

UNIT 3 - REGRESSION

Introduction

- ↳ Regression is Supervised Learning algo. (ie learns from labelled data) to make predictions.
- ↳ Focuses on predicting continuous outcomes based on relationship between independent & dependent (target) variables.

Independent → influence the dependent variable (size of house)

dependent → outcome that model aims to predict (house price, temp.)

$$y = f(x)$$

(GRE score) Dependent \longleftrightarrow independent (CGPA)

Types -

Linear Regression - models the relationship in straight line.

Simple Linear Regression - has only 1 independent variable & 1 dependent variable

Multiple Linear Regression - has multiple independent var & 1 dependent var

SLR

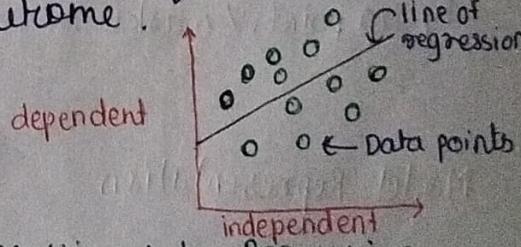
MLR

- | | |
|--|---|
| <p>① 1 independent var & 1 dependent var.</p> <p>② $y = \beta_0 + \beta_1 x + \epsilon$</p> <p>③ 1 independent var</p> <p>④ Less complex</p> <p>⑤ Analyse effect of single factor on outcome.</p> | <p>① many independent var & 1 dependent var.</p> <p>② $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \epsilon$</p> <p>③ 2 or more independent var.</p> <p>④ More complex</p> <p>⑤ Analyse effect of multiple factors on outcome.</p> |
|--|---|

Polynomial Regression

Logistic Regression

Decision Tree Regression, etc.



class 33 Linear Regression cont. as Univariate & Multivariate Regression

Univariate

+ Bivariate

Multivariate

① examines single variable at a time (dependent) independent).

① examines multiple variable at a time.

② less complex.

② More complex.

similar to simple & multiple regression

Univariate Regression

- ↳ goal is to find relationship betⁿ single feature & target variable
- ↳ most common technique used is Simple Linear Regression, that fits a straight line to the data.

Least Square Method.

- ↳ used to find the best-fitting line for set of data points.
- ↳ minimizes the difference betⁿ actual & predicted values / data.

Steps

- ① We have pairs of data points (x, y)
 $x \rightarrow$ independent
 $y \rightarrow$ dependent.
- ② Calculate difference betⁿ actual y and predicted y
- ③ Sum of all squares = $\sum (y_i - \bar{y}_i)^2$
 $y_i \rightarrow$ actual
 $\bar{y}_i \rightarrow$ predicted
- ④ Find the best fit line by slope & intercept calculation
- ⑤ Final equation in form $y = mx + c$
 $m \rightarrow$ slope
 $c \rightarrow$ y-intercept

Eg. X - hours studied
 Y - exam scores

\therefore for every additional hour studied score increases by 10 points / marks.

- ① calculate $\frac{\text{sum}}{n}$ of X & Y , XY & X^2
- ② find slope & intercept
- ③ $y = mx + c$ or $y = \beta_1 x + \beta_0$
eg $y = 50x + 10x$

$$\beta_1 = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2}$$

$$\beta_0 = \frac{\sum y - \beta_1 \sum x}{n}$$

Model Representation

Cost functions / Loss function

- ↳ calculates difference b/w predicted values & the actual values.
- ↳ minimizing cost function improves accuracy. ↓ cost fn ↑ accuracy

MSE (Mean Squared Error)

- ↳ average of squared differences b/w actual & predicted values.

$$MSE = \frac{1}{n} \sum (y_{\text{actual}} - y_{\text{predicted}})^2$$

↓ MSE = better fit Appl'n - evaluate model performance & make comparisons b/w different models

MAE (Mean Absolute Error)

- ↳ average magnitude of errors

↓ MAE = better fit

$$MAE = \frac{1}{n} \sum |y_{\text{actual}} - y_{\text{predicted}}|$$

MAE treat all errors equally Appl'n - model accuracy, esp when not sensitive to outliers

MAE - model accuracy, esp when outliers are present.

R² / Coefficient of Determination

- ↳ proportion of variance in dependent variable which can be explained by independent variable (ie can be predicted)

$$R^2 = 1 - \frac{\sum (y_{\text{actual}} - y_{\text{predicted}})^2}{\sum (y_{\text{actual}} - \bar{y})^2}$$

Ranges from 0 to 1

mean of actual values

↑ value = better fit

SSE (Sum of squared error)

$$SSE = \sum (y_{\text{actual}} - y_{\text{predicted}})^2$$

optimizing SLR with Gradient Descent Algorithm

↳ Iterative optimization algorithm used to minimize the cost/loss function

↳ updates parameters (β_0 & β_1) in direction that reduce loss

Steps -

① Initialize $\beta_0 = \beta_1 = 0$

② Define learning rate 'L' to a small value like 0.01

L → determines how big of a step we take during each iteration.

③ calculate gradient

$$D\beta_0 = -\frac{2}{n} \sum (y_{\text{actual}} - (\beta_1 x_i + \beta_0)) = -\frac{2}{n} \sum (y_{\text{actual}} - y_{\text{pred}})$$

$= -\frac{2}{n} \sum (y_{\text{actual}} - y_{\text{pred}})$

$$D\beta_1 = -\frac{2}{n} \sum (y_{\text{actual}} - (\beta_1 x_i + \beta_0)) x_i$$

$$= -\frac{2}{n} \sum x_i (y_{\text{actual}} - y_{\text{pred}})$$

④ Update β_0 & β_1 using gradients

$$\text{update } \beta_0 = \beta_0 - L D\beta_0$$

$$\text{update } \beta_1 = \beta_1 - L D\beta_1$$

⑤ Repeat 3 & 4 until cost function is small (ideally 0)

Multivariate Regression

↳ models the relationship between multiple independent & ~~multiple /~~ one dependent variables

$$\text{dependent } Y = \beta X + \epsilon \rightarrow \text{matrix of error terms}$$

↳ independent
matrix coefficient that represent relation b/w predictors & outcomes

$$\begin{bmatrix} Y_{11} & Y_{12} & \dots & Y_{1p} \\ Y_{21} & & \ddots & \\ \vdots & & & \\ Y_{n1} & & \dots & Y_{np} \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1q} \\ 1 & & \ddots & & \\ \vdots & & & & \\ 1 & x_{n1} & x_{n2} & \dots & x_{nq} \end{bmatrix} \begin{bmatrix} \beta_{01} & & \beta_{0p} \\ \beta_{11} & & \\ \vdots & & \\ \beta_{q1} & & \beta_{qp} \end{bmatrix} + \begin{bmatrix} \epsilon_{11} & \epsilon_{1p} \\ \epsilon_{21} & \dots \\ \vdots & \\ \epsilon_{n1} & \epsilon_{np} \end{bmatrix}$$

Advantages -

Captures complex Relationship - models how multiple predictors affect 1 outcome

Improve predictive power

Identify important factors

Control for other variable

Flexible

Disadvantages - complexity

Risk of overfitting

Data Requirements

Assumes linearity

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \epsilon$$

Polynomial Regression

↳ models relationship between independent & dependent variable as 'nth'- degree polynomial.

↳ it is extension of linear regression by introducing polynomial terms, allows model to fit non-linear relationships.

Model Representation

$$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \dots + \beta_n X^n + \epsilon$$

$$\begin{bmatrix} Y_1 \\ Y_2 \\ Y_3 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^m \\ 1 & x_2 & & & \\ 1 & x_3 & & & \\ \vdots & & & & \\ 1 & x_n & & & x_n^m \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_m \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

Generalization

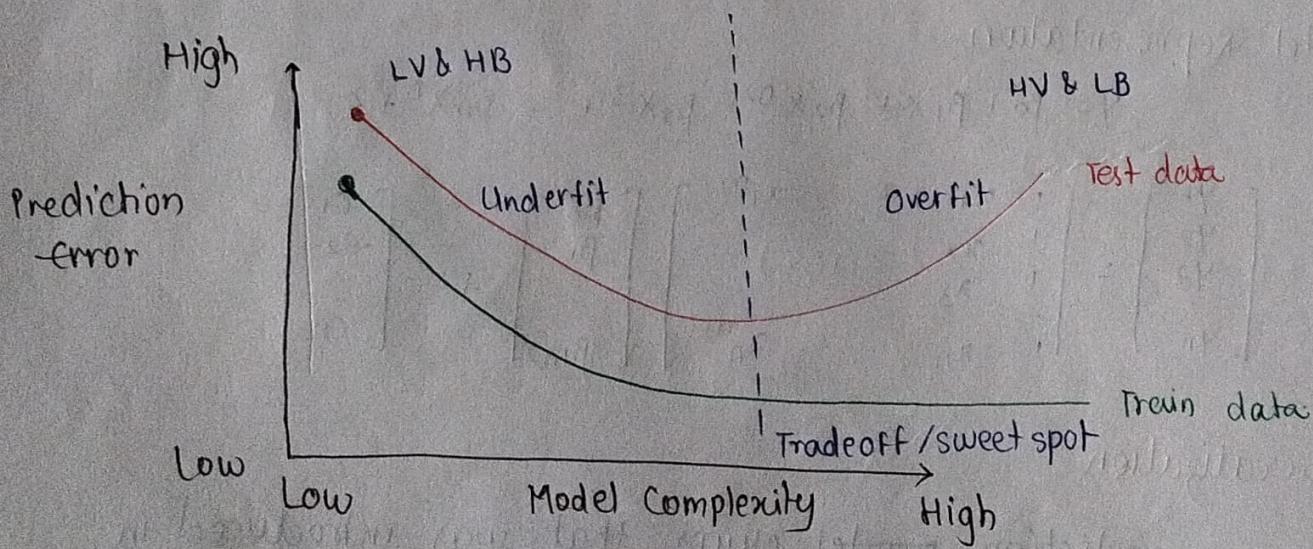
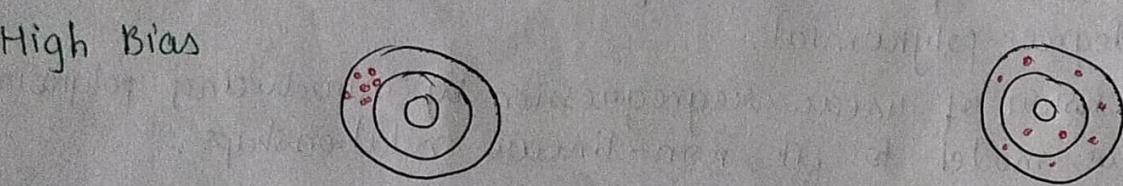
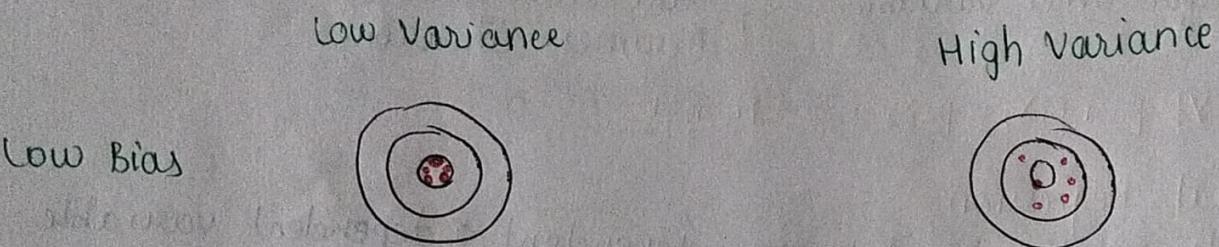
Bias - prediction model error that was introduced in model due to oversimplifying ML algorithm.
 gap b/w predicted & actual - represents difference between predicted and actual value

- High Bias = model is too simple and underfits the data

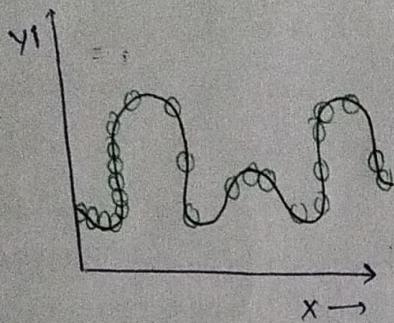
- Variance** - if model perform well with training dataset, but does not perform well with testing dataset, variance occurs
- it refers to models sensitivity to fluctuations in training data.
 - high variance = model is overfitting training data & not perform well on test/unseen data.
- how much scattered predicted values are wrt each other

Bias - Variance Dilemma / Tradeoff

- ↳ trade-off b/w 2 types of errors that affects the performance of predictive models.



Overfitting & Underfitting



error occurs when model is too complex & captures noise in training data

low training error

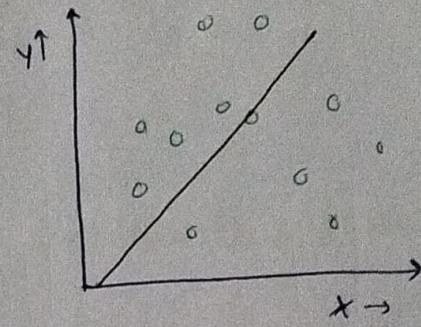
High testing error

Low bias

High variance

Too complex model complexity

Excellent performance on training data



error occurs when model is too simple to capture the underlying trend of data

High training error

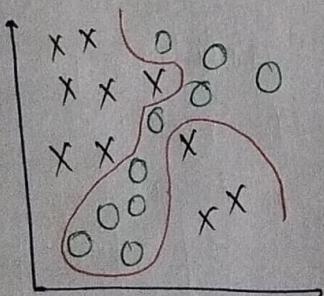
High testing error

High Bias

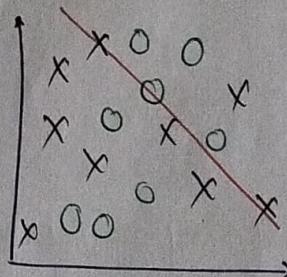
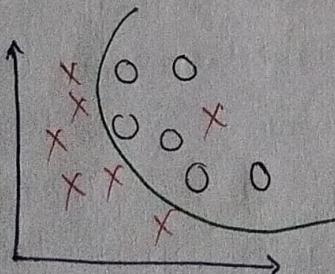
Low variance

Too simple model complexity

Low/poor performance on training data



Appropriate fit - not underfit, neither overfit



Ways to prevent overfitting

- ↳ early stopping
- ↳ train with more data
- ↳ feature selection
- ↳ cross validation
- ↳ data augmentation
- ↳ Regularization