

### Question #3:

According to the example above, normalize the column "height".

```
df['height'] = df['height']/df['height'].max()
```

```
df[["length","width","height"]].head()
```

## Question #4:

**Similar to before, create an indicator variable for the column "aspiration"**

```
dummy_variable_2 = pd.get_dummies(df['aspiration'])
```

```
dummy_variable_2.rename(columns={'std':'aspiration-std', 'turbo': 'aspiration-turbo'},  
inplace=True)
```

```
dummy_variable_2.head()
```

## Question #5:

**Merge the new dataframe to the original dataframe, then drop the column 'aspiration'.**

```
df = pd.concat([df, dummy_variable_2], axis=1)
```

```
df.drop('aspiration', axis = 1, inplace=True)
```

---

### Question 1

1/1 point (graded)

Consider the dataframe `df`. What is the result of the following operation: `df['symboling'] = df['symboling'] + 1` ?

Every element in the column "symboling" will increase by one.

Every element in the row "symboling" will increase by one.

Every element in the dataframe will increase by one.



[Save](#) | [Show answer](#)

[Submit](#)

You have used 1 of 2 attempts

---

✓ Correct (1/1 point)

---

### Question 2

1/1 point (graded)

Consider the dataframe `df`. What does the command `df.rename(columns={'a':'b'})` change about the dataframe `df` ?

Renames column "a" of the dataframe to "b".

Renames row "a" to "b".

Nothing. You must set the parameter "inplace = True".



[Save](#) | [Show answer](#)

### Question 3

1/1 point (graded)

Consider the dataframe "df". What is the result of the following operation `df['price'] = df['price'].astype(int)` ?

 Convert or cast the row 'price' to an integer value. Convert or cast the column 'price' to an integer value. Convert or cast the entire dataframe to an integer value.[Save](#) | [Show answer](#)[Submit](#)

You have used 1 of 2 attempts

Correct (1/1 point)

### Question 4

1/1 point (graded)

Consider the column of the dataframe `df['a']`. The column has been standardized. What is the standard deviation of the values as a result of applying the following operation: `df['a'].std()` ?

 1 0 3[Show answer](#)[Submit](#)

You have used 2 of 2 attempts

### Question 5 a)

1/1 point (graded)

Consider the column of the dataframe, `df['Fuel']`, with two values: 'gas' and 'diesel'. What will be the name of the new columns `pd.get_dummies(df['Fuel'])` ?

 1 and 0 Just 'diesel' Just 'gas' 'gas' and 'diesel'[Save](#) | [Show answer](#)[Submit](#)

You have used 1 of 2 attempts

Correct (1/1 point)

### Question 5 b)

1/1 point (graded)

What are the values of the new columns from part 5a)?

 1 and 0 Just 'diesel'

Just 'gas'

'gas' and 'diesel'

✓

[Save](#) | [Show answer](#)

**Submit** You have used 1 of 2 attempts

---

✓ Correct (1/1 point)

## Module 3 - Exploratory Data Analysis

### Exploratory Data Analysis (EDA)

- Preliminary step in data analysis to:
  - Summarize main characteristics of the data
  - Gain better understanding of the data set
  - Uncover relationships between variables
  - Extract important variables
- Question:  
"What are the characteristics that have the most impact on the car price?"

### Learning Objectives

In this module you will learn about:

- Descriptive Statistics
- GroupBy
- ANOVA
- Correlation
- Correlation - Statistics

### Descriptive Statistics

- Describe basic features of data
- Giving short summaries about the sample and measures of the data

### Descriptive Statistics- Describe()

• Summarizing statistics using pandas describe() method

### \* Summarize statistics using pandas `describe()` method

```
df.describe()
```

	Unnamed: 0	symboling	normalized-losses	wheel-base	length	width	height	curb-weight	engine-size	bore	stroke
count	201.000000	201.000000	164.000000	201.000000	201.000000	201.000000	201.000000	201.000000	201.000000	201.000000	201.000000
mean	100.000000	0.840796	122.000000	98.797015	174.200995	65.889055	53.766667	2555.666667	126.875622	3.319154	3.256766
std	58.167861	1.254802	35.442168	6.066366	12.322175	2.101471	2.447822	517.296727	41.546834	0.280130	0.316049
min	0.000000	-2.000000	65.000000	86.600000	141.100000	60.300000	47.800000	1488.000000	61.000000	2.540000	2.070000
25%	50.000000	0.000000	NaN	94.500000	166.800000	64.100000	52.000000	2169.000000	98.000000	3.150000	3.110000
50%	100.000000	1.000000	NaN	97.000000	173.200000	65.500000	54.100000	2414.000000	120.000000	3.310000	3.290000
75%	150.000000	2.000000	NaN	102.400000	183.500000	66.600000	55.500000	2926.000000	141.000000	3.580000	3.410000
max	200.000000	3.000000	256.000000	120.900000	208.100000	72.000000	59.800000	4066.000000	326.000000	3.940000	4.170000

## Descriptive Statistics - Value\_Counts()

- summarize the categorical data is by using the `value_counts()` method

```
drive_wheels_counts=df["drive-wheels"].value_counts()  
  
drive_wheels_counts.rename(columns={'drive-wheels':'value_counts' inplace=True}  
drive_wheels_counts.index.name= 'drive-wheels'
```

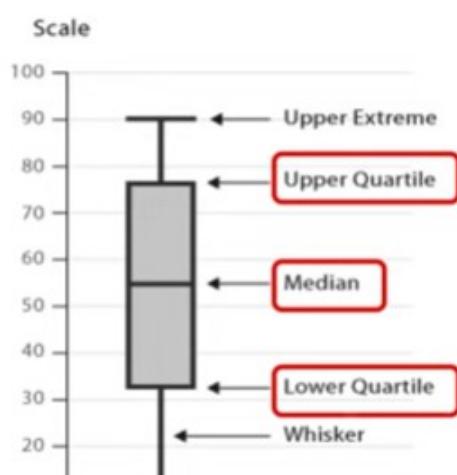
	value_counts
drive-wheels	
fwd	118
rwd	75
4wd	8



4



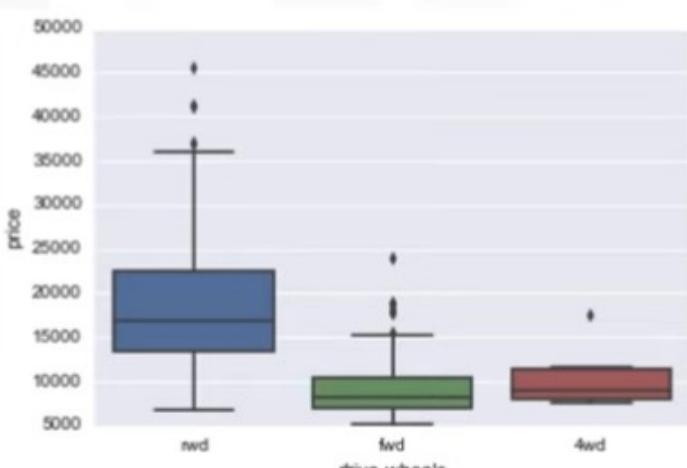
## Descriptive Statistics - Box Plots





## Box Plot - Example

```
sns.boxplot(x= "drive-wheels", y= "price", data=df)
```



**Lab3.ipynb se encuentra dentro de la carpeta.**

## Questions

### Question 1

1/1 point (graded)

Consider the dataframe "df". What method provides the summary statistics?

df.describe()

df.head()

df.tail()



[Save](#) | [Show answer](#)

[Submit](#)

You have used 1 of 2 attempts

Correct (1/1 point)

### Question 2

1/1 point (graded)

Consider the following dataframe:

```
df_test = df[['body-style', 'price']]
```

The following operation is applied:

```
df_grp = df_test.groupby(['body-style'], as_index=False).mean()
```

What are resulting values of `df_grp['price']`?

The average price for each body style.

The average price.

The average body style.



[Save](#) | [Show answer](#)

[Submit](#)

You have used 1 of 2 attempts

Correct (1/1 point)

### Question 3

1/1 point (graded)

Correlation implies causation:

False

True



[Save](#) | [Show answer](#)

[Submit](#)

You have used 1 of 2 attempts

Correct (1/1 point)

### Question 4

1/1 point (graded)

What is the minimum possible value of Pearson's Correlation?

1

-100

-1



[Save](#) | [Show answer](#)

[Submit](#)

You have used 1 of 2 attempts

Correct (1/1 point)

### Question 5

1/1 point (graded)

What is the Pearson correlation between variables X and Y if X=Y:

-1

1

0

✓

[Save](#) | [Show answer](#)

## Module 4 - Model Development

### Learning Objectives

In this module you will learn about:

1. Simple and Multiple Linear Regression
2. Model Evaluation using Visualization
3. Polynomial Regression and Pipelines
4. R-squared and MSE for In-Sample Evaluation
5. Prediction and Decision Making

- Question
  - How can you determine a fair value for a used car?

### Model Development

- A model can be thought of as a mathematical equation used to predict a value given one or more other values
- Relating one or more independent variables to dependent variables.

independent variables or features

'highway-mpg'

55 mpg



dependent variables

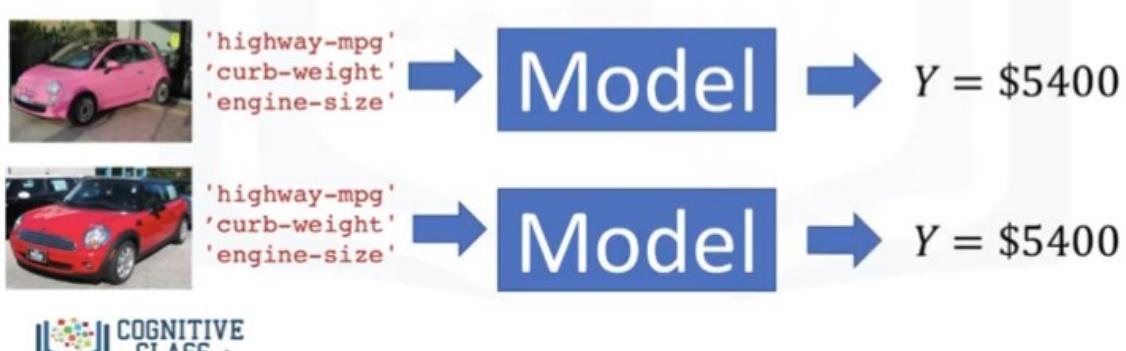
'predicted price'

\$5000

# Model Development

- To understand why more data is important consider the following situation:

- you have two almost identical cars
- Pink cars sell for significantly less



## Model Development

In addition to getting more data you can try different types of models.

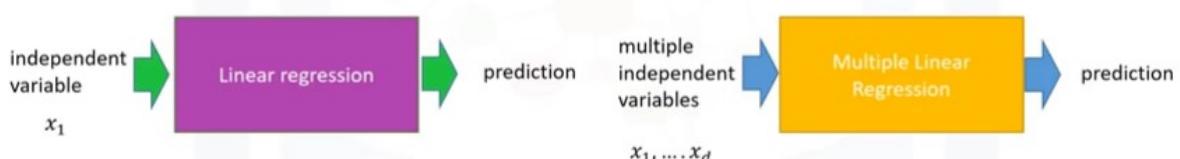
In this course you will learn about:

- Simple Linear Regression
- Multiple Linear Regression

## Linear and Multiple Linear Regression

### Introduction

- Linear regression will refer to one independent variable to make a prediction
- Multiple Linear Regression will refer to multiple independent variables to make a prediction



# Simple Linear Regression

1. The predictor (independent) variable -  $x$
2. The target (dependent) variable -  $y$

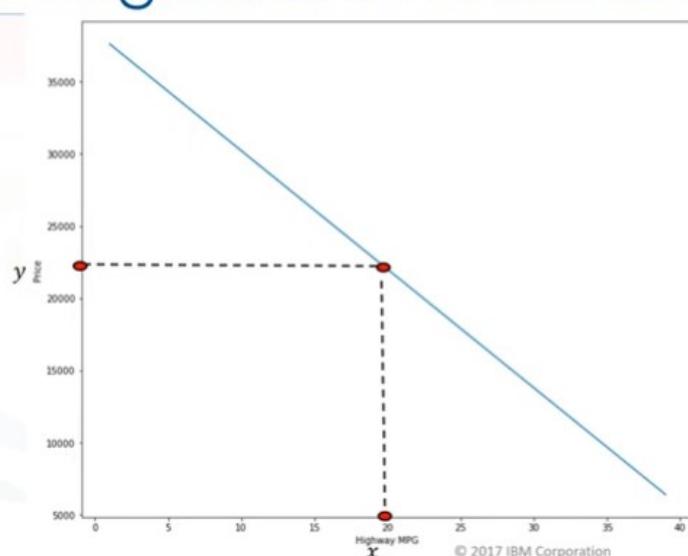
$$y = b_0 + b_1 x$$


- $b_0$ : the intercept
- $b_1$  : the slope

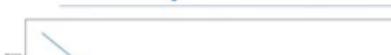


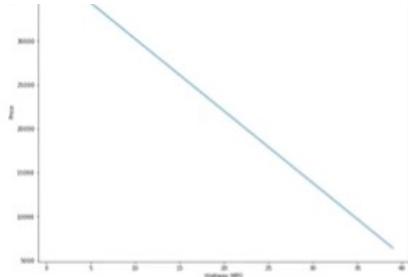
## Simple Linear Regression: Prediction

$$\begin{aligned}y &= 38423 - 821x \\&= 38423 - 821(20) \\&= 22\,003\end{aligned}$$



## Simple Linear Regression: Fit

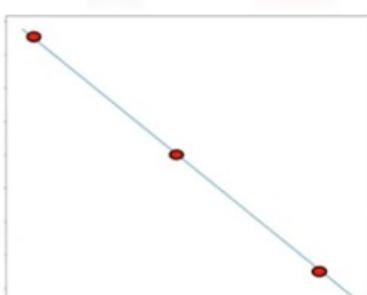




Fit

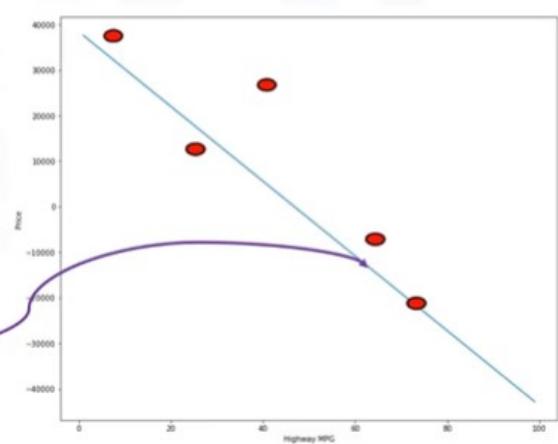
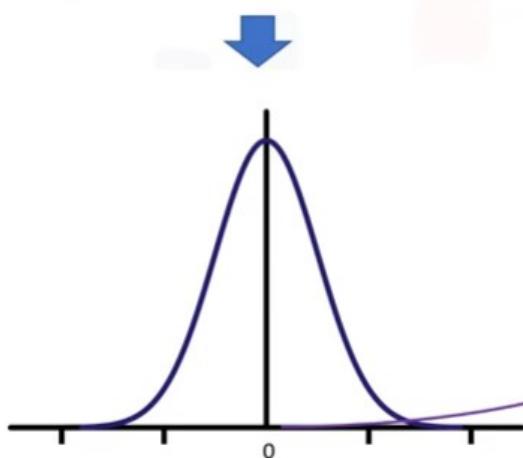
 $(b_0, b_1)$ 

## Simple Linear Regression: Fit

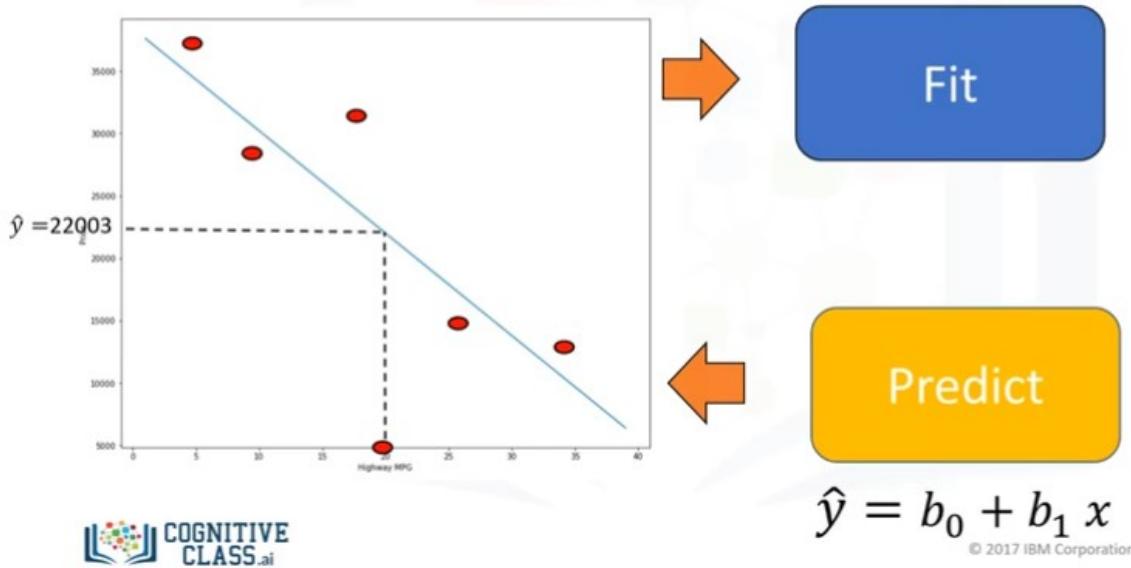


$$X = \begin{bmatrix} 0 \\ 20 \\ 40 \end{bmatrix} \quad Y = \begin{bmatrix} 38423 \\ 22003 \\ 5583 \end{bmatrix}$$

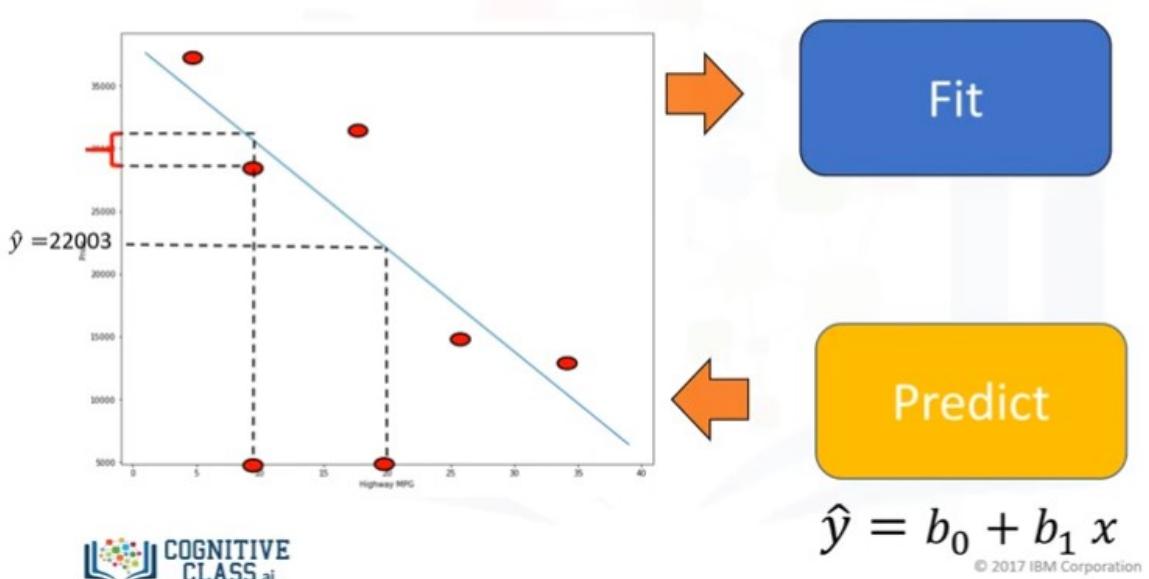
## Simple Linear Regression: Fit



## Simple Linear Regression



## Simple Linear Regression



## Fitting a Simple Linear Model Estimator

- X : Predictor variable
- Y: Target variable

1. Import linear\_model from scikit-learn

```
from sklearn.linear_model import LinearRegression
```

2. Create a Linear Regression Object using the constructor :

```
lm=LinearRegression()
```

## Fitting a Simple Linear Model

- We define the predictor variable and target variable

```
X = df[['highway-mpg']]  
Y = df['price']
```

- Then use lm.fit (X, Y) to fit the model , i.e fine the parameters  $b_0$  and  $b_1$

```
lm.fit(X, Y)
```

- We can obtain a prediction

```
Yhat=lm.predict(X)
```

Yhat	X
2	5
:	
3	4

10



## SLR – Estimated Linear Model

- We can view the intercept ( $b_0$ ):

```
lm.intercept_  
38423.305858
```

- We can also view the slope ( $b_1$ ):

```
lm.coef_  
-821.73337832
```

- The Relationship between Price and Highway MPG is given by:

Price = 38423.31 - 821.73 \* highway-mpg

$$\hat{Y} = b_0 + b_1 x$$



## Multiple Linear Regression (MLR)

This method is used to explain the relationship between:

- One continuous target (Y) variable
- Two or more predictor (X) variables

## Multiple Linear Regression (MLR)

$$\hat{Y} = b_0 + b_1x_1 + b_2x_2 + b_3x_3 + b_4x_4$$

- $b_0$  : intercept ( $X=0$ )
- $b_1$ : the coefficient or parameter of  $x_1$
- $b_2$ : the coefficient of parameter  $x_2$  and so on..

## Multiple Linear Regression (MLR)

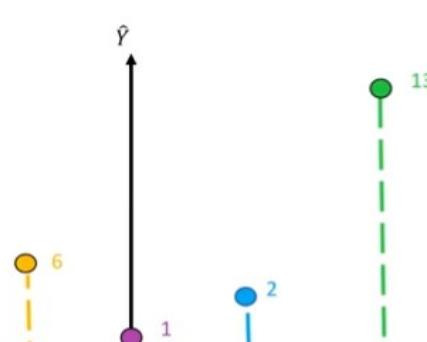
$$\hat{Y} = 1 + 2x_1 + 3x_2$$

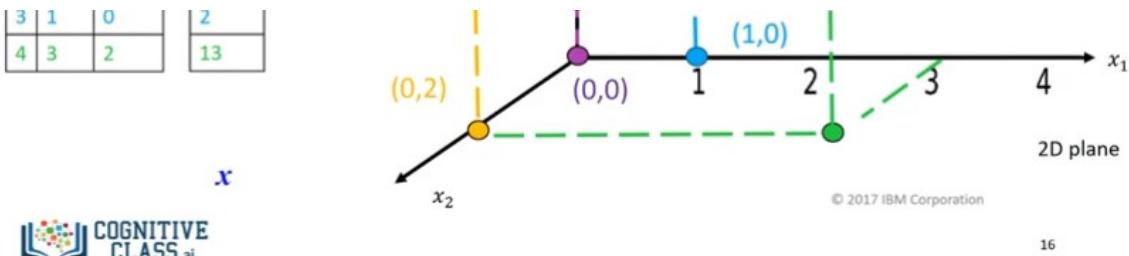
- The variables  $x_1$  and  $x_2$  can be visualized on a 2D plane, lets do an example on the next slide

- This is shown below where

$$\hat{Y} = 1 + 2x_1 + 3x_2$$

n	$x_1$	$x_2$	$\hat{Y}$
1	0	0	1
2	0	2	6





## Fitting a Multiple Linear Model Estimator

1. We can extract the for 4 predictor variables and store them in the variable Z

```
Z = df[['horsepower', 'curb-weight', 'engine-size', 'highway-mpg']]
```

2. Then train the model as before:

```
lm.fit(Z, df['price'])
```

3. We can also obtain a prediction

```
Yhat=lm.predict(X)
```

A diagram showing a matrix multiplication from  $X$  to  $\hat{Y}$ . An arrow points from the matrix  $X$  to the vector  $\hat{Y}$ .

$x_1$	$x_2$	$x_3$	$x_4$	$\hat{Y}$
3	5	-4	3	2
:	:	:	:	:
2	4	2	-4	3

COGNITIVE

## MLR – Estimated Linear Model

1. Find the intercept ( $b_0$ )

```
lm.intercept_
-15678.742628061467
```

2. Find the coefficients ( $b_1, b_2, b_3, b_4$ )

```
lm.coef_
array([52.65851272 ,4.69878948,81.95906216 , 33.58258185])
```

The Estimated Linear Model:

- Price =  $-15678.74 + (52.66) * \text{horsepower} + (4.70) * \text{curb-weight} + (81.96) * \text{engine-size} + (33.58) * \text{highway-mpg}$

$$\hat{Y} = b_0 + b_1 x_1 + b_2 x_2 + b_3 x_3 + b_4 x_4$$

## Model Evaluation using Visualization

### Regression Plot

## Why use regression plot?

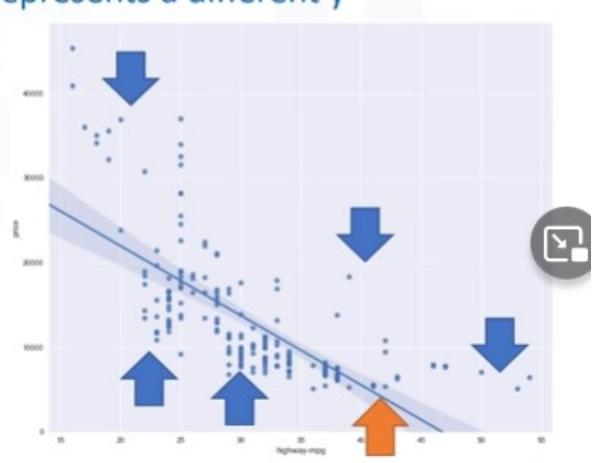
It gives us a good estimate of:

1. The relationship between two variables
2. The strength of the correlation
3. The direction of the relationship (positive or negative)

# Regression Plot

Regression Plot shows us a combination of:

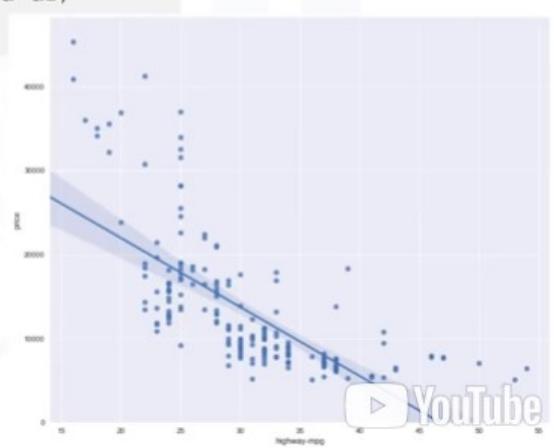
- The scatterplot: where each point represents a different y
- The fitted linear regression line ( $\hat{y}$ ).



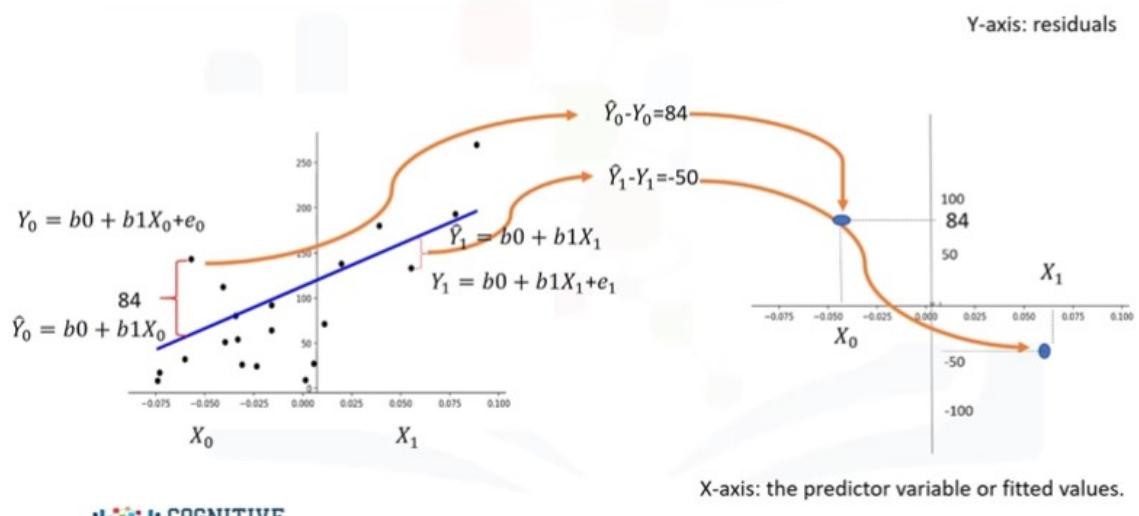
## Regression Plot

Ver más ta... Compartir

```
import seaborn as sns  
sns.regplot(x="highway-mpg", y="price", data=df)  
plt.ylim(0,)
```

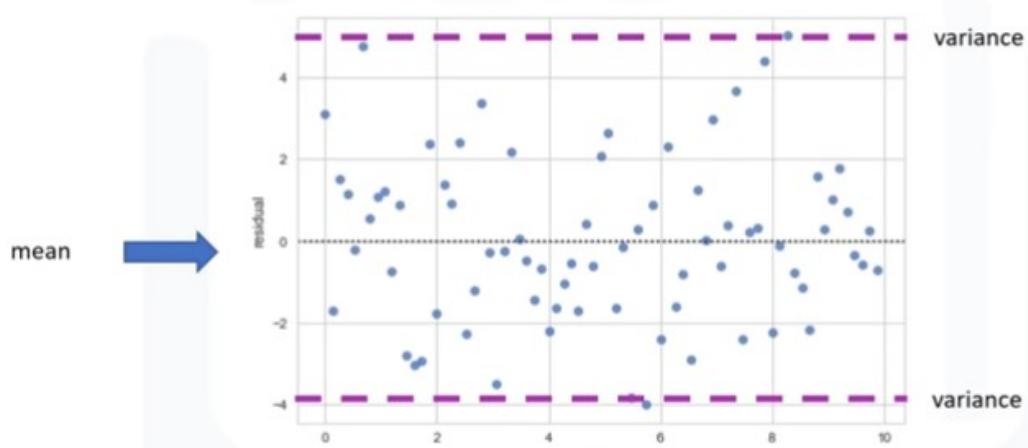


## Residual Plot



5

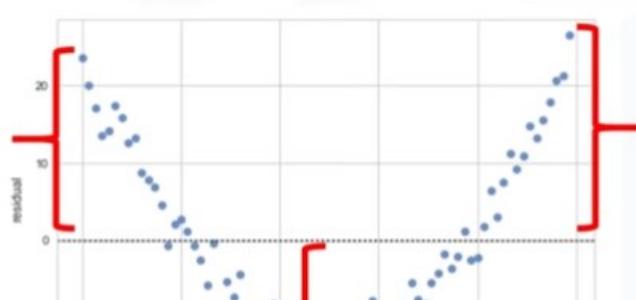
## Residual Plot

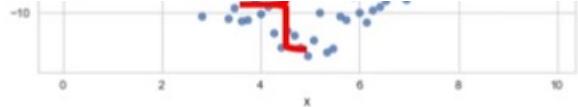


- Look at the **spread of the residuals**:
  - Randomly spread out around x-axis then a linear model is appropriate.



## Residual Plot

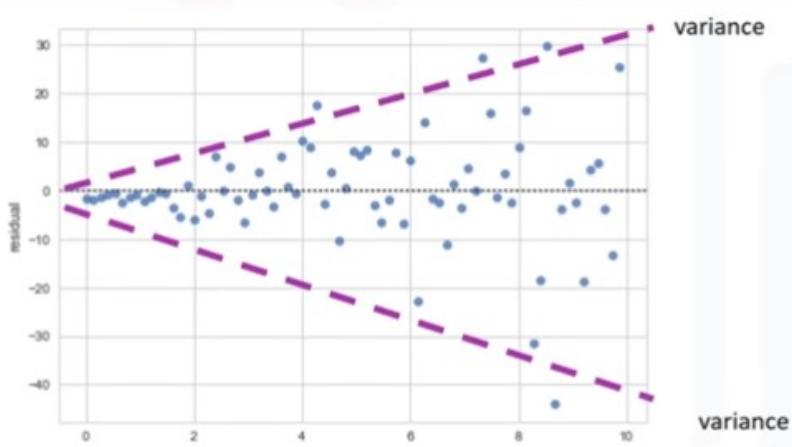




- Not randomly spread out around the x-axis
- Nonlinear model may be more appropriate



## Residual Plot

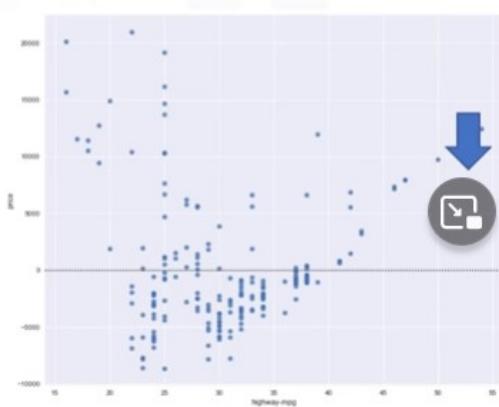


- Not randomly spread out around the x-axis
- Variance appears to change with x axis

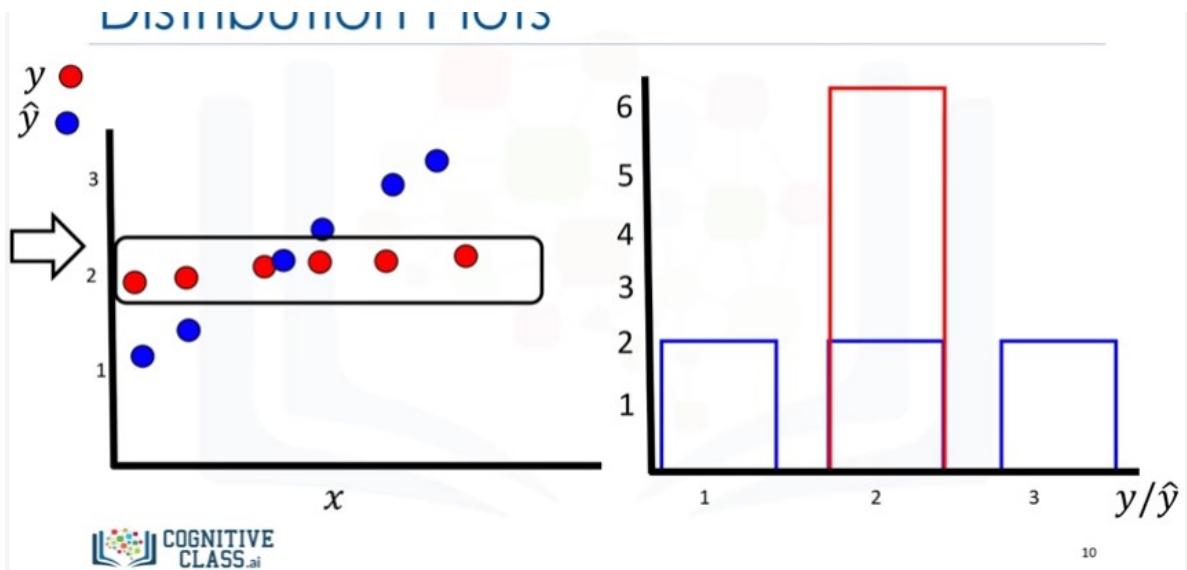


## Residual Plot

```
import seaborn as sns  
sns.residplot(df['highway-mpg'], df['price'])
```



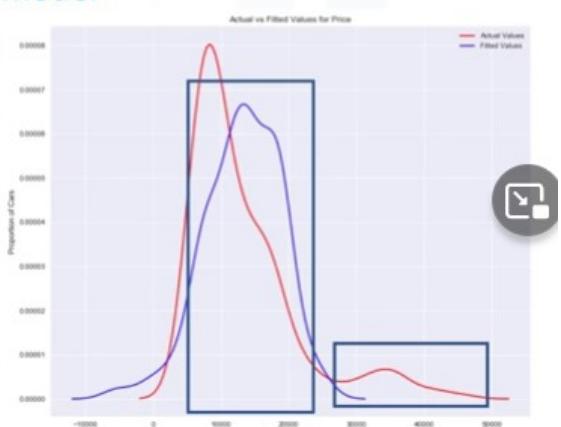
## Distribution Plots



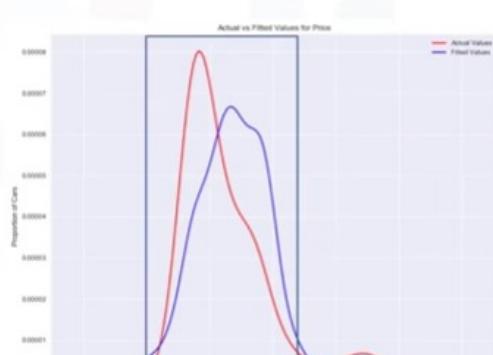
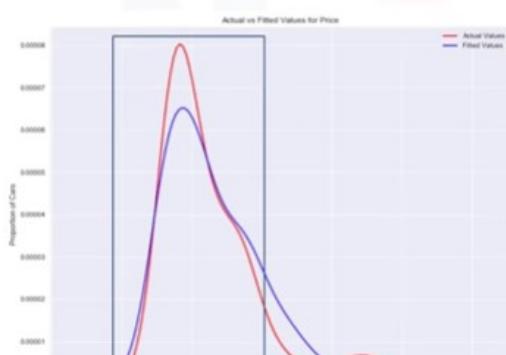
## Distribution Plots

Compare the distribution plots:

- The fitted values that result from the model
- The actual values



## MLR – Distribution Plots





13

## Distribution Plots

```
import seaborn as sns  
ax1 = sns.distplot(df['price'], hist=False, color="r", label="Actual Value")  
  
sns.distplot(Yhat, hist=False, color="b", label="Fitted Values" , ax=ax1)
```



© 2017 IBM Corporation

14

## Polynomial Regression and Pipelines

### Distribution Plots

```
import seaborn as sns  
ax1 = sns.distplot(df['price'], hist=False, color="r", label="Actual Value")  
  
sns.distplot(Yhat, hist=False, color="b", label="Fitted Values" , ax=ax1)
```



© 2017 IBM Corporation

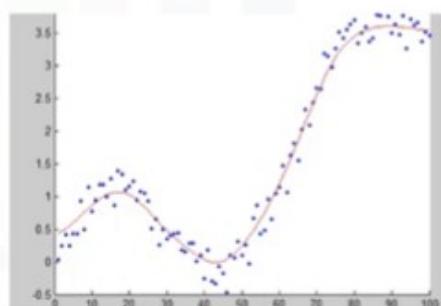
14

# Polynomial Regressions

- A special case of the general linear regression model
- Useful for describing curvilinear relationships.

## Curvilinear relationships:

By squaring or setting higher-order terms of the predictor variables.



2



# Polynomial Regression

- Quadratic – 2<sup>nd</sup> order

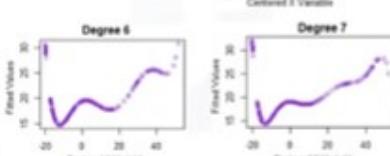
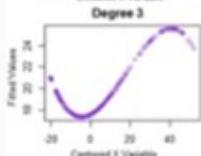
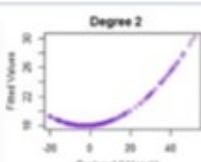
$$\hat{Y} = b_0 + b_1 x_1 + b_2 (x_1)^2$$

- Cubic – 3<sup>rd</sup> order

$$\hat{Y} = b_0 + b_1 x_1 + b_2 (x_1)^2 + b_3 (x_1)^3$$

- Higher order

$$\hat{Y} = b_0 + b_1 x_1 + b_2 (x_1)^2 + b_3 (x_1)^3 + ..$$



3

# Polynomial Regression

1. Calculate Polynomial of 3<sup>rd</sup> order

```
f=np.polyfit(x,y,3)
```

```
p=np.poly1d(f)
```

2. We can print out the model

```
print (p)
```

$$-1.557(x_1)^3 + 204.8(x_1)^2 + 8965x_1 + 1.37 \times 10^5$$

## Polynomial Regression with More than One Dimension

- We can also have multi dimensional polynomial linear regression

$$\hat{Y} = b_0 + b_1 X_1 + b_2 X_2 + b_3 X_1 X_2 + b_4 (X_1)^2 + b_5 (X_2)^2 + \dots$$

## Polynomial Regression with More than One Dimension

- The "preprocessing" library in scikit-learn,

```
from sklearn.preprocessing import PolynomialFeatures
pr=PolynomialFeatures(degree=2, include_bias=False)

x_polly=pr.fit_transform(x[['horsepower', 'curb-weight']])
```

## Polynomial Regression with More than One Dimension

```
pr=PolynomialFeatures(degree=2)
```

$X_1$	$X_2$
1	2

```
pr.fit_transform([1,2], include_bias=False)
```



$X_1$	$X_2$	$X_1 X_2$	$X_1^2$	$X_2^2$
1	2	(1) 2	1	(2) <sup>2</sup>
1	2	2	1	4

## Pre-processing

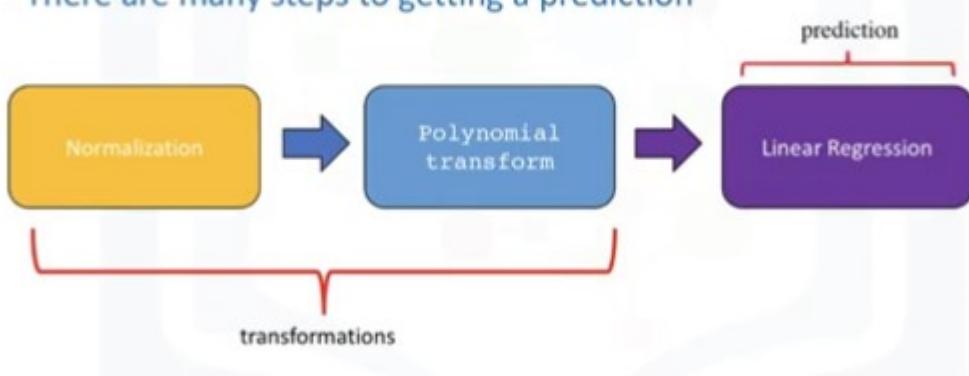
- For example we can Normalize the each feature simultaneously:

```
from sklearn.preprocessing import StandardScaler
SCALE=StandardScaler()
SCALE.fit(x_data[['horsepower', 'highway-mpg']])

x_scale=SCALE.transform(x_data[['horsepower', 'highway-mpg']])
```

## Pipelines

- There are many steps to getting a prediction



## Pipeline Constructor

```
Input=[('scale',StandardScaler()),('polynomial',PolynomialFeatures(degree=2),-.
('model',LinearRegression()))
pipe=Pipeline(Input)
```

• Pipeline constructor

- We can train the pipeline object

```
Pipe.train(X['horsepower', 'curb-weight', 'engine-size', 'highway-mpg'],y)
yhat=Pipe.predict(X[['horsepower', 'curb-weight', 'engine-size', 'highway-mpg']])
```



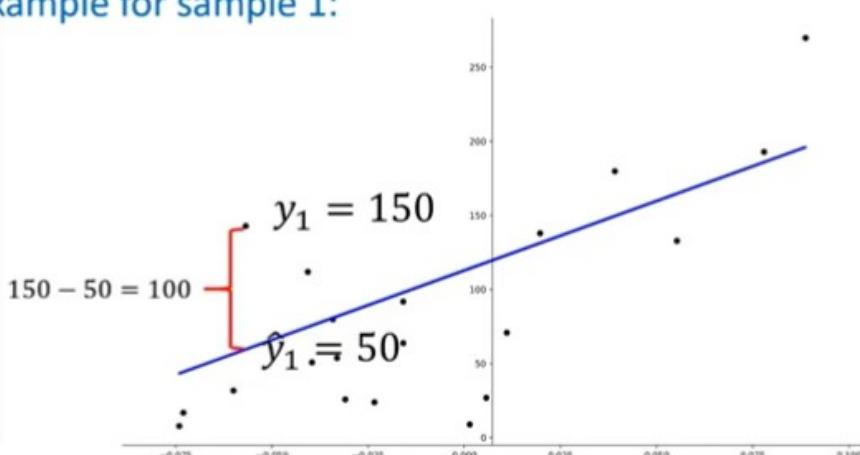
## Measures for In-Sample Evaluation

### Measures for In-Sample Evaluation

- A way to numerically determine how good the model fits on dataset.
- Two important measures to determine the fit of a model:
  - Mean Squared Error (MSE)
  - R-squared (R<sup>2</sup>)

## Mean Squared Error (MSE)

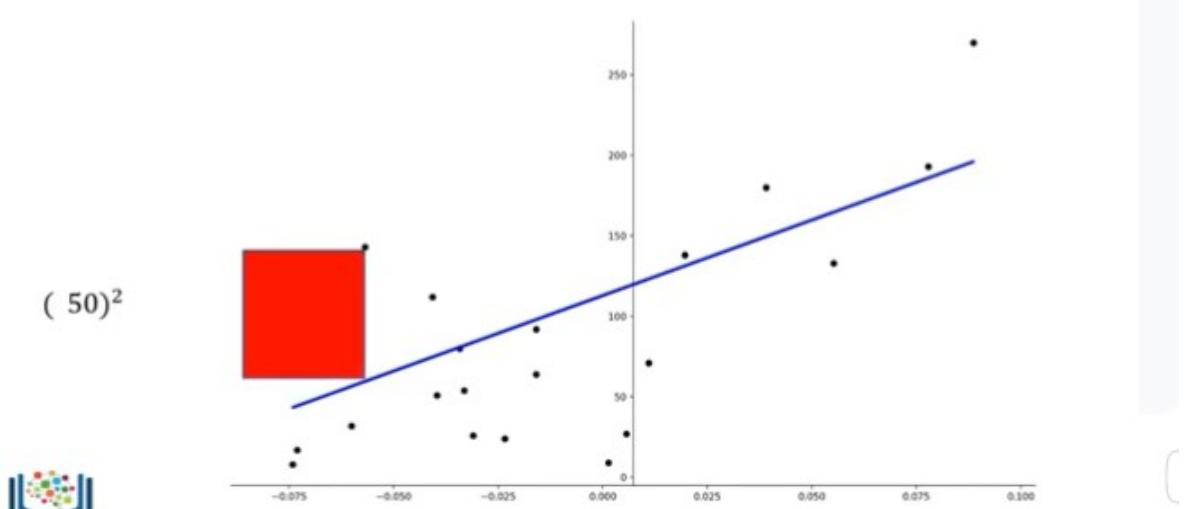
- For Example for sample 1:



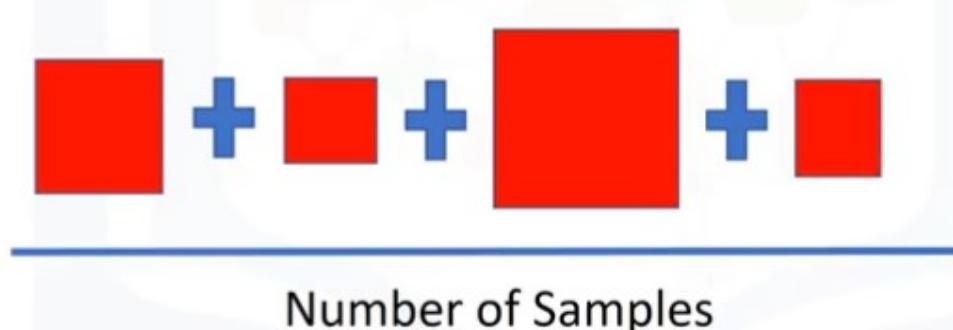
## Mean Squared Error (MSE)

Ver más t

- To make all the values positive we square it



## Mean Squared Error (MSE)



# Mean Squared Error (MSE)

- In python we can measure the MSE as follows

```
from sklearn.metrics import mean_squared_error
mean_squared_error(df['price'], Y_predict_simple_fit)
3163502.944639888
```

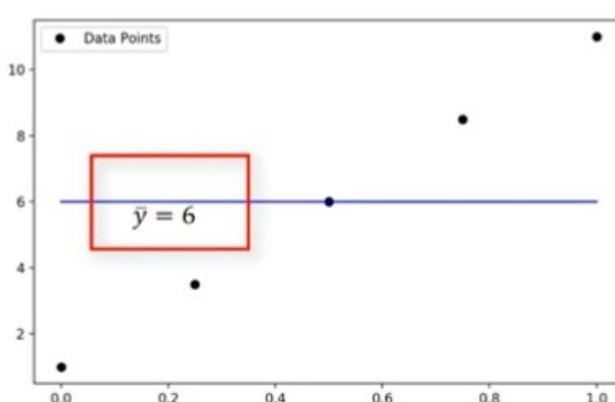
## R-squared/ R<sup>2</sup>

[Ver más tareas](#) [Compartir](#)

- The Coefficient of Determination or R squared (R<sup>2</sup>)
- Is a measure to determine how close the data is to the fitted regression line.
- R<sup>2</sup>: the percentage of variation of the target variable (Y) that is explained by the linear model.
- Think about as comparing a regression model to a simple model i.e the mean of the data points

## Coefficient of Determination (R<sup>2</sup>)

- In this example the average of the data points  $\bar{y}$  is 6

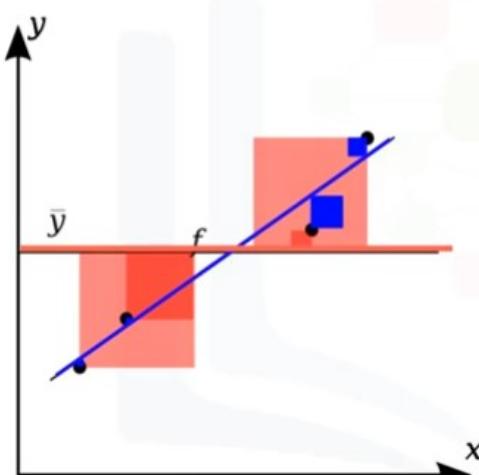


## Coefficient of Determination (R<sup>2</sup>)

## COEFFICIENT OF DETERMINATION ( $R^2$ )

$$R^2 = \left( 1 - \frac{\text{MSE of regression line}}{\text{MSE of the average of the data}} \right)$$

## Coefficient of Determination ( $R^2$ )



- The blue line represents the regression line
- The blue squares represents the MSE of the regression line
- The red line represents the average value of the data points
- The red squares represent the MSE of the red line
- We see the area of the blue squares is much smaller then the area of the red squares



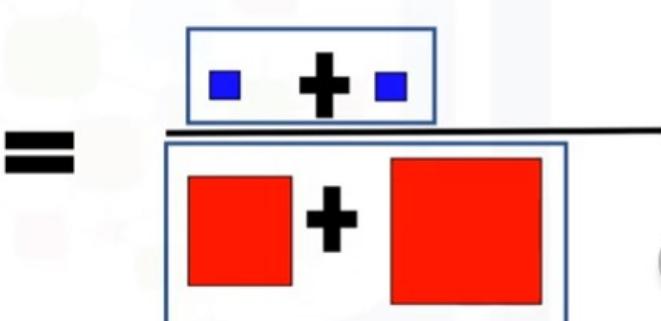
Source: [https://en.wikipedia.org/wiki/Coefficient\\_of\\_determination](https://en.wikipedia.org/wiki/Coefficient_of_determination)

10

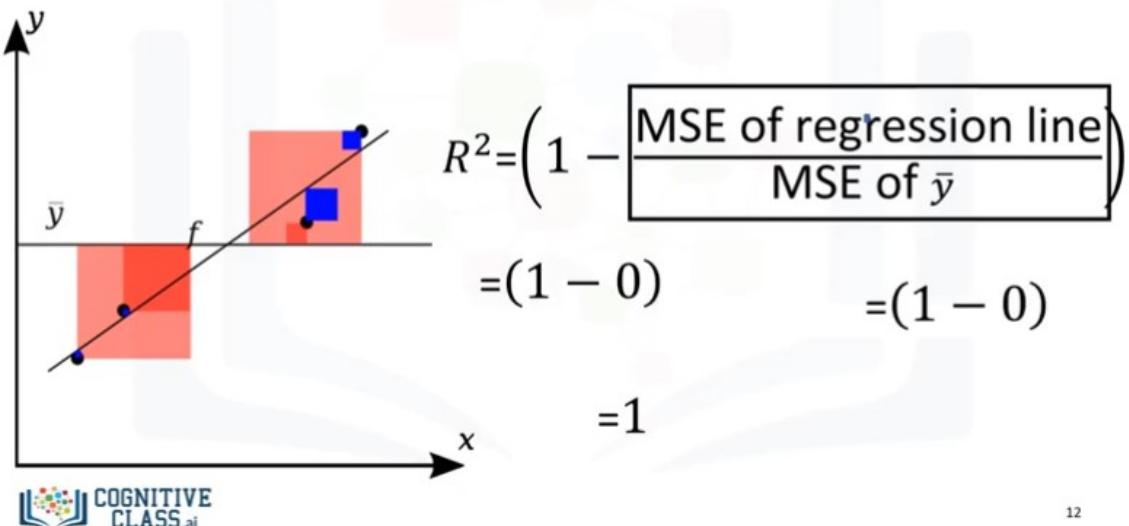
## DATOTYPEN - Measures for In Sample Evaluation (3.50) Coefficient of Determination ( $R^2$ ) Ver más la... Com

- In this case ratio of the areas of MSE is close to zero

$$\frac{\text{MSE of regression line}}{\text{MSE of } \bar{y}} = 0$$

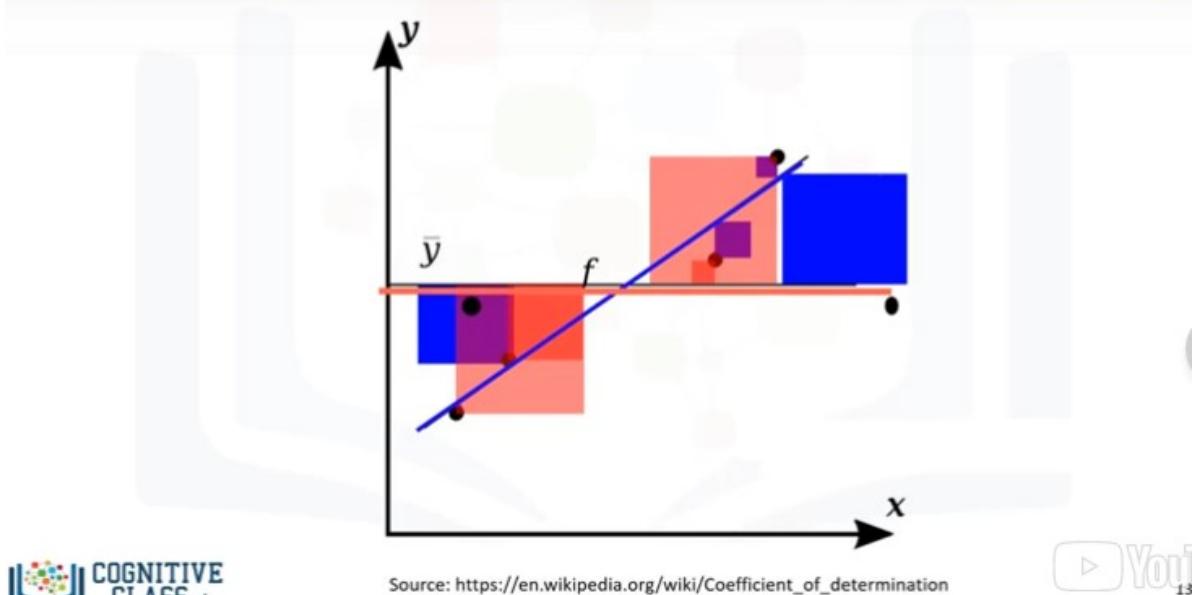


## Coefficient of Determination ( $R^2$ )



12

## Coefficient of Determination ( $R^2$ )

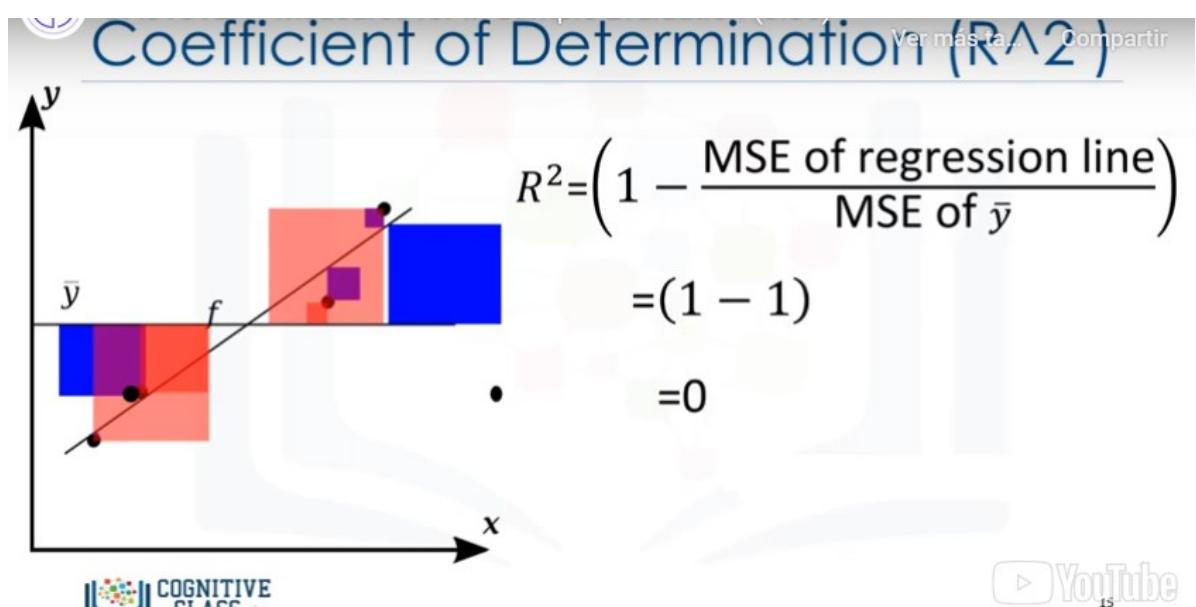


## Coefficient of Determination ( $R^2$ )

$$\frac{\text{MSE of regression line}}{\text{MSE of } \bar{y}}$$

$$= \frac{\text{Blue Box} + \text{Blue Box}}{\text{Red Box} + \text{Red Box}}$$

$$= 1$$



## R-squared/ $R^2$

- Generally the values of the MSE are between 0 and 1.
- WE can calculate the the  $R^2$  as follows

```
X = df[['highway-mpg']]
Y = df['price']
```

```
lm.fit(X, Y)
```

```
lm.score(X, y)
0.496591188
```

## Prediction and Decision Making

## Prediction and Decision Making

## Decision Making: Determining a Good Model Fit

To determine final best fit, we look at a combination of:

- Do the predicted values make sense
- Visualization
- Numerical measures for evaluation
- Comparing Models

### Do the predicted values make sense

- First we train the model

```
lm.fit(df['highway-mpg'], df['prices'])
```

- Let's predict the price of a car with 30 highway-mpg.

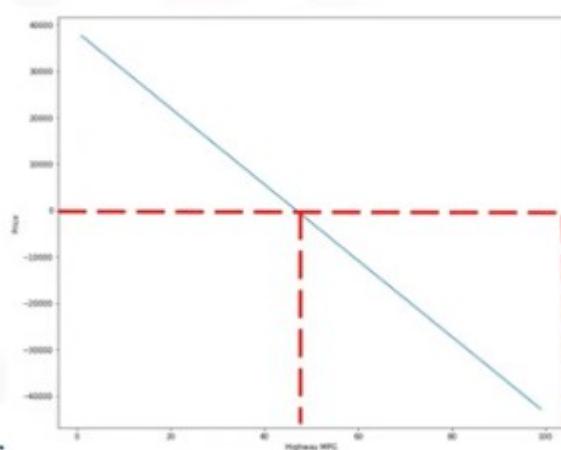
```
lm.predict(30)
```

- Result: \$ 13771.30

```
lm.coef  
-821.73337832
```

- Price = 38423.31 - 821.73 \* highway-mpg

### Do the predicted values make sense



## DO THE PREDICTED VALUES MAKE SENSE

- First we import numpy

```
import numpy as np
```

- We use the numpy function `arange` to generate a sequence from 1 to 100

```
new_input=np.arange(1,101,1).reshape(-1,1)
```



## DO THE PREDICTED VALUES MAKE SENSE

- We can predict new values

```
yhat=lm.predict(new_input)
```

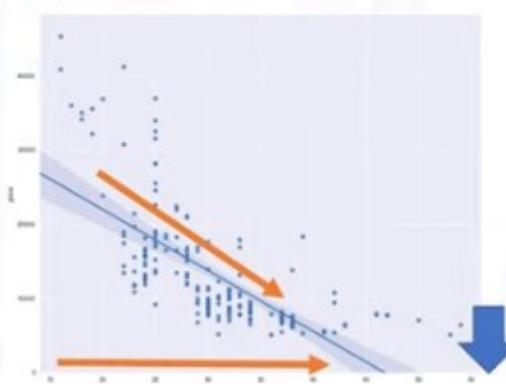
```
array([-37601.57247984, 34779.83910151, 35958.10572339, 35136.37234487,
       34314.8391005, 33615.8391005, 32801.10572339, 31884.10572339,
       31097.77133997, 30201.81207494, 29384.23869462, 28562.5951829,
       27748.77133997, 26918.61864145, 26097.39518333, 25375.57180501,
       24553.83842668, 23632.10504836, 22810.37167004, 21988.63829172,
       21166.90493134, 20345.17153506, 19523.43815675, 18701.70477843,
       17819.97146001, 17058.21802179, 16234.50444347, 15414.77128514,
       14593.83788682, 13771.30459085, 12949.57133818, 12127.83775186,
       11306.10437353, 10484.37099521, 9662.63761689, 8840.98423857,
       8019.17086625, 7297.43748192, 6375.70410389, 5553.91972528,
       4732.23734494, 3918.50294884, 3088.77059631, 2267.03721199,
```

COGNITIVE



## Visualization

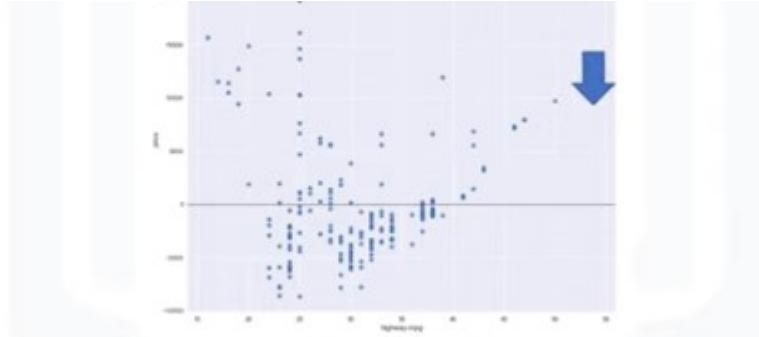
- Simply visualizing your data with a regression



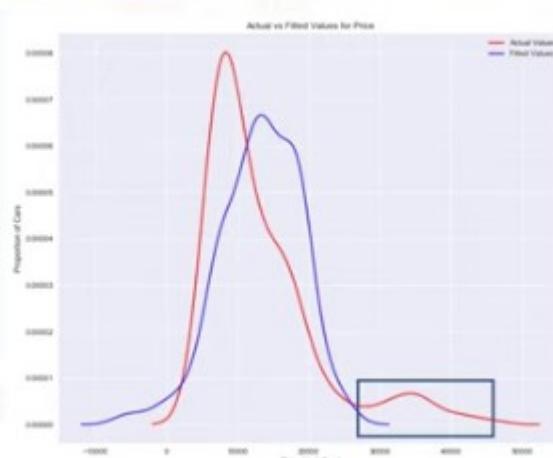
COGNITIVE  
CLASS.ai

© 2017 IBM Corporation

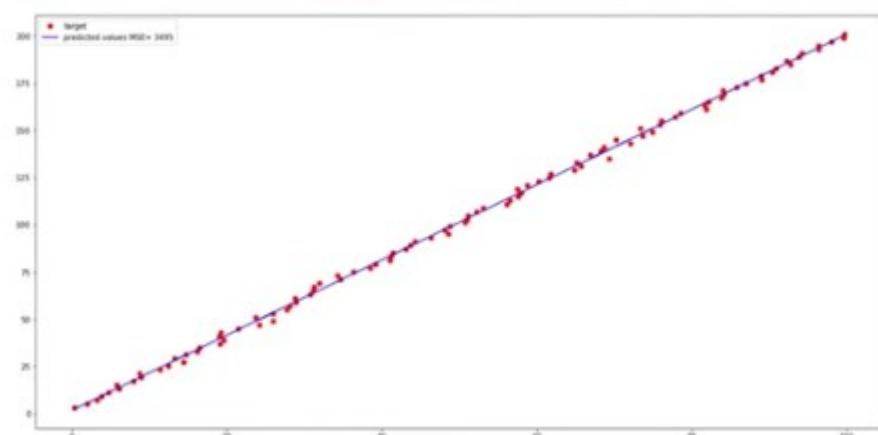
## Residual Plot



## Visualization

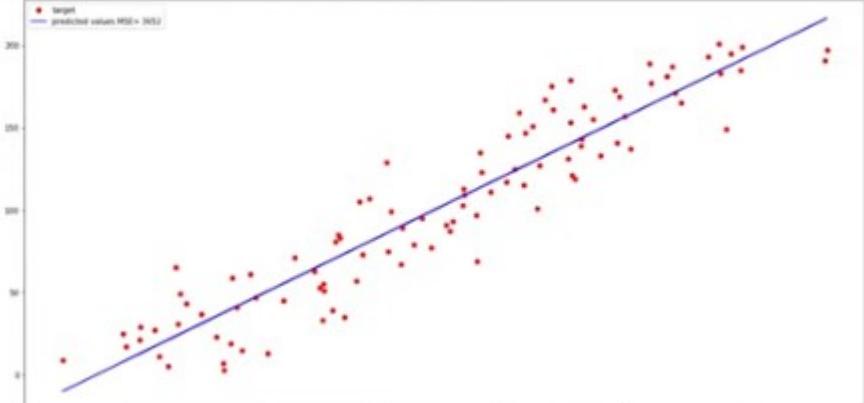


## Numerical measures for Evaluation



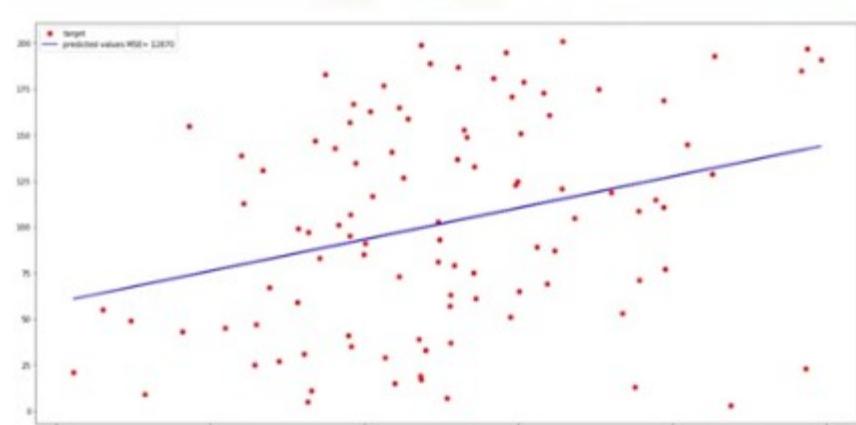
10

## Numerical measures for Evaluation



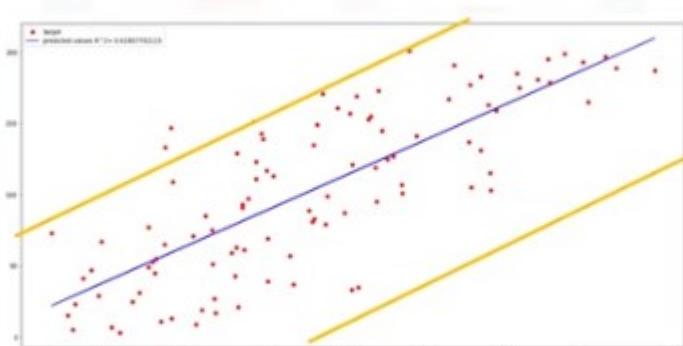
10

## Numerical measures for Evaluation



10

## Numerical measures for Evaluation



11

## Comparing MI R and SI R

## Comparing MLR with SLR

1. Is a lower MSE always implying a better fit?
  - Not necessarily.
2. MSE for an MLR model will be smaller than the MSE for an SLR model, since the errors of the data will decrease when more variables are included in the model
3. Polynomial regression will also have a smaller MSE than regular regression
4. A similar inverse relationship holds for R<sup>2</sup>
5. In the next section we will look at better ways to evaluate the model

## Graded Review Questions

### Question 1

1/1 point (graded)

Let `X` be a dataframe with 100 rows and 5 columns. Let `y` be the target with 100 samples. Assuming all the relevant libraries and data have been imported, the following line of code has been executed:

```
LR = LinearRegression()  
LR.fit(X, y)  
yhat = LR.predict(X)
```

How many samples does `yhat` contain?

 5 500 100 0

[Save](#) | [Show answer](#)

[Submit](#)

You have used 1 of 2 attempts

Correct (1/1 point)

### Question 2

1/1 point (graded)

What value of R<sup>2</sup> (coefficient of determination) indicates your model performs best?

 -100 -1 0 1

### Question 3

1/1 point (graded)

Which statement is true about polynomial linear regression?

- Polynomial linear regression is not linear in any way.
- Although the predictor variables of polynomial linear regression are not linear, the relationship between the parameters or coefficients is linear.
- Polynomial linear regression uses wavelets.

[Save](#) | [Show answer](#)[Submit](#)

You have used 1 of 2 attempts

Correct (1/1 point)

#### Question 4

1/1 point (graded)

The larger the mean squared error, the better your model performs:

- False
- True

[Show answer](#)[Submit](#)

You have used 1 of 1 attempt

Correct (1/1 point)

#### Question 5

1/1 point (graded)

Assume all the libraries are imported.  $y$  is the target and  $X$  is the features or dependent variables. Consider the following lines of code:

```
Input=[('scale',StandardScaler()),('model',LinearRegression())]
pipe=Pipeline(Input)
pipe.fit(X,y)
ypipe=pipe.predict(X)
```

What is the result of  $ypipe$ ?

- Polynomial transform, standardize the data, then perform a prediction using a linear regression model.
- Standardize the data, then perform prediction using a linear regression model.
- Polynomial transform, then standardize the data.

[Save](#) | [Show answer](#)[Submit](#)

You have used 1 of 2 attempts

Correct (1/1 point)

## Module 5 - Model Evaluation

### Model Evaluation and Refinement

#### Model Evaluation and Refinement

### Model Evaluation

#### Model Evaluation

#### Model Evaluation

- In-sample evaluation tells us how well our model will fit the data used to train it
- Problem?
  - It does not tell us how well the trained model can be used to predict new data
- Solution?
  - In- sample data or training data
  - Out-of-sample evaluation or test set

#### Training/Testing Sets

Data:

- Split dataset into:
  - Training set (70%), 
  - Testing set (30%) 

- Build and train the model with a training set
- Use testing set to assess the performance of a predictive model
- When we have completed testing our model we should use all the data to train the model to get the best performance

## Function `train_test_split()`

- Split data into random train and test subsets

```
from sklearn.model_selection import train_test_split
```

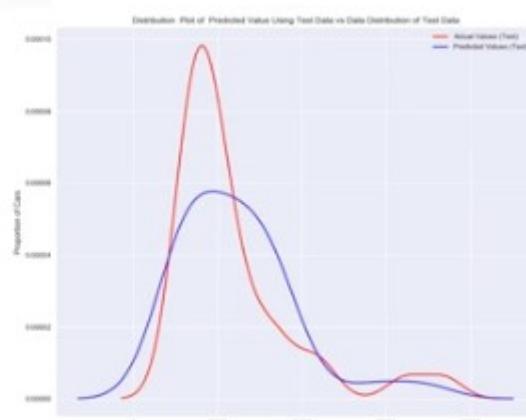
```
x_train, x_test, y_train, y_test = train_test_split(x_data, y_data, test_size=0.3, random_state=0)
```

- `x_data`: features or independent variables
- `y_data`: dataset target: `df['price']`
- `x_train, y_train`: parts of available data as training set
- `x_test, y_test`: parts of available data as testing set
- `test_size`: percentage of the data for testing (here 30% )
- `random_state`: number generator used for random sampling



4

## Generalization Error



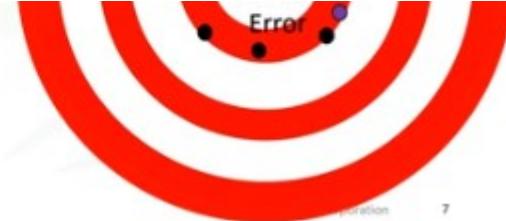
ID 2017 IBM Corporation

6

## Lots of Training Data

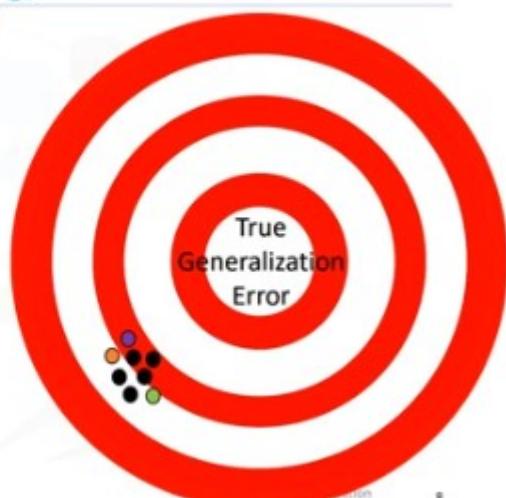
Model →





7

## Lots of Training Data

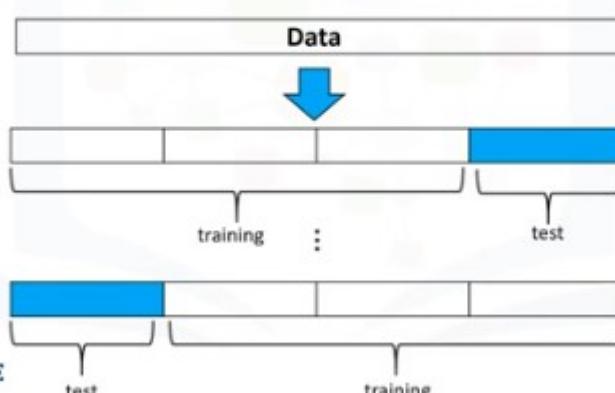


8



## Cross Validation

- Most common out-of-sample evaluation metrics
- More effective use of data (each observation is used for both training and testing)



9

## Function cross\_val\_score()

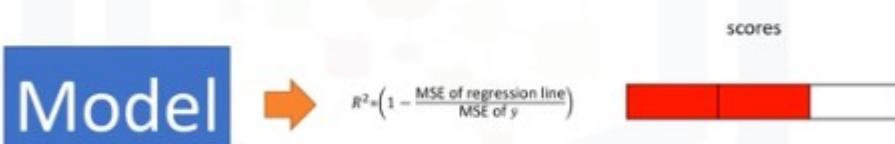
```
from sklearn.model_selection import cross_val_score
```

```
scores= cross_val_score(lr, x_data, y_data, cv=3)
```



```
np.mean(scores)
```

## Function cross\_val\_score()



11

## Function cross\_val\_predict()

- It returns the prediction that was obtained for each element when it was in the test set
- Has a similar interface to cross\_val\_score()

```
from sklearn.model_selection import cross_val_predict
```

```
yhat= cross_val_predict (lr2e, x_data, y_data, cv=3) ←
```

## Function cross\_val\_predict()



Model



## Function cross\_val\_predict()

x <sub>1</sub>	x <sub>2</sub>	x <sub>3</sub>	x <sub>4</sub>	x <sub>5</sub>	x <sub>6</sub>	x <sub>7</sub>	x <sub>8</sub>	x <sub>9</sub>
----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------

Model

$\hat{y}_1$	$\hat{y}_2$	$\hat{y}_3$
$\hat{y}_4$	$\hat{y}_5$	$\hat{y}_6$
$\hat{y}_7$	$\hat{y}_8$	$\hat{y}_9$

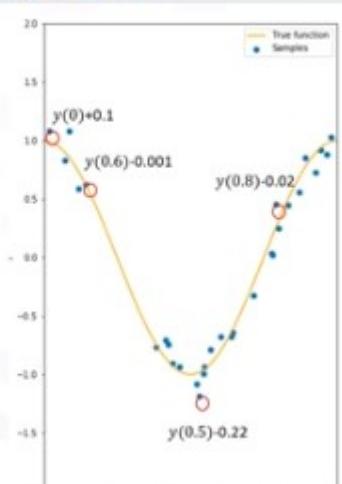


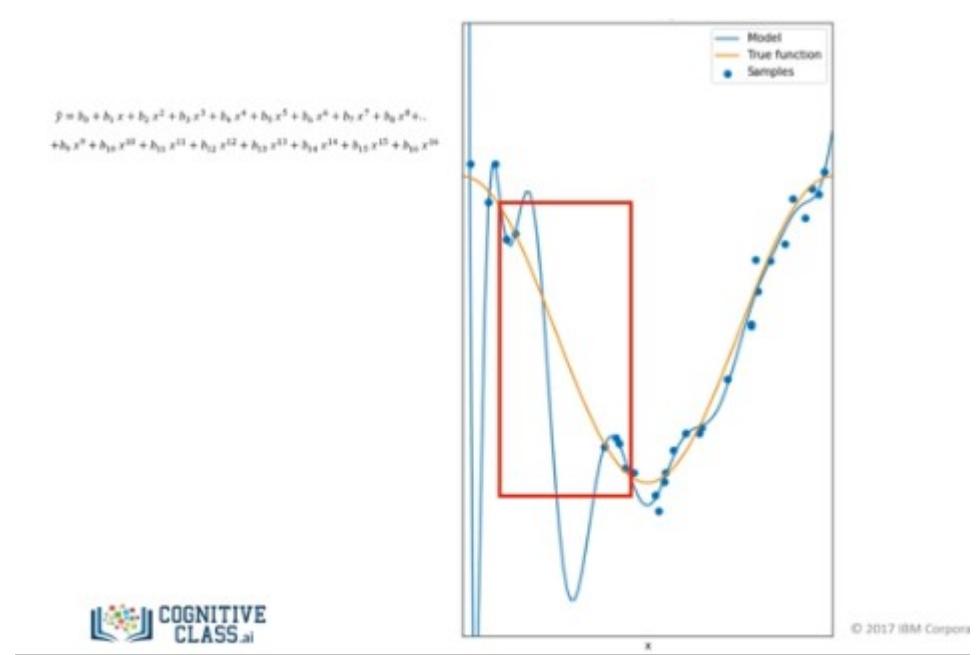
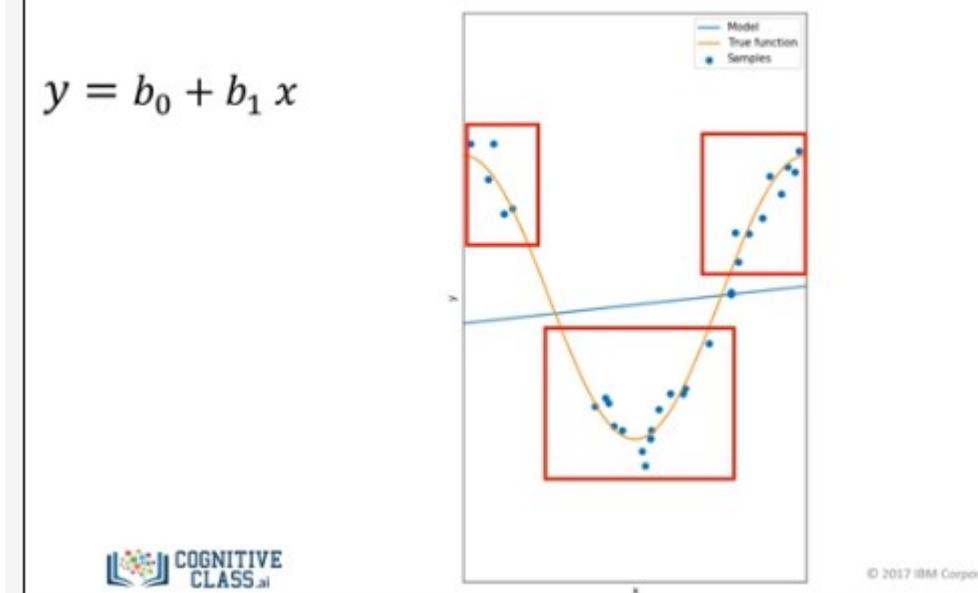
## Overfitting, Underfitting and Model Selection

### Over-fitting, Under-fitting and Model Selection

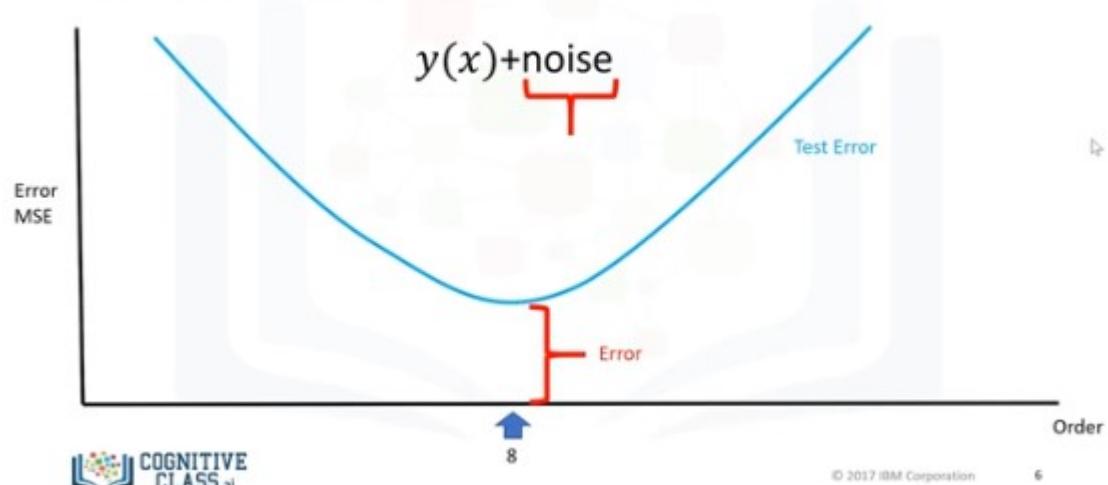
#### Model Selection

y(x)+noise

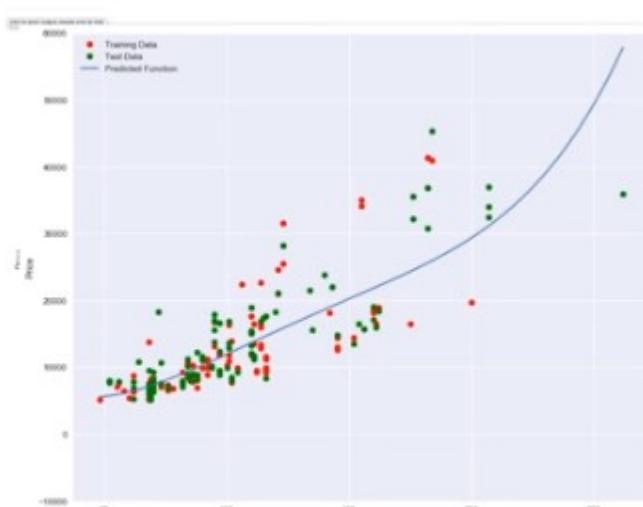
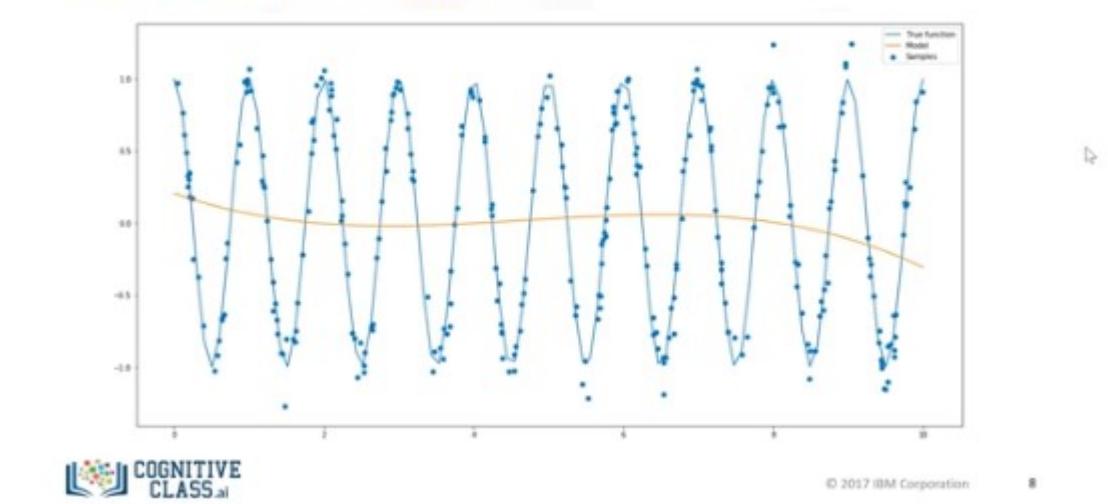




## Model Selection

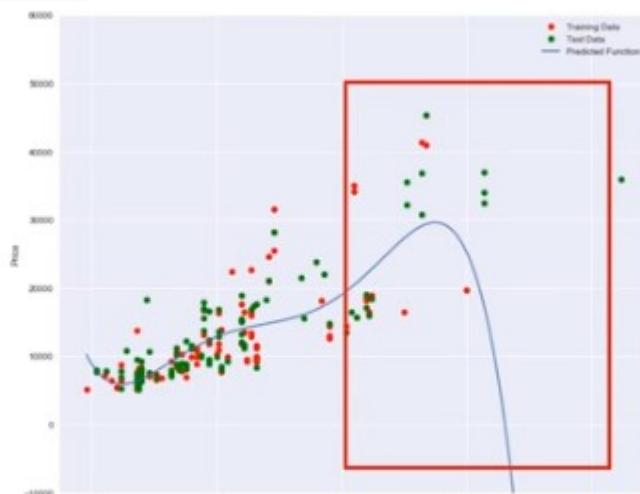


## Model Selection





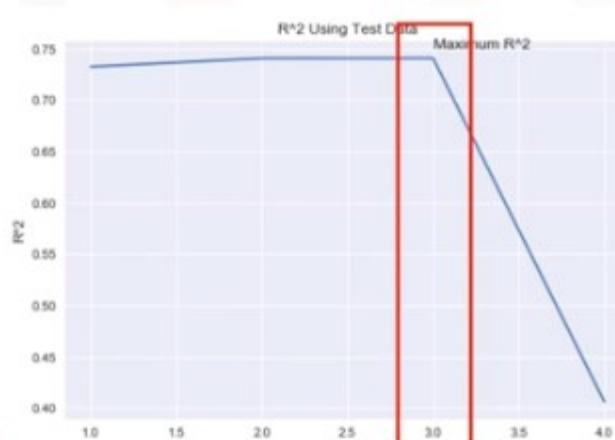
© 2017 IBM Cor



© 2017 IBM Corporation

9

## Model Selection



Corporation

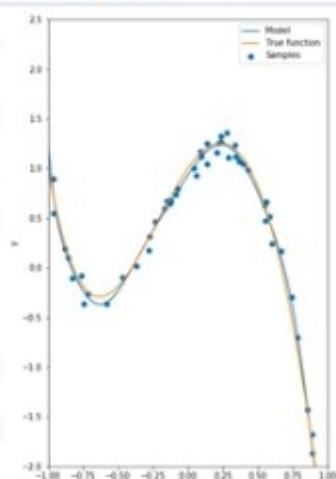
I:\contoso\_DNC1\attachment\contoso\_DNC\

## Ridge Regression

## Ridge Regression

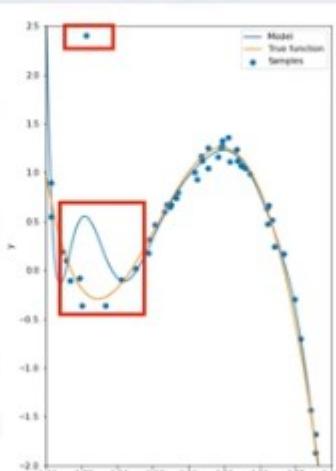
## Ridge Regression

$$y = 1 + 2x - 3x^2 - 4x^3 + x^4$$



## Ridge Regression

$$1 + 2x - 3x^2 - 4x^3 + x^4$$



© 2017 IBM Corporation

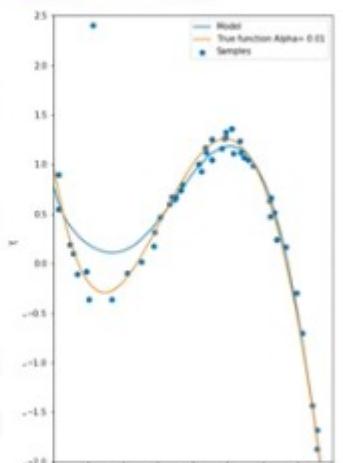
## Ridge Regression

$$\hat{y} = 1 + 2x - 3x^2 - 2x^3 - 12x^4 - 40x^5 + 80x^6 + 71x^7 - 141x^8 - 38x^9 + 75x^{10}$$

Alpha	$x$	$x^2$	$x^3$	$x^4$	$x^5$	$x^6$	$x^7$	$x^8$	$x^9$	$x^{10}$
0	2	-3	-2	-12	-40	80	71	-141	-38	75
0.001	2	-3	-7	5	4	-6	4	-4	4	6
0.01	1	-2	-5	-0.04	0.15	-1	1	-0.5	0.3	1
1	0.5	-1	-1	-0.614	0.70	-0.38	-0.56	-0.21	-0.5	-0.1
10	0	-0.5	-0.3	-0.37	-0.30	-0.30	-0.22	-0.22	-0.22	-0.17

## Ridge Regression

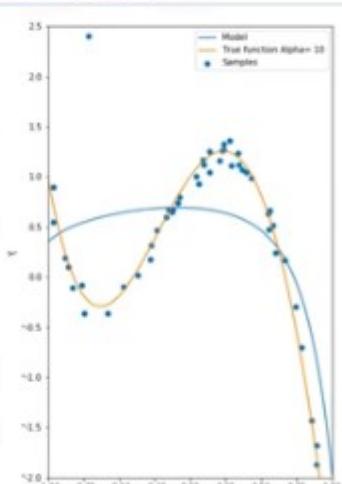
alpha
0
0.001
0.01
1
10



© 2017

## Ridge Regression

alpha
0
0.001
0.01
1
10



© 2017 IBM

## Ridge Regression

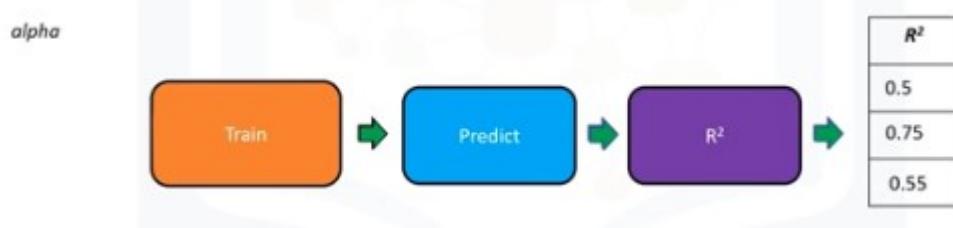
```
from sklearn.linear_model import Ridge
```

```
RidgeModel=Ridge(alpha=0.1)
```

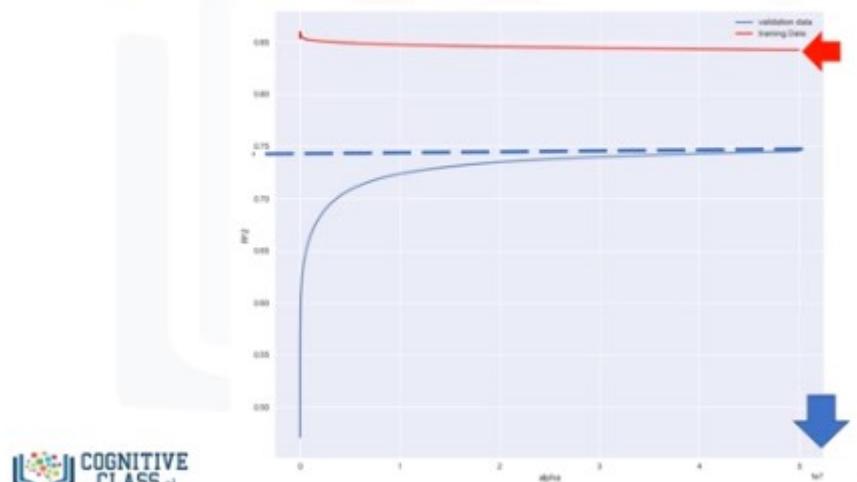
```
RidgeModel.fit(X,y)
```

```
Yhat=RidgeModel.predict(X)
```

## Ridge Regression



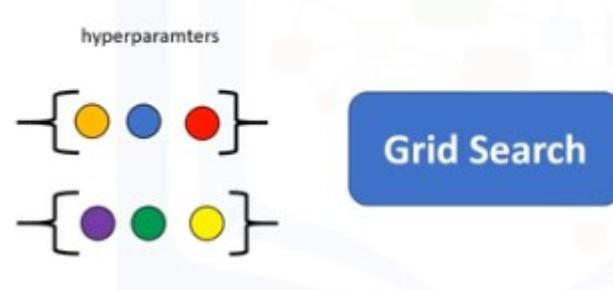
## Ridge Regression



## Grid Search

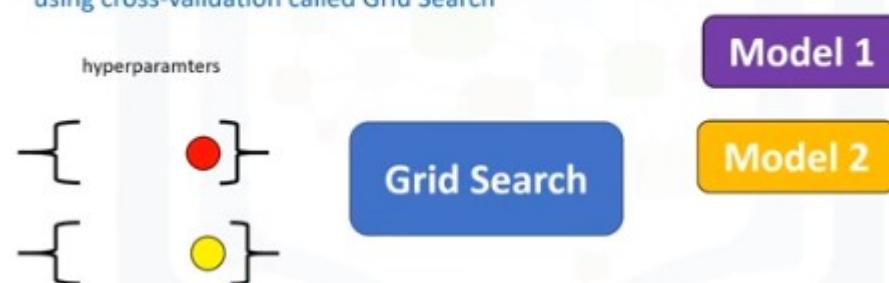
### Hyperparameters

- In the last section, the term alpha in Ridge regression is called a hyperparameter
- Scikit-learn has a means of automatically iterating over these hyperparameters using cross-validation called Grid Search



### Hyperparameters

- In the last section, the term alpha in Ridge regression is called a hyperparameter
- Scikit-learn has a means of automatically iterating over these hyperparameters using cross-validation called Grid Search



## Grid Search

Data



Training

Validation

Test

## Hyperparameters

`class sklearn.linear_model.Ridge`

`class sklearn.linear_model.Ridge(alpha=1.0, fit_intercept=True, normalize=False, copy_X=True, max_iter=None, tol=0.001, solver='auto', random_state=None)`

Linear least squares with l2-regularization.

This model solves a regression model where the loss function is the linear least squares function and regularization is given by the l2-norm. Also known as Ridge Regression or Tikhonov regularization. This estimator has built-in support for multi-variate regression (i.e., when y is a 2d-array of shape [n\_samples, n\_targets]).

Read more in the User Guide.

**Parameters**

`alpha : float, array-like, shape (n_targets)`

Regularization strength; must be a positive float. Regularization improves the conditioning of the problem and reduces the variance of the estimates. Larger values specify stronger regularization. Alpha corresponds to  $\gamma^{-1}$  in other linear models such as LogisticRegression or LinearSVC. If an array is passed, penalties are assumed to be specific to the targets. Hence they must correspond in number.

`fit_intercept : boolean`

Whether to calculate the intercept for this model. If set to false, no intercept will be used in calculations (e.g. data is expected to be already centered).

`normalize : boolean, optional, default=False`

This parameter is ignored when `fit_intercept` is set to False. If True, the regressors X will be normalized before regression by subtracting the mean and dividing by the l2-norm. If you wish to standardize, please use `sklearn.preprocessing.StandardScaler` before calling `Ridge` on an estimator with `normalize=False`.

`copy_X : boolean, optional, default=True`

If True, X will be copied; else, it may be overwritten.

`max_iter : int, optional`

© 2017 IBM Corp.

## Grid Search

parameters = [{ 'alpha': [1, 10, 100, 1000] }]

Alpha	1	10	100	1000
-------	---	----	-----	------

Ridge()

## Grid Search

Grid Search CV

Alpha	1	10	100	1000
R^2	0.74	0.35	0.073	0.008

COGNITIVE

```

from sklearn.linear_model import Ridge
from sklearn.model_selection import GridSearchCV

parameters1= [{"alpha": [0.001, 0.1, 1, 10, 100, 1000, 10000, 100000, 1000000]}]

RR=Ridge()

Grid1 = GridSearchCV(RR, parameters1, cv=4)

Grid1.fit(x_data[['horsepower', 'curb-weight', 'engine-size', 'highway-mpg']], y_data)

Grid1.best_estimator_

scores = Grid1.cv_results_
scores['mean_test_score']

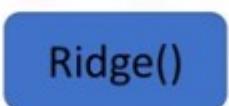
array([ 0.66549413, 0.66554568, 0.66602936, 0.66896822, 0.67334636, 0.65781884, 0.65781884])

```

## Grid Search

parameters = [{ 'alpha' : [1, 10, 100, 1000], 'normalize' : [[True, False]] }]

Alpha	1	10	100	1000
Normalize	True	True	True	True
	False	False	False	False



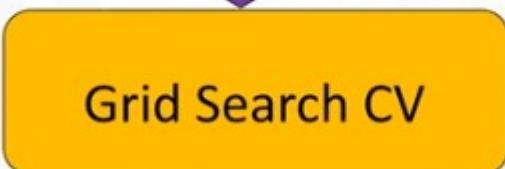

## Grid Search




Scoring



Number  
of Folds

Alpha	1	10	100	1000
True	0.69	0.32	0.17	0.17
False	0.67	0.66	0.66	0.64

```

from sklearn.linear_model import Ridge
from sklearn.model_selection import GridSearchCV

parameters2= [{"alpha": [0.001,0.1,1, 10, 100], 'normalize': [True, False] }]

RR=Ridge()

Grid1 = GridSearchCV(RR, parameters2, cv=4)

Grid1.fit(x_data[['horsepower', 'curb-weight', 'engine-size', 'highway-mpg']], y_data)

Grid1.best_estimator_

scores = Grid1.cv_results_

```

**scores = Grid1.cv\_results\_**

```

for param, mean_val, mean_test in zip(scores['params'], scores['mean_test_score'], scores['mean_train_score']):
    print(param, "R^2 on test data:", mean_val, "R^2 on train data:", mean_test)

```

{'alpha': 0.001, 'normalize': True} R^2 on test data: 0.66605547293 R^2 on train data: 0.814001968709  
{'alpha': 0.001, 'normalize': False} R^2 on test data: 0.665488366584 R^2 on train data: 0.814002698797  
{'alpha': 0.1, 'normalize': True} R^2 on test data: 0.694175625356 R^2 on train data: 0.810546768311  
{'alpha': 0.1, 'normalize': False} R^2 on test data: 0.665488937796 R^2 on train data: 0.814002698794  
{'alpha': 1, 'normalize': True} R^2 on test data: 0.690486934584 R^2 on train data: 0.749104440368  
{'alpha': 1, 'normalize': False} R^2 on test data: 0.665494127178 R^2 on train data: 0.814002698472  
{'alpha': 10, 'normalize': True} R^2 on test data: 0.321376875232 R^2 on train data: 0.341856042902  
{'alpha': 10, 'normalize': False} R^2 on test data: 0.665545680812 R^2 on train data: 0.8140026666  
{'alpha': 100, 'normalize': True} R^2 on test data: 0.0170551710263 R^2 on train data: 0.0496044796826  
{'alpha': 100, 'normalize': False} R^2 on test data: 0.666029359996 R^2 on train data: 0.813999791851  
{'alpha': 1000, 'normalize': True} R^2 on test data: -0.0301961745066 R^2 on train data: 0.005184451599  
{'alpha': 1000, 'normalize': False} R^2 on test data: 0.668968215369 R^2 on train data: 0.813870488264  
{'alpha': 10000, 'normalize': True} R^2 on test data: -0.0351687400461 R^2 on train data: 0.000520784757979  
{'alpha': 10000, 'normalize': False} R^2 on test data: 0.673346359342 R^2 on train data: 0.812583743226  
{'alpha': 100000, 'normalize': True} R^2 on test data: -0.0356685844558 R^2 on train data: 5.2101975528e-05  
{'alpha': 100000, 'normalize': False} R^2 on test data: 0.657818838432 R^2 on train data: 0.789541446486  
{'alpha': 100000, 'normalize': True} R^2 on test data: -0.0356685844558 R^2 on train data: 5.2101975528e-05  
{'alpha': 100000, 'normalize': False} R^2 on test data: 0.657818838432 R^2 on train data: 0.789541446486

## lab 5 esta en el folder

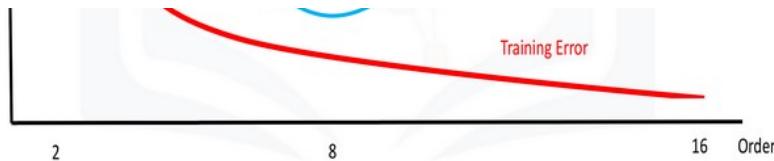
## Graded Review Questions

### Question 1

1/1 point (graded)

In the following plot, the vertical axis shows the mean square error and the horizontal axis represents the order of the polynomial. The red line represents the training error the blue line is the test error. What is the best order of the polynomial given the possible choices in the horizontal axis?





- 2
- 8
- 16

[Save](#) | [Show answer](#)[Submit](#)

You have used 1 of 2 attempts

Correct (1/1 point)

### Question 2

1/1 point (graded)

What is the correct use of the "train\_test\_split" function such that 40% of the data samples will be utilized for testing; the parameter "random\_state" is set to zero; and the input variables for the features and targets are `x_data`, `y_data` respectively?

- `train_test_split(x_data, y_data, test_size=0, random_state=0.4)`
- `train_test_split(x_data, y_data, test_size=0.4, random_state=0)`
- `train_test_split(x_data, y_data)`

[Save](#) | [Show answer](#)[Submit](#)

You have used 1 of 2 attempts

Correct (1/1 point)

### Questions 3

1/1 point (graded)

What is the output of `cross_val_score(lre, x_data, y_data, cv=2)`?

- The predicted values of the test data using cross-validation.
- The average R^2 on the test data for each of the two folds.
- This function finds the free parameter alpha.

[Save](#) | [Show answer](#)[Submit](#)

You have used 1 of 2 attempts

Correct (1/1 point)

### Question 4

1/1 point (graded)

What is the code to create a ridge regression object "RR" with an alpha term equal 10?

RR=LinearRegression(alpha=10)

RR=Ridge(alpha=10)

RR=Ridge(alpha=1)

[Save](#) | [Show answer](#)[Submit](#)

You have used 1 of 2 attempts

Correct (1/1 point)

### Question 5

1/1 point (graded)

What dictionary value would we use to perform a grid search for the following values of alpha: 1,10, 100? No other parameter values should be tested.

alpha=[1,10,100]

[{"alpha": [1,10,100]}]

[{"alpha": [0.001,0.1,1, 10, 100, 1000,10000,100000,100000],'normalize':[True,False]} ]

[Save](#) | [Show answer](#)[Submit](#)

You have used 1 of 2 attempts

Correct (1/1 point)

# Examen Final

## Exam

 [Bookmark this page](#)

### Question 1

1/1 point (graded)

What does the following command do?

```
df.dropna(subset=["price"], axis=0)
```

Drop the "not a number" values from the column "price".

Drop the row "price".

Rename the dataframe "price".

[Show answer](#)[Submit](#)

You have used 2 of 2 attempts

Correct (1/1 point)

### Question 2

1/1 point (graded)

How would you provide many of the summary statistics for all the columns in the dataframe "df"?

df.describe(include = "all")

df.head()

type(df)

df.shape



[Save](#) | [Show answer](#)

[Submit](#)

You have used 1 of 2 attempts

### Question 3

1/1 point (graded)

How would you find the shape of the dataframe df?

df.describe()

df.head()

type(df)

df.shape



[Save](#) | [Show answer](#)

[Submit](#)

You have used 1 of 2 attempts

Correct (1/1 point)

### Question 4

1/1 point (graded)

What task does the following command, df.to\_csv("A.csv"), perform:

Change the name of the column to "A.csv".

Load the data from a csv file called "A" into a dataframe.

Save the dataframe df to a csv file called "A.csv".



[Save](#) | [Show answer](#)

[Submit](#)

You have used 1 of 2 attempts

Correct (1/1 point)

### Question 5

0/1 point (graded)

What task does the following line of code perform?

```
result = np.linspace(min(df["city-mpg"]), max(df["city-mpg"]), 5)
```

Builds a bin array ranging from the smallest value to the largest value of "city-mpg" in order to build 4 bins of equal length.

Builds a bin array ranging from the smallest value to the largest value of "city-mpg" in order to build 5 bins of equal length.

Determines which bin each value of "city-mpg" belongs to.

[Show answer](#)

[Submit](#) You have used 2 of 2 attempts

---

Incorrect (0/1 point)

### Question 6

1/1 point (graded)

What task does the following line of code perform:

```
df['peak-rpm'].replace(np.nan, 5,inplace=True)
```

- Replace the "not a number" values with 5 in the column 'peak-rpm'.

- Rename the column 'peak-rpm' to 5.

- Add 5 to the dataframe.



[Save](#) | [Show answer](#)

[Submit](#)

You have used 1 of 2 attempts

- 
- Correct (1/1 point)

### Question 7

1/1 point (graded)

How do you "one-hot encode" the column 'fuel-type' in the dataframe df?

- pd.get\_dummies(df["fuel-type"])

- df.mean(["fuel-type"])

- df[df["fuel-type"]==1]=1



[Save](#) | [Show answer](#)

[Submit](#)

You have used 1 of 2 attempts

- 
- Correct (1/1 point)

### Question 8

1/1 point (graded)

What does the vertical axis on a scatterplot represent?

- Independent variable

- Dependent variable



[Show answer](#)

[Submit](#)

You have used 1 of 1 attempt

- 
- Correct (1/1 point)

### Question 9

1/1 point (graded)

What does the horizontal axis on a scatterplot represent?

Independent variable

Dependent variable



[Show answer](#)

[Submit](#)

You have used 1 of 1 attempt

Correct (1/1 point)

### Question 10

1/1 point (graded)

If we have 10 columns and 100 samples, how large is the output of df.corr()?

10 x 100

10 x 10

100x100

100x100



[Save](#) | [Show answer](#)

[Submit](#)

You have used 1 of 2 attempts

Correct (1/1 point)

### Question 11

1/1 point (graded)

What is the largest possible element resulting in the following operation "df.corr()"?

100

1000

1



[Save](#) | [Show answer](#)

[Submit](#)

You have used 1 of 2 attempts

Correct (1/1 point)

### Question 12

1/1 point (graded)

If the Pearson Correlation of two variables is zero:

The two variable have zero mean.

The two variables are not correlated.

✓

[Show answer](#)

You have used 1 of 1 attempt

---

✓ Correct (1/1 point)

### Question 13

1/1 point (graded)

If the p-value of the Pearson Correlation is 1:

- The variables are correlated.
- The variables are not correlated.
- None of the above.

✓

[Save](#) [Show answer](#)

You have used 1 of 2 attempts

---

✓ Correct (1/1 point)

### Question 14

1/1 point (graded)

What does the following line of code do: lm = LinearRegression()?

- Fit a regression object "lm".
- Create a linear regression object.
- Predict a value.

✓

[Save](#) [Show answer](#)

You have used 1 of 2 attempts

---

✓ Correct (1/1 point)

### Question 15

0/1 point (graded)

If the predicted function is:

$$\hat{Y} = a + b_1 X_1 + b_2 X_2 + b_3 X_3 + b_4 X_4$$

The method is:

- Polynomial Regression
- Multiple Linear Regression ✓

[Show answer](#)

You have used 1 of 1 attempt

---

Answers are displayed within the problem

Advanced options displayed in the problem

### Question 16

1/1 point (graded)

What steps do the following lines of code perform:

```
Input=[('scale',StandardScaler()),('model',LinearRegression())]  
pipe=Pipeline(Input)  
pipe.fit(Z,y)  
ypipe=pipe.predict(Z)
```

- Standardize the data, then perform a polynomial transform on the features Z.
- Find the correlation between Z and y.
- Standardize the data, then perform a prediction using a linear regression model using the features Z and targets y.



[Save](#) | [Show answer](#)

---

### Question 17

1/1 point (graded)

What is the maximum value of R^2 that can be obtained?

10

1

0



[Save](#) | [Show answer](#)

[Submit](#)

You have used 1 of 2 attempts

---

Correct (1/1 point)

---

### Question 18

1/1 point (graded)

We create a polynomial feature PolynomialFeatures(degree=2). What is the order of the polynomial?

0

1

2



[Save](#) | [Show answer](#)

[Submit](#)

You have used 1 of 2 attempts

---

Correct (1/1 point)

---

### Question 19

1/1 point (graded)

You have a linear model. The average  $R^2$  value on your training data is 0.5. You perform a 100th order polynomial transform on your data, then use these values to train another model. Your average  $R^2$  is 0.99. Which comment is correct?

- 100th order polynomial will work better on unseen data.
- You should always use the simplest model.
- The results on your training data is not the best indicator of how your model performs. You should use your test data to get a better idea.

[Save](#) | [Show answer](#)[Submit](#)

You have used 1 of 2 attempts

Correct (1/1 point)

### Question 20

1/1 point (graded)

You train a ridge regression model. You get a  $R^2$  of 1 on your validation data and you get a  $R^2$  of 0.5 on your training data. What should you do?

- Nothing. Your model performs flawlessly on your validation data.
- Your model is under fitting perform a polynomial transform.
- Your model is overfitting, increase the parameter alpha.

[Save](#) | [Show answer](#)[Submit](#)

You have used 1 of 2 attempts

Correct (1/1 point)

## Badge

In recognition of the commitment to achieve professional excellence



# Victor Hugo Avila Felipe

Has successfully satisfied the requirements for:

Data Analysis Using Python



Issued on: 11 NOV 2022



Issued by IBM

Verify: <https://www.credly.com/go/nPoi4xdm>

