

# **Paper 060010907: Machine Learning**

## **Unit-2 Regression**

Dharmendra Bhatti

1

1

### **Regression Analysis**



- It is a statistical modeling technique.
- It is the process of investigating relationships between dependent and independent variables.
- Aims to predict outputs on a continuous scale rather than categorical class labels.

Prof.(Dr.) Dharmendra Bhatti

2

2

1

## Regression Analysis



- The dependent or target variable is estimated as a function of independent or predictor variables.
- The estimation function is called the *regression function*.

Prof.(Dr.) Dharmendra Bhatti

3

3

## Regression Analysis



- According to Wikipedia,
  - In statistical modeling, regression analysis is a set of statistical processes for estimating the **relationships between a dependent variable** (often called the 'outcome variable') **and one or more independent variables** (often called 'predictors', 'covariates', or 'features').
- [https://en.wikipedia.org/wiki/Regression\\_analysis](https://en.wikipedia.org/wiki/Regression_analysis)

Prof.(Dr.) Dharmendra Bhatti

4

4

2

## Regression Analysis



- It includes,
  - The independent variable X
  - The dependent or target variable Y
  - Unknown parameter(s), denoted as  $\beta$
- $Y = f(X, \beta)$
- The function  $f()$  needs to learn from the training dataset.

Prof.(Dr.) Dharmendra Bhatti

5

## Why Regression Analysis??



- Regression estimates the relationship between the target and the independent variable.
- It is used to find the trends in data.
- It helps to predict real/continuous values.
- It is used to determine the most/least important factor.

Prof.(Dr.) Dharmendra Bhatti

6

6

3

# Types of Regression



## Regression

Linear Regression

Polynomial Regression

Decision Tree Regression

Random Forest Regression

Logistic Regression

Prof.(Dr.) Dharmendra Bhatti

7

# Linear Regression



- Linear regression is the analysis of relationship between the dependent and independent variables.
- Linear regression assumes linear relationship between the two variables.

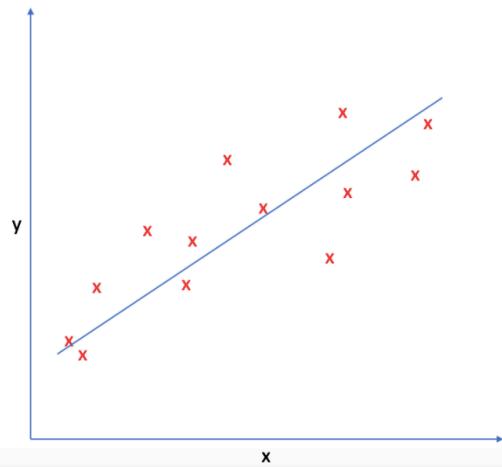
Prof.(Dr.) Dharmendra Bhatti

8

8

4

# Linear Regression



Prof.(Dr.) Dharmendra Bhatti

9

9

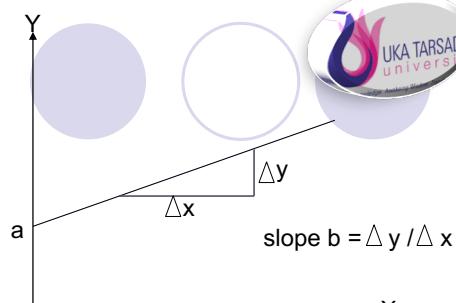
# Linear Regression



- $Y = a + bX$

- Where,

- Y is the dependent variable (that's the variable that goes on the Y axis)
- X is the independent variable (i.e. it is plotted on the X axis)
- b is the slope of the line
- a is the y-intercept



Prof.(Dr.) Dharmendra Bhatti

10

10

## Linear Regression



- According to Wikipedia,
  - In statistics, linear regression is a linear approach to modeling the relationship between a scalar response (or dependent variable) and one or more explanatory variables (or independent variables).
- [https://en.wikipedia.org/wiki/Linear\\_regression](https://en.wikipedia.org/wiki/Linear_regression)

Prof.(Dr.) Dharmendra Bhatti

11

11

## Polynomial Regression



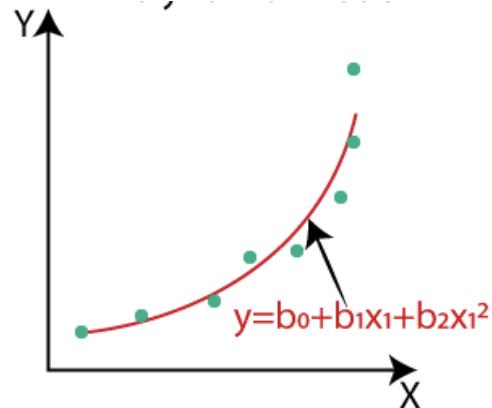
- Dependent variable is related polynomially to independent variable
- Datapoints are best fitted using a polynomial line

Prof.(Dr.) Dharmendra Bhatti

12

12

## Polynomial Regression



Prof.(Dr.) Dharmendra Bhatti

13

13

## Decision Tree Regression

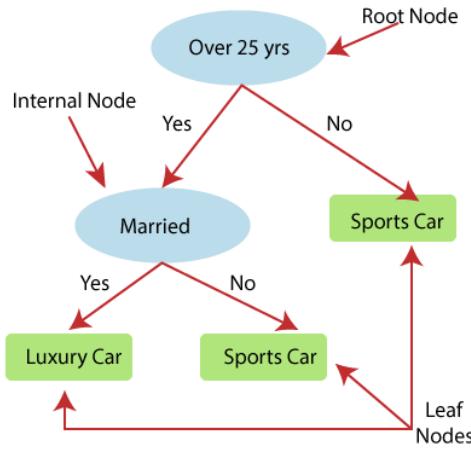
- Builds a tree-like structure in which each internal node represents the "test" for an attribute, each branch represents the result of the test, and each leaf node represents the final decision or result.

Prof.(Dr.) Dharmendra Bhatti

14

14

## Decision Tree Regression



Prof.(Dr.) Dharmendra Bhatti

15

## Random Forest Regression

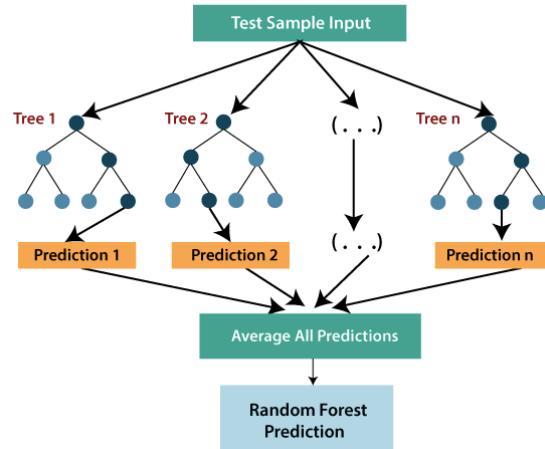
- Combines multiple decision trees and predicts the final output based on the average of each tree output.

Prof.(Dr.) Dharmendra Bhatti

16

16

## Random Forest Regression



Prof.(Dr.) Dharmendra Bhatti

17

17

## Logistic Regression

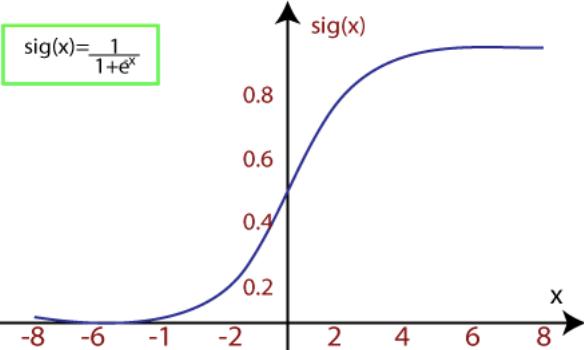
- Used to solve the classification problems
- Works with the categorical variable such as 0 or 1, Yes or No, True or False, Spam or not spam, etc.

Prof.(Dr.) Dharmendra Bhatti

18

18

## Logistic Regression



Prof.(Dr.) Dharmendra Bhatti

19

19

## LINEAR REGRESSION

Prof.(Dr.) Dharmendra Bhatti

20

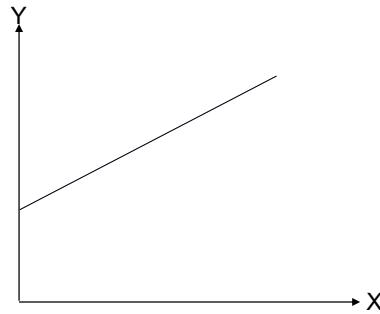
20

10

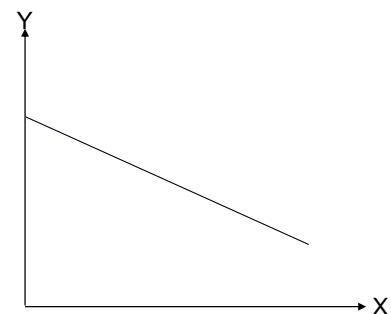
## Linear Regression



- A linear line showing the relationship between the dependent and independent variables.



Prof.(Dr.) Dharmendra Bhatti



Negative Line

21

21

## Linear Regression



- Objective is to find the best fit line.
- Error between predicted values and actual values should be minimized.

Prof.(Dr.) Dharmendra Bhatti

22

22

11

## Linear Regression

### Linear Regression

#### Simple Linear Regression

#### Multiple Linear Regression

Prof.(Dr.) Dharmendra Bhatti

23

23

## Simple Linear Regression

- When we have a single input attribute ( $x$ ) and we want to use linear regression, this is called simple linear regression.

Prof.(Dr.) Dharmendra Bhatti

24

24

12

## Simple Linear Regression



- According to Wikipedia,
  - In statistics, simple linear regression is a linear regression model with a single explanatory variable.
- [https://en.wikipedia.org/wiki/Simple\\_linear\\_regression](https://en.wikipedia.org/wiki/Simple_linear_regression)

Prof.(Dr.) Dharmendra Bhatti

25

25

## Simple Linear Regression



- It models the relationship between a dependent variable and a single independent variable.
- Dependent variable must be continuous/numerical value.

Prof.(Dr.) Dharmendra Bhatti

26

26

13

## Simple Linear Regression



- Examples,

- Increase in website traffic vs SEO efforts
- Product sell vs marketing expenditure
- Reducing air from tyre vs time

Prof.(Dr.) Dharmendra Bhatti

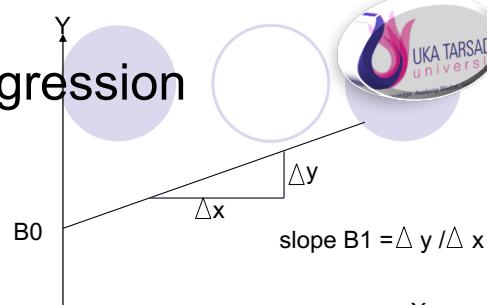
27

27

## Simple Linear Regression



- $Y = B_0 + B_1 X$



- Where,

- Y is the dependent variable (output variable we want to predict)
- X is the independent variable (we know/have it)
- $B_1$  is the slope of the line
- $B_0$  is the y-intercept

Prof.(Dr.) Dharmendra Bhatti

28

28

## Simple Linear Regression



- $Y = B_0 + B_1 X$
- $x$  is the input variable we know and  $B_0$  and  $B_1$  are coefficients that we need to estimate
- The objective is to find the best estimates for the coefficients  $B_0$  and  $B_1$  to minimize the errors in predicting  $y$  from  $x$ .

Prof.(Dr.) Dharmendra Bhatti

29

29

## Simple Linear Regression



- Estimate the slope  $B_1$ 
  - $B_1 = \frac{\sum((x_i - \text{mean}(x)) * (y_i - \text{mean}(y)))}{\sum((x_i - \text{mean}(x))^2)}$
- Calculate  $B_0$ 
  - $B_0 = \text{mean}(y) - B_1 * \text{mean}(x)$

Prof.(Dr.) Dharmendra Bhatti

30

30

15

## Simple Linear Regression

- Sample calculation

- $\text{mean}(x) = 3$

- $\text{mean}(y) = 2.8$

x	y
1	1
2	3
4	3
3	2
5	5

Prof.(Dr.) Dharmendra Bhatti

31

31

## Simple Linear Regression

- Calculate the error of each variable from the mean

x -		
x	mean(x)	mean(x)
1	3	-2
2	3	-1
4	3	1
3	3	0
5	3	2

y -		
y	mean(y)	mean(y)
1	2.8	-1.8
3	2.8	0.2
3	2.8	0.2
2	2.8	-0.8
5	2.8	2.2

Prof.(Dr.) Dharmendra Bhatti

32

32

16

## Simple Linear Regression



- Now, multiply the error for each x with the error for each y

x - mean(x)	y - mean(y)	multiplication
-2	-1.8	3.6
-1	0.2	-0.2
1	0.2	0.2
0	-0.8	0
2	2.2	4.4

- calculate the sum of these multiplications
  - Sum of these multiplications is 8

Prof.(Dr.) Dharmendra Bhatti

33

33

## Simple Linear Regression



- $B1 = \frac{\sum((xi - \text{mean}(x)) * (yi - \text{mean}(y)))}{\sum((xi - \text{mean}(x))^2)}$
- Calculate denominator (sum of the squared differences of each x value from the mean)

- Sum is 10

x - mean(x)	squared
-2	4
-1	1
1	1
0	0
2	4

Prof.(Dr.) Dharmendra Bhatti

34

34

17

## Simple Linear Regression



- Now calculate slope
- $B1 = 8 / 10$
- $B1 = 0.8$

Prof.(Dr.) Dharmendra Bhatti

35

35

## Simple Linear Regression



- Calculate interceptor  $B0$
- $B0 = \text{mean}(y) - B1 * \text{mean}(x)$
- $B0 = 2.8 - 0.8 * 3$
- $B0 = 0.4$

Prof.(Dr.) Dharmendra Bhatti

36

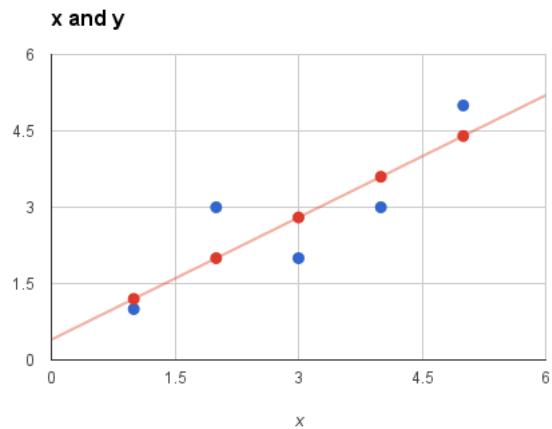
36

18

## Simple Linear Regression

- Now, Predict

- $y = B_0 + B_1 * x$
- $y = 0.4 + 0.8 * x$



Prof.(Dr.) Dharmendra Bhatti

37

37

## Simple Linear Regression

- Calculate error in prediction

- Root Mean Squared Error =  $\sqrt{\sum (y_i - \hat{y}_i)^2 / n}$

predicted y	y	error
1.2	1	0.2
2	3	-1
3.6	3	0.6
2.8	2	0.8
4.4	5	-0.6

Prof.(Dr.) Dharmendra Bhatti

38

38

19

## Simple Linear Regression



- calculate the square of each of these error values

error	squared error
0.2	0.04
-1	1
0.6	0.36
0.8	0.64
-0.6	0.36

- The sum of these errors is 2.4

Prof.(Dr.) Dharmendra Bhatti

39

39

## Simple Linear Regression



- Dividing by n and taking the square root

$$\text{Root Mean Square Error} = \sqrt{2.4 / 5}$$

$$\text{Root Mean Square Error} = 0.692$$

- Each prediction is on average wrong by about 0.692

Prof.(Dr.) Dharmendra Bhatti

40

40

20

## Multiple Linear Regression



- It is the extension of simple linear regression that predicts a response using two or more features.
- It models the linear relationship between a single dependent continuous variable and more than one independent variables.

Prof.(Dr.) Dharmendra Bhatti

41

41

## Multiple Linear Regression



- The dependent or target variable(Y) must be the continuous/real, but the predictor or independent variable may be of continuous or categorical form.

Prof.(Dr.) Dharmendra Bhatti

42

42

## Multiple Linear Regression



- Each feature variable must model the linear relationship with the dependent variable.
- It tries to fit a regression line through a multidimensional space of data-points.

Prof.(Dr.) Dharmendra Bhatti

43

43

## Multiple Linear Regression



- $y = b_0 + b_1x_1 + b_2x_2 + \dots + b_nx_n$
- Where,
  - $y$ = Output/Response variable
  - $b_0, b_1, b_2, b_3, b_n, \dots$ = Coefficients of the model.
  - $x_1, x_2, x_3, x_4, \dots$ = Various Independent/feature variable

Prof.(Dr.) Dharmendra Bhatti

44

44

## Multiple Linear Regression



### ● Assumptions

- A linear relationship should exist between the Target and predictor variables.
- The regression residuals must be normally distributed.
- No multicollinearity (correlation between the independent variable) in data exists.

Prof.(Dr.) Dharmendra Bhatti

45

45

## Multiple Linear Regression



### ● Example,

- Company is spending money on R&D, Administration, Marketing
- Companies operate in different states
- Profit for a financial year is available

Prof.(Dr.) Dharmendra Bhatti

46

46

23

## Multiple Linear Regression



- Importing Libraries

```
11  # Importing the libraries
12  import numpy as np
13  import matplotlib.pyplot as plt
14  import pandas as pd
```

Prof.(Dr.) Dharmendra Bhatti

47

47

## Multiple Linear Regression



- Importing Dataset

- # Importing the dataset
- dataset = pd.read\_csv('Invest2Profit.csv')

Prof.(Dr.) Dharmendra Bhatti

48

48

24

## Multiple Linear Regression

Independent Variables

Dependent Variable

Index	R&D Spend	Administration	Marketing Spend	State	Profit
0	165349	136898	471784	New York	192262
1	162598	151378	443899	California	191792
2	153442	101146	407935	Florida	191050
3	144372	118672	383200	New York	182902
4	142107	91391.8	366168	Florida	166188
5	131877	99814.7	362861	New York	156991
6	134615	147199	127717	California	156123
7	130298	145530	323877	Florida	155753
8	120543	148719	311613	New York	152212
9	123335	108679	304982	California	149760
10	101913	110594	229161	Florida	146122
11	100672	91790.6	249745	California	144259
12	93863.8	127320	249839	Florida	141586

Prof.(Dr.) Dharmendra Bhatti

49

49

## Multiple Linear Regression

- #Extracting Independent and dependent Variable
- X= dataset.iloc[:, :-1].values
- y= dataset.iloc[:, 4].values
- print(X)

Prof.(Dr.) Dharmendra Bhatti

50

50

25

# Multiple Linear Regression



## ● Encoding Dummy Variables

```
○ # Encoding categorical data
○ from sklearn.compose import ColumnTransformer
○ from sklearn.preprocessing import OneHotEncoder
○ ct = ColumnTransformer(
○     transformers=[
○         ('encoder', OneHotEncoder(), [3]),
○         remainder='passthrough'
○     )
○ X = np.array(ct.fit_transform(X))
○ print(X)
```

Prof.(Dr.) Dharmendra Bhatti

51

51

# Multiple Linear Regression



## ● X - array

X - NumPy object array (read only)

	0	1	2	3	4	5
0	0.0	0.0	1.0	165349.2	136897.8	471784.1
1	1.0	0.0	0.0	162597.7	151377.59	443898.53
2	0.0	1.0	0.0	153441.51	101145.55	407934.54
3	0.0	0.0	1.0	144372.41	118671.85	383199.62
4	0.0	1.0	0.0	142107.34	91391.77	366168.42
5	0.0	0.0	1.0	131876.9	99814.71	362861.36
6	1.0	0.0	0.0	134615.46	147198.87	127716.82
7	0.0	1.0	0.0	130298.13	145530.06	323876.68
8	0.0	0.0	1.0	120542.52	148718.95	311613.29

Format    Resize     Background color

Close

52

52

26

## Multiple Linear Regression

- We **should not** use all the dummy variables at the same time
- So it must be 1 less than the total number of dummy variables, else it will create a dummy variable trap.

Prof.(Dr.) Dharmendra Bhatti

53

53

## Multiple Linear Regression

- #avoiding the dummy variable trap:
- $X = X[:, 1:]$

	0	1	2	3	4
0	0.0	1.0	165349.2	136897.8	471784.1
1	0.0	0.0	162597.7	151377.59	443898.53
2	1.0	0.0	153441.51	101145.55	407934.54
3	0.0	1.0	144372.41	118671.85	383199.62
4	1.0	0.0	142107.34	91391.77	366168.42
5	0.0	1.0	131876.9	99814.71	362861.36
6	0.0	0.0	134615.46	147198.87	127716.82
7	1.0	0.0	130298.13	145530.06	323876.68
8	0.0	1.0	120542.52	148718.95	311613.29

54

54

27

## Multiple Linear Regression



- # Splitting the dataset into training and test set.
- from sklearn.model\_selection import train\_test\_split
- X\_train, X\_test, y\_train, y\_test= train\_test\_split(X, y, test\_size= 0.2, random\_state=0)

Prof.(Dr.) Dharmendra Bhatti

55

55

## Multiple Linear Regression



- #Fitting the MLR model to the training set:
- from sklearn.linear\_model import LinearRegression
- regressor= LinearRegression()
- regressor.fit(X\_train, y\_train)

Prof.(Dr.) Dharmendra Bhatti

56

56

## Multiple Linear Regression

- #Predicting the Test set result;
- `y_pred= regressor.predict(X_test)`

Prof.(Dr.) Dharmendra Bhatti

57

57

## Multiple Linear Regression

- Now, compare `y_pred` and `y_test`

0	
0	103015
1	132582
2	132448
3	71976.1
4	178537
5	116161
6	67851.7
7	98791.7
8	113969

0	
0	103282
1	144259
2	146122
3	77798.8
4	191050
5	105008
6	81229.1
7	97483.6
8	110352

58

29

## Multiple Linear Regression

- `print("Train Score: ", regressor.score(X_train, y_train))`
- `print("Test Score: ", regressor.score(X_test, y_test))`

- **Output:**

```
Train Score:  0.9501847627493607  
Test Score:  0.9347068473282446
```

- Model is 95% accurate with the training dataset and 93% accurate with the test dataset.

Prof.(Dr.) Dharmendra Bhatti

59

## Multiple Linear Regression – Backward Elimination

- Feature selection is used to remove features that do not have a significant effect on the dependent variable or prediction of output.

Prof.(Dr.) Dharmendra Bhatti

60

60

30

## Multiple Linear Regression – Backward Elimination



- Feature selection algorithms:

- All-in
- Backward Elimination
- Forward Selection
- Bidirectional Elimination
- Score Comparison

Prof.(Dr.) Dharmendra Bhatti

61

61

## Multiple Linear Regression – Backward Elimination



Which is  
Most/least  
significant  
variable?

Independent Variables

Dependent Variable

	R&D Spend	Administration	Marketing Spend	State	Profit
1	165349	136898	471784	New York	192262
2	162598	151378	443899	California	191792
3	153442	101146	407935	Florida	191050
4	144372	118672	383200	New York	182902
5	142107	91391.8	366168	Florida	166188
6	131877	99814.7	362861	New York	156991
7	134615	147199	127717	California	156123
8	130298	145530	323877	Florida	155753
9	120543	148719	311613	New York	152212
10	123335	108679	304982	California	149760
11	101913	110594	229161	Florida	146122
12	100672	91790.6	249745	California	144259
	93863.8	127320	249839	Florida	141586

Format Resize Background color Column min/max

Save and Close Close

Prof.(Dr.) Dharmendra Bhatti

62

62

## Multiple Linear Regression – Backward Elimination



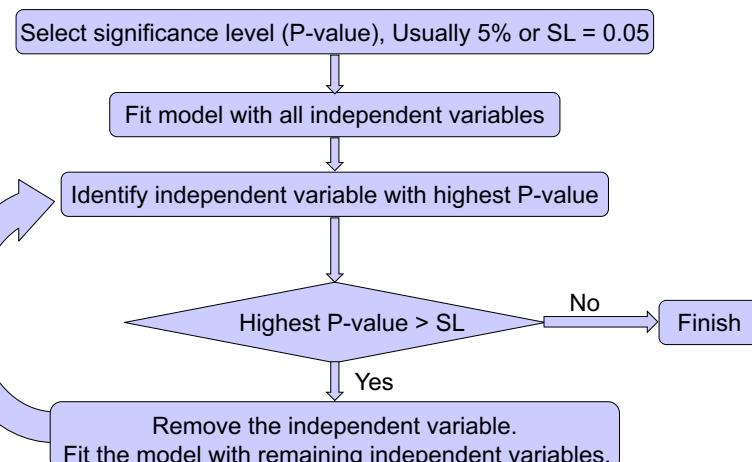
- Objective:
  - Use only statistically significant independent variable(s)
  - Find optimal matrix of features
- Each of the selected independent variable is having great impact on the dependent variable.

Prof.(Dr.) Dharmendra Bhatti

63

63

## Multiple Linear Regression – Backward Elimination



Prof.(Dr.) Dharmendra Bhatti

64

64

## Multiple Linear Regression – Backward Elimination



- # Backward Elimination
- import statsmodels.api as sm

Prof.(Dr.) Dharmendra Bhatti

65

65

## Multiple Linear Regression – Backward Elimination



$$y = b_0 + b_1x_1 + b_2x_2 + \dots + b_nx_n$$

$$y = b_0x_0 + b_1x_1 + b_2x_2 + \dots + b_nx_n$$

- Where,
  - $y$ = Output/Response variable
  - $b_0, b_1, b_2, b_3, b_n, \dots$ = Coefficients of the model.
  - $x_1, x_2, x_3, x_4, \dots$ = Various Independent/feature variable
  - $x_0 = 1$

Prof.(Dr.) Dharmendra Bhatti

66

66

33

## Multiple Linear Regression – Backward Elimination

- Add column of 1s at position 0
- $X = np.append(arr = np.ones((50, 1)).astype(int), values = X, axis = 1)$

Prof.(Dr.) Dharmendra Bhatti

67

67

## Multiple Linear Regression – Backward Elimination

- Original X

	0	1	2	3	4	
0	0.0	1.0	165349.2	136897.8	471784.1	
1	0.0	0.0	162597.7	151377.59	443898.53	
2	1.0	0.0	153441.51	101145.55	407934.54	
3	0.0	1.0	144372.41	118671.85	383199.62	
4	1.0	0.0	142107.34	91391.77	366168.42	
5	0.0	1.0	131876.9	99814.71	362861.36	
6	0.0	0.0	134615.46	147198.87	127716.82	
7	1.0	0.0	130298.13	145530.06	323876.68	
8	0.0	1.0	120542.52	148718.95	311613.29	

Close

68

68

## Multiple Linear Regression – Backward Elimination

- After adding ones

X - NumPy object array (read only)						
0	1	2	3	4	5	
0	1	0.0	1.0	165349.2	136897.8	471784.1
1	1	0.0	0.0	162597.7	151377.59	443898.53
2	1	1.0	0.0	153441.51	101145.55	407934.54
3	1	0.0	1.0	144372.41	118671.85	383199.62
4	1	1.0	0.0	142107.34	91391.77	366168.42
5	1	0.0	1.0	131876.9	99814.71	362861.36
6	1	0.0	0.0	134615.46	147198.87	127716.82
7	1	1.0	0.0	130298.13	145530.06	323876.68
8	1	0.0	1.0	120542.52	148718.95	311613.29

Close

69

69

## Multiple Linear Regression – Backward Elimination

- Use all independent variables
- $X_{opt} = X[:, [0, 1, 2, 3, 4, 5]]$
- $X_{opt} = X_{opt}.astype(np.float64)$
- `regressor_OLS = sm.OLS(endog = y, exog = X_opt).fit()`
- `regressor_OLS.summary()`

Prof.(Dr.) Dharmendra Bhatti

70

70

## Multiple Linear Regression – Backward Elimination

- $X_{\text{opt}} = X[:, [0, 1, 2, 3, 4, 5]]$

	coef	std err	t	P> t	[0.025	0.975]
const	5.013e+04	6884.820	7.281	0.000	3.62e+04	6.4e+04
x1	198.7888	3371.007	0.059	0.953	-5595.030	6992.607
x2	-41.8870	3256.039	-0.013	0.990	-5604.003	6520.229
x3	0.8060	0.046	17.369	0.000	0.712	0.900
x4	-0.0270	0.052	-0.517	0.608	-0.132	0.078
x5	0.0270	0.017	1.574	0.123	-0.008	0.062
<hr/>						
Omnibus:		14.782	Durbin-Watson:		1.283	
Prob(Omnibus):		0.001	Jarque-Bera (JB):		21.266	
Skew:		-0.948	Prob(JB):		2.41e-05	
Kurtosis:		5.572	Cond. No.		1.45e+06	
<hr/>						

Prof.(Dr.) Dharmendra Bhatti

71

71

## Multiple Linear Regression – Backward Elimination

- Highest P-value is 0.990 for independent variable 2
- So remove 2nd independent variable from the X

X - NumPy object array (read only)

	0	1	2	3	4	5
0	1	0.0	1.0	15349.2	136897.8	471784.1
1	1	0.0	0.0	162597.7	151377.59	443898.53
2	1	1.	0.0	15341.51	101145.55	407934.54
3	1	0.0	1.0	144372.41	118671.85	383199.62
4	1	1.0	0.0	142167.34	91391.77	366168.42
5	1	0.0	1.0	13186.9	99814.71	362861.36
6	1	0.	0.0	13465.46	147198.87	127716.82
7	1	1.0	0.0	13098.13	145530.86	323876.68
8	1	0.0	1.0	10542.52	148718.95	311613.29

Format Resize Background color Close

72

## Multiple Linear Regression – Backward Elimination

- $X_{\text{opt}} = X[:, [0, 1, 3, 4, 5]]$

	coef	std err	t	P> t	[0.025	0.975]
const	5.011e+04	6647.870	7.537	0.000	3.67e+04	6.35e+04
x1	220.1585	2900.536	0.076	0.940	-5621.821	6062.138
x2	0.8060	0.046	17.606	0.000	0.714	0.898
x3	-0.0270	0.052	-0.523	0.604	-0.131	0.077
x4	0.0270	0.017	1.592	0.118	-0.007	0.061

Prof.(Dr) Dharmendra Bhatti

73

73
74

## Multiple Linear Regression – Backward Elimination

- Highest P-value is 0.940 for independent variable 1
- So remove 1st independent variable from the X

X - NumPy object array (read only)

	0	1	2	3	4	5
0	1	0.0	1.0	165349.2	136897.8	471784.1
1	1	0.0	0.0	162597.7	151377.59	443898.53
2	1	1.0	0.0	153441.51	101145.55	407934.54
3		0.0	1.0	144372.41	118671.85	383199.62
4		1.0	0.0	142107.34	91391.77	366168.42
5		0.0	1.0	131876.9	99814.71	362861.36
6		1	0.0	134615.46	147198.87	127716.82
7		1	1.0	130298.13	145530.86	323876.68
8		1	0.0	120542.52	148718.95	311613.29

Format    Resize    Background color    Close

37

## Multiple Linear Regression – Backward Elimination

- $X_{\text{opt}} = X[:, [0, 3, 4, 5]]$

	coef	std err	t	P> t	[0.025	0.975]
const	5.012e+04	6572.353	7.626	0.000	3.69e+04	6.34e+04
x1	0.8057	0.045	17.846	0.000	0.715	0.897
x2	-0.0268	0.051	-0.526	0.602	-0.130	0.076
x3	0.0272	0.016	1.655	0.105	-0.006	0.060

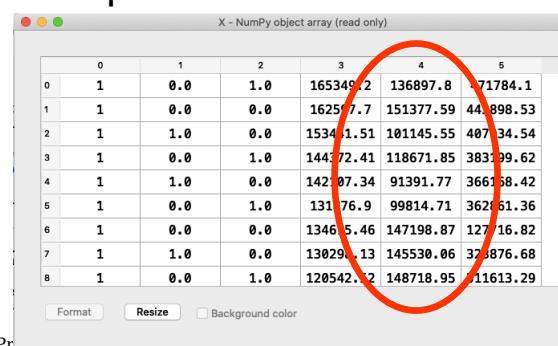
Prof.(Dr.) Dharmendra Bhatti

75

75

## Multiple Linear Regression – Backward Elimination

- Highest P-value is 0.602 for independent variable 4
- So remove 4th independent variable from the X



X - NumPy object array (read only)

	0	1	2	3	4	5
0	1	0.0	1.0	165349.2	136897.8	71784.1
1	1	0.0	0.0	162577.7	151377.59	443898.53
2	1	1.0	0.0	15341.51	101145.55	40734.54
3	1	0.0	1.0	14472.41	118671.85	383399.62
4	1	1.0	0.0	14207.34	91391.77	366158.42
5	1	0.0	1.0	13176.9	99814.71	36251.36
6	1	0.0	0.0	13465.46	147198.87	12716.82
7	1	1.0	0.0	130291.13	145530.86	373876.68
8	1	0.0	1.0	120542.2	148718.95	11613.29

Format Resize Background color Close

76

## Multiple Linear Regression – Backward Elimination

- X\_opt = X[:, [0, 3, 5]]

	coef	std err	t	P> t	[0.025	0.975]
const	4.698e+04	2689.933	17.464	0.000	4.16e+04	5.24e+04
x1	0.7966	0.041	19.266	0.000	0.713	0.880
x2	0.0299	0.016	1.927	0.060	-0.001	0.061

Prof.(Dr.) Dharmendra Bhatti

77

77

## Multiple Linear Regression – Backward Elimination

- Highest P-value is 0.060 for independent variable 5
- So remove 5th independent variable from the X

X - NumPy object array (read only)

	0	1	2	3	4	5
0	1	0.0	1.0	165349.2	136897.8	471784.1
1	1	0.0	0.0	162597.7	151377.59	443898.53
2	1	1.0	0.0	153441.51	10115.55	407934.54
3	1	0.0	1.0	144372.41	11861.85	383199.62
4	1	1.0	0.0	142107.34	91351.77	366168.42
5	1	0.0	1.0	131876.9	99811.71	362861.36
6	1	0.0	0.0	134615.46	147151.87	127716.82
7	1	1.0	0.0	130298.13	145536.06	323876.68
8	1	0.0	1.0	120542.52	148718.55	311613.29

Format Resize Background color Close

78

39

## Multiple Linear Regression – Backward Elimination

- $X_{opt} = X[:, [0, 3]]$

	coef	std err	t	P> t	[0.025	0.975]
const	4.903e+04	2537.897	19.320	0.000	4.39e+04	5.41e+04
x1	0.8543	0.029	29.151	0.000	0.795	0.913

- Highest P-value is 0.000 which is < significance level 0.05 so our model is optimal and ready

Prof.(Dr.) Dharmendra Bhatti

79

## Polynomial Regression

- Simple Linear Regression is good if linear relationship is present between dependent variable and independent variable.
- What if the relationship is non-linear or curvilinear??

Prof.(Dr.) Dharmendra Bhatti

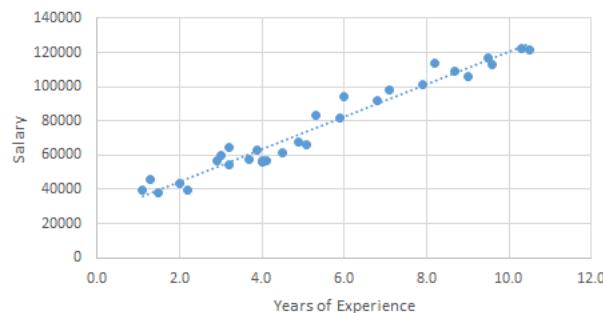
80

80

40

# Polynomial Regression

## ● Simple Linear Regression



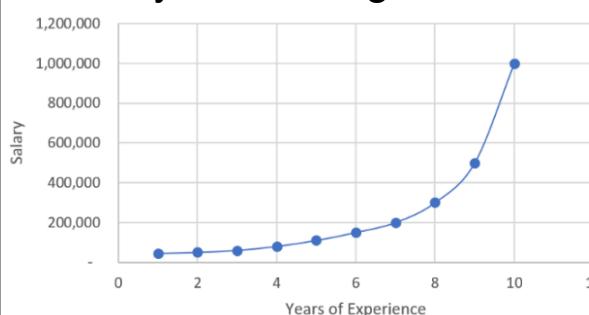
Prof.(Dr.) Dharmendra Bhatti

81

81

# Polynomial Regression

## ● Polynomial Regression



A	B	C
Position	Level	Salary
Business Analyst	1	45000
Junior Consultant	2	50000
Senior Consultant	3	60000
Manager	4	80000
Country Manager	5	110000
Region Manager	6	150000
Partner	7	200000
Senior Partner	8	300000
C-level	9	500000
CEO	10	1000000

Prof.(Dr.) Dharmendra Bhatti

82

82

## Polynomial Regression



- Polynomial regression is a special case of linear regression where we fit a polynomial equation on the data with a curvilinear relationship between the target variable and the independent variables.

Prof.(Dr.) Dharmendra Bhatti

83

83

## Polynomial Regression



- Simple Linear Regression

$$y = b_0 + b_1 x_1$$

- Multiple Linear Regression

$$y = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_n x_n$$

- Polynomial Regression

$$y = b_0 + b_1 x_1 + b_2 x_1^2 + \dots + b_n x_1^n$$

84

84

42

## Polynomial Regression



- Polynomial regression, like linear regression, uses the relationship between the variables x and y to find the best way to draw a line through the data points.

Prof.(Dr.) Dharmendra Bhatti

85

85

## Polynomial Regression



- Position\_Salaries.csv

	A	B	C
1	Position	Level	Salary
2	Business Analyst	1	45000
3	Junior Consultant	2	50000
4	Senior Consultant	3	60000
5	Manager	4	80000
6	Country Manager	5	110000
7	Region Manager	6	150000
8	Partner	7	200000
9	Senior Partner	8	300000
10	C-level	9	500000
11	CEO	10	1000000

86

86

43

# Polynomial Regression



- Importing Libraries

```
9  # Polynomial Regression
10
11 # Importing the libraries
12 import numpy as np
13 import matplotlib.pyplot as plt
14 import pandas as pd
```

Prof.(Dr.) Dharmendra Bhatti

87

87

# Polynomial Regression



- Importing Dataset

```
16 # Importing the dataset
17 dataset = pd.read_csv('Position_Salaries.csv')
18 X = dataset.iloc[:, 1:-1].values
19 y = dataset.iloc[:, -1].values
```

Prof.(Dr.) Dharmendra Bhatti

88

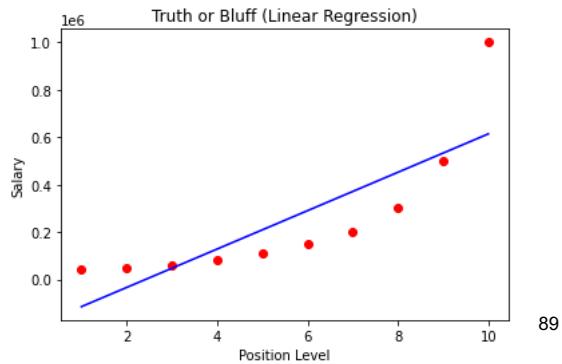
88

44

# Polynomial Regression

## • Training Linear Model

```
21 # Training the Linear Regression model on the whole dataset
22 from sklearn.linear_model import LinearRegression
23 lin_reg = LinearRegression()
24 lin_reg.fit(X, y)
```

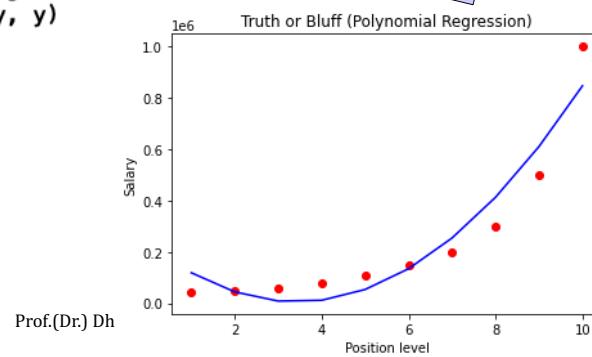


89

# Polynomial Regression

## • Training Polynomial Model

```
26 # Training the Polynomial Regression model on the whole dataset
27 from sklearn.preprocessing import PolynomialFeatures
28 poly_reg = PolynomialFeatures(degree = 2)
29 X_poly = poly_reg.fit_transform(X)
30 lin_reg_2 = LinearRegression()
31 lin_reg_2.fit(X_poly, y)
```

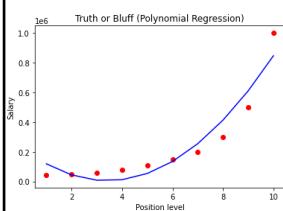


90

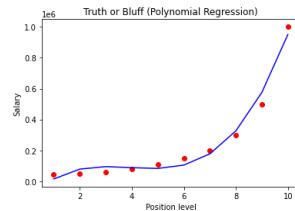
# Polynomial Regression

## • Training Polynomial Model

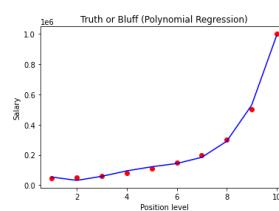
Degree = 2



Degree = 3



Degree = 4



Prof.(Dr.) Dharmendra Bhatti

91

91

# Polynomial Regression

## • # Visualizing the Polynomial Regression results (for higher resolution and smoother curve)

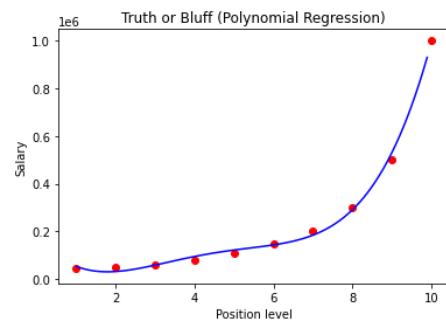
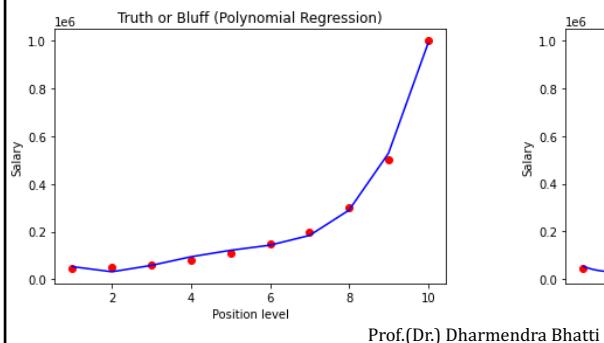
```
49  # Visualising the Polynomial Regression results
50  # (for higher resolution and smoother curve)
51  X_grid = np.arange(min(X), max(X), 0.1)
52  X_grid = X_grid.reshape((len(X_grid), 1))
53  plt.scatter(X, y, color = 'red')
54  plt.plot(X_grid, lin_reg_2.predict(
55      poly_reg.fit_transform(X_grid)), color = 'blue')
56  plt.title('Truth or Bluff (Polynomial Regression)')
57  plt.xlabel('Position level')
58  plt.ylabel('Salary')
59  plt.show()
```

92

46

## Polynomial Regression

- # Visualizing the Polynomial Regression results (for higher resolution and smoother curve)



93

## Polynomial Regression

- Prediction – Linear vs Polynomial

```
61 # Predicting a new result with Linear Regression
62 lin_reg.predict([[6.5]])
```

```
In [119]: lin_reg.predict([[6.5]])
Out[119]: array([330378.79])
```

```
64 # Predicting a new result with Polynomial Regression
65 lin_reg_2.predict(poly_reg.fit_transform([[6.5]]))
```

```
In [120]:
lin_reg_2.predict(poly_reg.fit_transform([[6.5]]))
Out[120]: array([158862.45])
```

94

94

## Decision Tree Regression



Prof.(Dr.) Dharmendra Bhatti

95

## Decision Tree Regression



- Decision tree learning is one of the predictive modelling approaches used in machine learning.

Prof.(Dr.) Dharmendra Bhatti

96

96

48

## Decision Tree Regression



- Decision trees are predictive models that use a set of binary rules to calculate a target value.
- Each individual tree has branches, nodes and leaves.

Prof.(Dr.) Dharmendra Bhatti

97

97

## Decision Tree Regression



- Tree models where the target variable can take a discrete set of values are called **classification trees** while if the target variable can take continuous values (typically real numbers) are called **regression trees**.

Prof.(Dr.) Dharmendra Bhatti

98

98

49

## Decision Tree Regression



### Decision Tree

#### Regression Tree

#### Classification Tree

Prof.(Dr.) Dharmendra Bhatti

99

99

## Decision Tree Regression



- Basic terminology in decision tree
  - **Root Node** represents entire population or dataset.
  - **Splitting** is a process of dividing a node into two or more sub-nodes.
  - When a sub-node splits into further sub-nodes, then it is called **decision node**.
  - Nodes, which do not split, is called **Leaf/Terminal node**.

Prof.(Dr.) Dharmendra Bhatti

100

100

50

## Decision Tree Regression



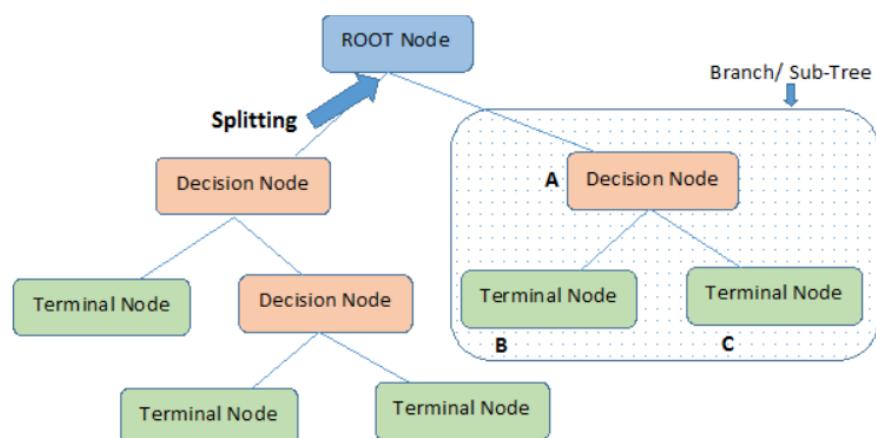
- Basic terminology in decision tree
  - When we remove sub-nodes of a decision node, this process is called **pruning**.
  - A sub section of entire tree is called **branch/subtree**.
  - A node, which is divided into sub-nodes is called **parent node** of sub-nodes whereas sub-nodes are the **child nodes** of parent node.

Prof.(Dr.) Dharmendra Bhatti

101

101

## Decision Tree Regression



Prof.(Dr.) Dharmendra Bhatti

102

102

## Decision Tree Regression



- A decision tree estimates by asking a series of questions, each question narrowing possible values until the model reaches to the final prediction.

Prof.(Dr.) Dharmendra Bhatti

103

103

## Decision Tree Regression



- The order of the question as well as their content are being determined by the model.
- The questions asked are all of True/False type.

Prof.(Dr.) Dharmendra Bhatti

104

104

## Decision Tree Regression



- During training, the model learns relationships between the data and the target variable.
- After the training phase, the decision tree produces a tree, calculating the best questions as well as their order to ask in order to make the most accurate estimates possible.

Prof.(Dr.) Dharmendra Bhatti

105

105

## Decision Tree Regression



- Decision tree regression normally use mean squared error (MSE) to decide to split a node in two or more sub-nodes.

Prof.(Dr.) Dharmendra Bhatti

106

106

## Decision Tree Regression



### • How to split?

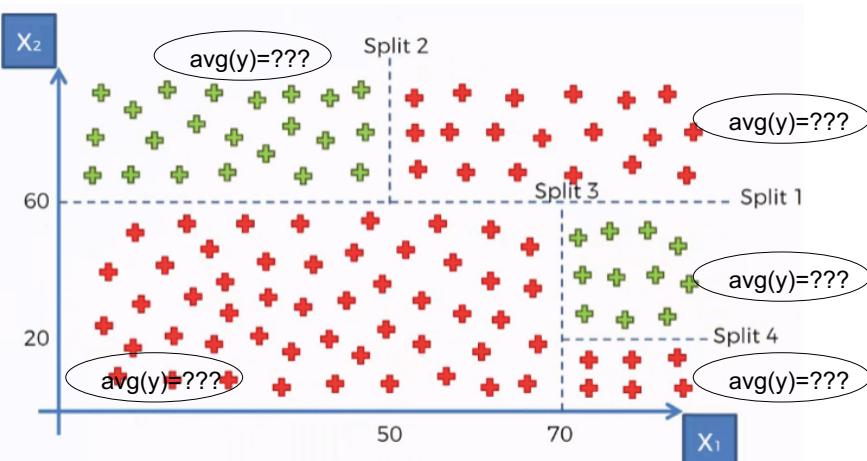
- Select a variable and the value to split such that the two groups are different from each other.
- For each variable, for each possible value of that variable see whether it is better (using mean square error).

Prof.(Dr.) Dharmendra Bhatti

107

107

## Decision Tree Regression

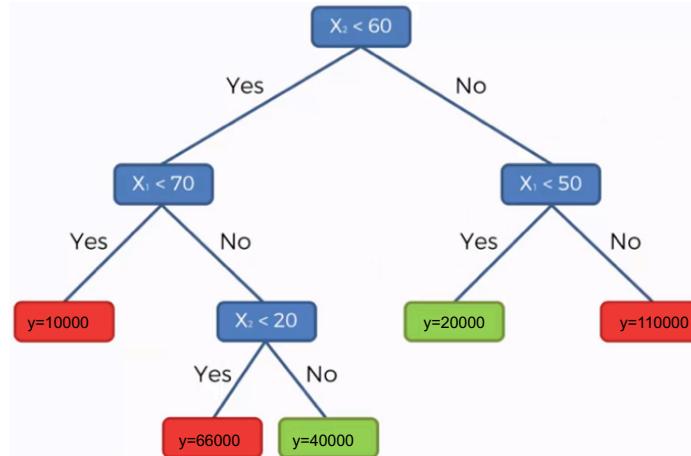


Prof.(Dr.) Dharmendra Bhatti

108

108

## Decision Tree Regression



Prof.(Dr.) Dharmendra Bhatti

109

109

## Decision Tree Regression

- Position\_Salaries.csv

	A	B	C
1	Position	Level	Salary
2	Business Analyst	1	45000
3	Junior Consultant	2	50000
4	Senior Consultant	3	60000
5	Manager	4	80000
6	Country Manager	5	110000
7	Region Manager	6	150000
8	Partner	7	200000
9	Senior Partner	8	300000
10	C-level	9	500000
11	CEO	10	1000000

110

110

## Decision Tree Regression



- Importing Libraries

```
9  # Decision Tree Regression
10
11 # Importing the libraries
12 import numpy as np
13 import matplotlib.pyplot as plt
14 import pandas as pd
```

Prof.(Dr.) Dharmendra Bhatti

111

111

## Decision Tree Regression



- Importing Dataset

```
16 # Importing the dataset
17 dataset = pd.read_csv('Position_Salaries.csv')
18 X = dataset.iloc[:, 1:-1].values
19 y = dataset.iloc[:, -1].values
```

Prof.(Dr.) Dharmendra Bhatti

112

112

## Decision Tree Regression



- Training the Decision Tree Regression model

```
21 # Training the Decision Tree Regression model on the whole dataset
22 from sklearn.tree import DecisionTreeRegressor
23 regressor = DecisionTreeRegressor(random_state = 0)
24 regressor.fit(X, y)
```

Prof.(Dr.) Dharmendra Bhatti

113

113

## Decision Tree Regression



- Predicting a new result

```
26 # Predicting a new result
27 regressor.predict([[6.5]])
```

Prof.(Dr.) Dharmendra Bhatti

114

114

## Decision Tree Regression

- Visualizing the Decision Tree Regression results

```
29 # Visualising the Decision Tree Regression results
30 plt.scatter(X, y, color = 'red')
31 plt.plot(X, regressor.predict(X), color = 'blue')
32 plt.title('Truth or Bluff (Decision Tree Regression)')
33 plt.xlabel('Position level')
34 plt.ylabel('Salary')
35 plt.show()
```

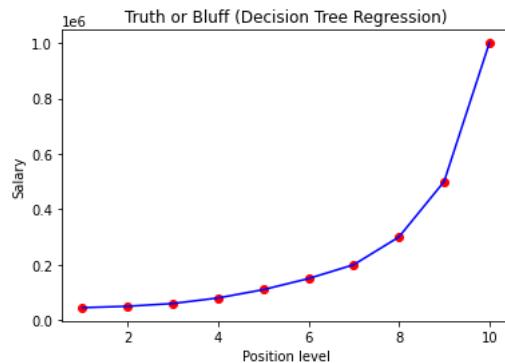
Prof.(Dr.) Dharmendra Bhatti

115

115

## Decision Tree Regression

- Visualizing the Decision Tree Regression results



Prof.(Dr.) Dharmendra Bhatti

116

116

## Decision Tree Regression



- BUT,

- Decision Tree Regression is non-linear and non-continuous

Prof.(Dr.) Dharmendra Bhatti

117

117

## Decision Tree Regression



- Visualizing the Decision Tree Regression results (high resolution)

```
37 # Visualising the Decision Tree Regression results (high resolution)
38 X_grid = np.arange(min(X), max(X), 0.1)
39 X_grid = X_grid.reshape((len(X_grid), 1))
40 plt.scatter(X, y, color = 'red')
41 plt.plot(X_grid, regressor.predict(X_grid), color = 'blue')
42 plt.title('Truth or Bluff (Decision Tree Regression)')
43 plt.xlabel('Position level')
44 plt.ylabel('Salary')
45 plt.show()
```

Prof.(Dr.) Dharmendra Bhatti

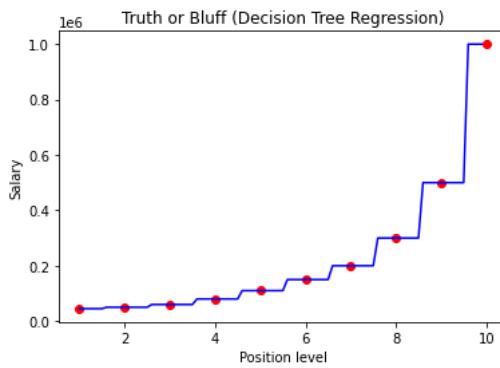
118

118

## Decision Tree Regression



- Visualizing the Decision Tree Regression results (high resolution)



Prof.(Dr.) Dharmendra Bhatti

119

119

## Decision Tree Regression



- Visualizing the Decision Tree Regression results (higher resolution)

```
47 # Visualising the Decision Tree Regression results (higher resolution)
48 X_grid = np.arange(min(X), max(X), 0.01)
49 X_grid = X_grid.reshape((len(X_grid), 1))
50 plt.scatter(X, y, color = 'red')
51 plt.plot(X_grid, regressor.predict(X_grid), color = 'blue')
52 plt.title('Truth or Bluff (Decision Tree Regression)')
53 plt.xlabel('Position level')
54 plt.ylabel('Salary')
55 plt.show()
```

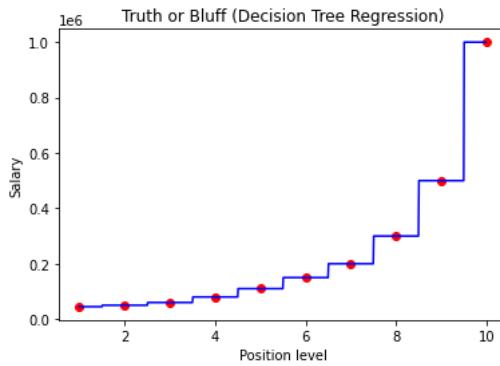
Prof.(Dr.) Dharmendra Bhatti

120

120

## Decision Tree Regression

- Visualizing the Decision Tree Regression results (higher resolution)



Prof.(Dr.) Dharmendra Bhatti

121

## Random Forest Regression

- Random Forest is an **ensemble technique** capable of performing both regression and classification tasks with the use of multiple decision trees.

Prof.(Dr.) Dharmendra Bhatti

122

122

## Random Forest Regression



- An **ensemble** method is a technique that combines the predictions from multiple machine learning algorithms together to make more accurate predictions than any individual model.

Prof.(Dr.) Dharmendra Bhatti

123

123

## Random Forest Regression



### Ensemble Learning

#### Boosting

#### Bootstrap Aggregation (Bagging)

Prof.(Dr.) Dharmendra Bhatti

124

124

## Random Forest Regression



- **Boosting** refers to a set of algorithms that utilize weighted averages to make weak learners into stronger learners.
- Boosting focuses on “teamwork”.
- Each model that runs, dictates what features the next model will focus on.

Prof.(Dr.) Dharmendra Bhatti

125

125

## Random Forest Regression



- **Bootstrap Aggregation (Bagging)** refers to random sampling with replacement.
- Bootstrap involves random sampling of small subset of data from the dataset.
- Each model run independently and then aggregates the outputs.

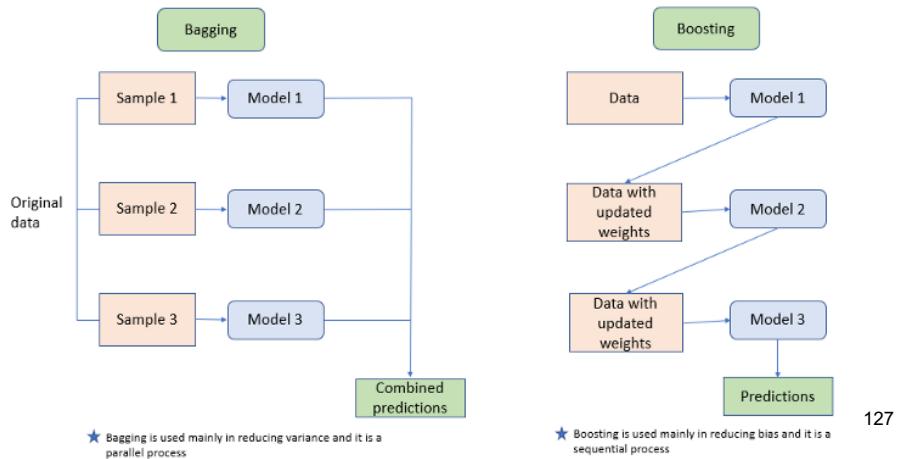
Prof.(Dr.) Dharmendra Bhatti

126

126

# Random Forest Regression

## • Bagging Vs Boosting



127

# Random Forest Regression

- Random forest is a **bagging** technique and **not a boosting** technique.

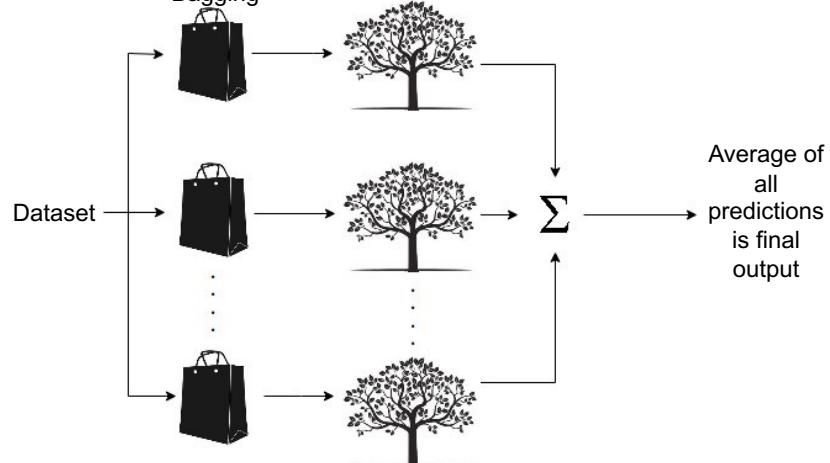
Prof.(Dr.) Dharmendra Bhatti

128

128

## Random Forest Regression

Bootstrap Aggregation /  
Bagging



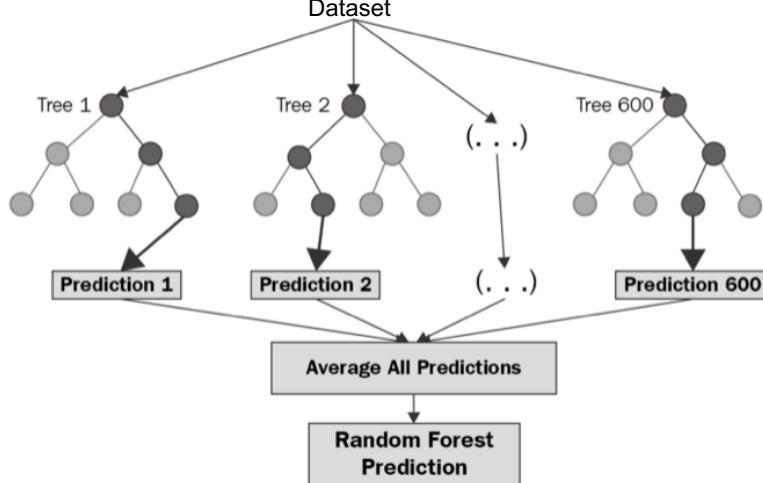
Prof.(Dr.) Dharmendra Bhatti

129

129

## Random Forest Regression

Dataset

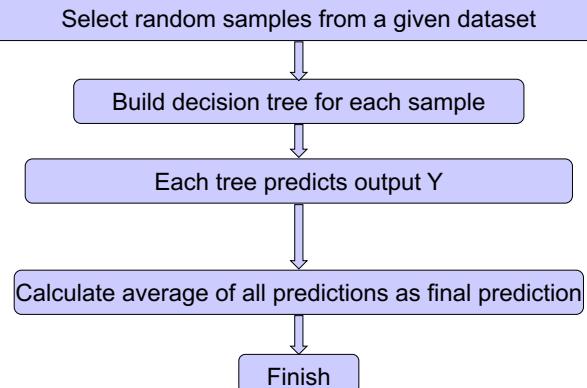


Prof.(Dr.) Dharmendra Bhatti

130

130

# Random Forest Regression



Prof.(Dr.) Dharmendra Bhatti

131

131

# Random Forest Regression



## • Position\_Salaries.csv

	A	B	C
1	Position	Level	Salary
2	Business Analyst	1	45000
3	Junior Consultant	2	50000
4	Senior Consultant	3	60000
5	Manager	4	80000
6	Country Manager	5	110000
7	Region Manager	6	150000
8	Partner	7	200000
9	Senior Partner	8	300000
10	C-level	9	500000
11	CEO	10	1000000

132

132

## Random Forest Regression



- Importing Libraries

```
9  # Random Forest Regression
10
11 # Importing the libraries
12 import numpy as np
13 import matplotlib.pyplot as plt
14 import pandas as pd
```

Prof.(Dr.) Dharmendra Bhatti

133

133

## Random Forest Regression



- Importing Dataset

```
16 # Importing the dataset
17 dataset = pd.read_csv('Position_Salaries.csv')
18 X = dataset.iloc[:, 1:-1].values
19 y = dataset.iloc[:, -1].values
```

Prof.(Dr.) Dharmendra Bhatti

134

134

## Random Forest Regression



- Training the Random Forest Regression model with 10, 100, 300 trees and predict

```
21 # Training the Random Forest Regression model on the whole dataset
22 from sklearn.ensemble import RandomForestRegressor
23
24 # Generate 10 trees for training
25 regressor = RandomForestRegressor(n_estimators = 10, random_state = 0)
26 regressor.fit(X, y)
27
28 # Predicting a new result with 10 trees
29 regressor.predict([[6.5]])
30
31 # Generate 100 trees for training
32 regressor = RandomForestRegressor(n_estimators = 100, random_state = 0)
33 regressor.fit(X, y)
34
35 # Predicting a new result with 100 trees
36 regressor.predict([[6.5]])
37
38 # Generate 300 trees for training
39 regressor = RandomForestRegressor(n_estimators = 300, random_state = 0)
40 regressor.fit(X, y)
41
42 # Predicting a new result with 300 trees
43 regressor.predict([[6.5]])
```

135

135

## Random Forest Regression



- Visualizing the Random Forest Regression results

```
45 # Visualising the Random Forest Regression results (higher resolution)
46 X_grid = np.arange(min(X), max(X), 0.01)
47 X_grid = X_grid.reshape((len(X_grid), 1))
48 plt.scatter(X, y, color = 'red')
49 plt.plot(X_grid, regressor.predict(X_grid), color = 'blue')
50 plt.title('Truth or Bluff (Random Forest Regression)')
51 plt.xlabel('Position level')
52 plt.ylabel('Salary')
53 plt.show()
```

Prof.(Dr.) Dharmendra Bhatti

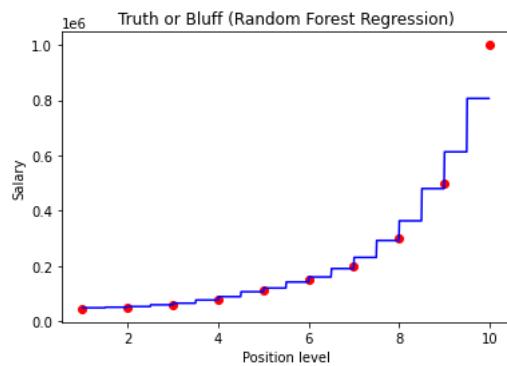
136

136

## Random Forest Regression



- Visualizing the Random Forest Regression results



Prof.(Dr.) Dharmendra Bhatti

137

137

Questions ???



Dharmendra Bhatti

138

138