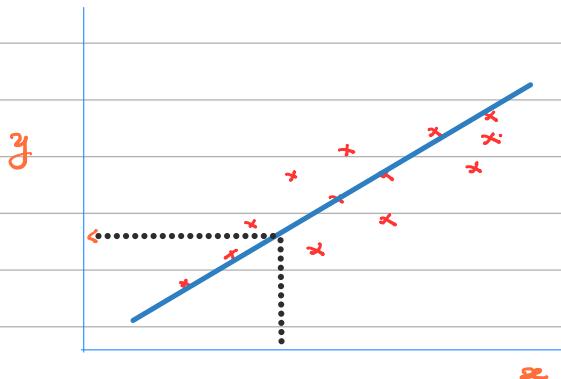


Simple Linear Regression

- Supervised Learning

- predict continuous numerical value
- finds best fitting straight line relation btw input and output
- used for a simple way to understand relationships and make prediction

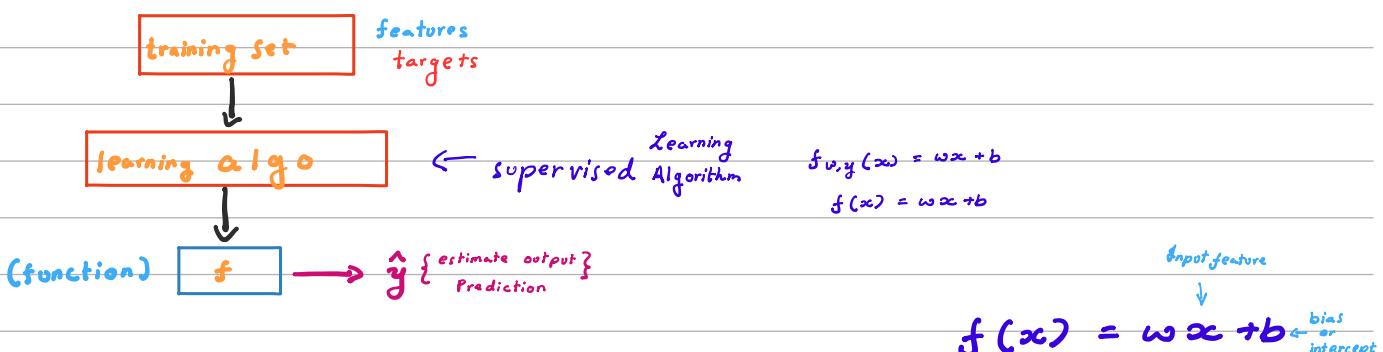


x = "input" value

y = "output" value

$(x^{(i)}, y^{(i)})$ = i^{th} training example.

$(x^{(i)}, y^{(i)}) = (812, 120)$



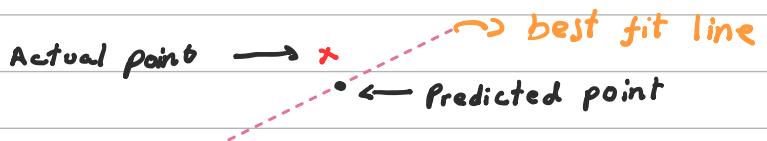
Inputs represents the features
Output represents the target

weight or coefficient

bias or intercept

The Linear regression algo learns from the data to generate best fit line

Residual Error: Difference btw actual point and predicted point

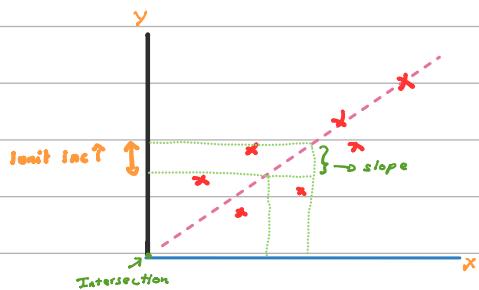


Aim: To find best fit line with minimal error

$$\hat{y} = m\alpha + c$$

m = slope or coefficient

c = Interception



Interception When the value of α is zero where is the best fit line intersecting y

$$h_{\theta}(\alpha) = \theta_0 + \theta_1 \alpha \rightarrow \text{coefficient}$$

\downarrow

intercept

Cost function formula

$$\sum_{i=1}^m \frac{1}{2m} (h_{\theta}(\alpha)^{(i)} - y^{(i)})^2$$

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(\alpha)^{(i)} - y^{(i)})^2$$

↓↓↓

↳ squared Error Function

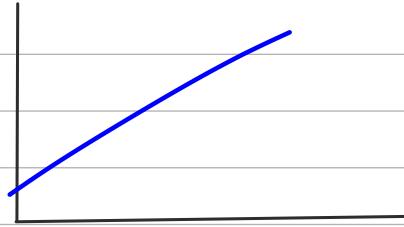
Model :

$$f_{w,b}(x) = w_0 + b$$

parameters:

$$w, b$$

cost function :



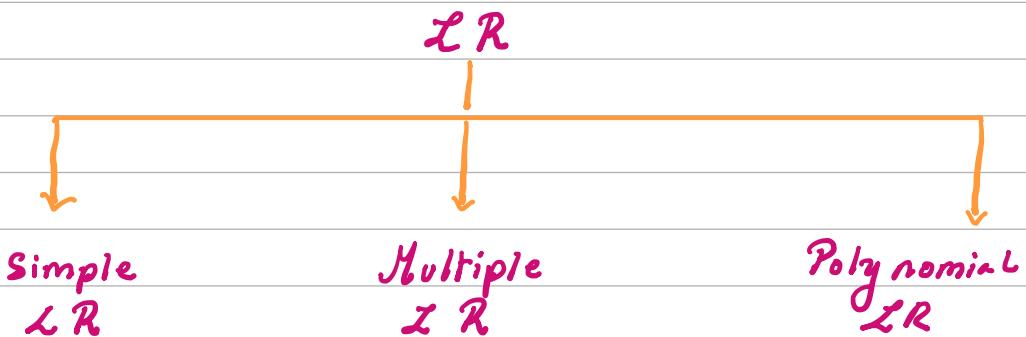
$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

goal :

Predicted Value ↑
actual value

minimize $J(w, b)$
 w, b

Types of Linear Regression



$x_1 | x_2 | x_3 | y$,
GPA | gen | grade | pa

$$y = m_0 + b$$

$$y = m_0 x_1 + n x_2 + b$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

$$y = \beta_0 + \beta_1 z_1 + \dots + \beta_n z_n$$

$$y = \beta_0 + \sum_{i=1}^n \beta_i x_i$$

Linear Regression

$$\text{OLS} \quad \beta = (x^T x)^{-1} x^T y$$

Gradient decent

Linear regression
can be found in
two diff ways

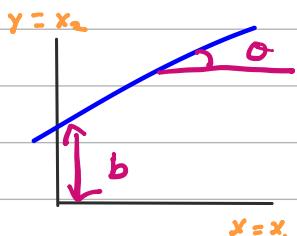
We use gradient decent most of the time because calculating inverse computationally cost a lot. And the time complexity of a inverse $O(n^3)$ is more time consuming when large sets of column data present

Hyperplane in N Dimentions

The line formed in 3D, 4D or ND is Known as hyperplane

Equation of Hyperplane

2D

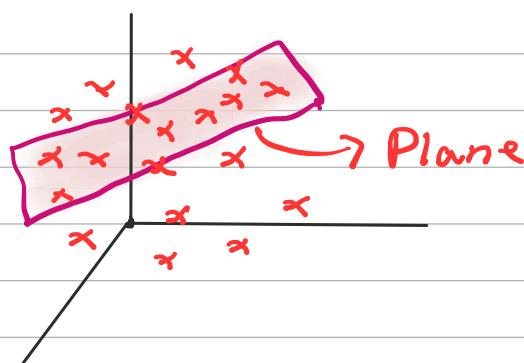


$$Y = mx + b$$

$$ax + by + c = 0$$

$$ax_1 + bx_2 + c = 0$$

$$w_1x_1 + w_2x_2 + w_0 = 0$$



Equation of Plane

$$w_1x_1 + w_2x_2 + \dots + w_nx_n + w_0 = 0$$

multiple Linear Regression (mathematical formulation)

$$\hat{Y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

$$\hat{Y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_{100} \end{bmatrix} = \begin{bmatrix} \beta_0 & \beta_1 x_{11} & \beta_2 x_{12} & \beta_3 x_{13} \\ \beta_0 & \beta_1 x_{21} & \beta_2 x_{22} & \beta_3 x_{23} \\ \vdots \\ \beta_0 & \beta_1 x_{1001} & \beta_2 x_{1002} & \beta_3 x_{1003} \end{bmatrix} \quad (100, 3)$$

$$\hat{Y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} \beta_0 & \beta_1 x_{11} & \beta_2 x_{12} & \beta_3 x_{13} & \dots & \beta_m x_{1m} \\ \beta_0 & \beta_1 x_{21} & \beta_2 x_{22} & \beta_3 x_{23} & \dots & \beta_m x_{2m} \\ \vdots \\ \beta_0 & \beta_1 x_{n1} & \beta_2 x_{n2} & \beta_3 x_{n3} & \dots & \beta_m x_{nm} \end{bmatrix}$$

$$\hat{Y} = \begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} & \dots & x_{1m} \\ 1 & x_{21} & x_{22} & x_{23} & \dots & x_{2m} \\ \vdots \\ 1 & x_{n1} & x_{n2} & x_{n3} & \dots & x_{nm} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_m \end{bmatrix}$$

$$\boxed{\hat{Y} = X \beta} \quad \text{--- ①}$$

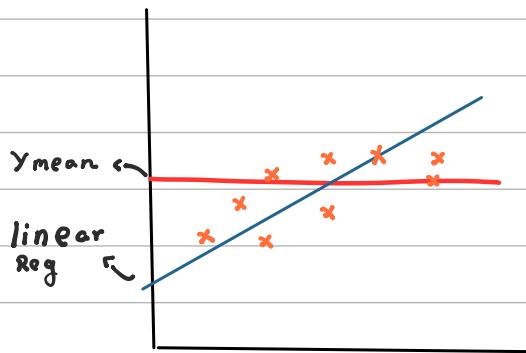
$$e = y - \hat{Y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} - \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix}$$

Regression Metric

To find the efficiency of Regression algorithm.

1. MAE - mean absolute error
 2. MSE - mean square Error
 3. RMSE - Root mean square error
 4. R² Score - coefficient of determination
 5. Adjusted R² score
- } ← loss function / error function
- } ← checking the accuracy
- } ← goodness of fit (literature)

R² Score



$$R^2 = 1 - \frac{SS_R}{SS_m} = \frac{\text{sum of Squared Error in Regression}}{\text{sum of Squared Error in mean line}}$$

$$R^2 = 1 - \frac{\left[\sum_{i=1}^n (y_i - \hat{y}_i)^2 \right]_{\text{Reg}}}{\left[\sum_{i=1}^n (y_i - \hat{y}_i)^2 \right]_m}$$

(for linear Reg)
(for mean line)

$R^2 = 0$ ← If the score is moving towards $R^2 = 0$ then its worst Reg line
 $R^2 = 1$ ← If the score is moving towards $R^2 = 1$ then its perfect lin line

$R^2 = -(\text{neg})$ ← It means that the lin line performs worst than the mean line

$R^2 = 0.80 \rightsquigarrow 80\% \text{ of Variance}$

Adjusted R₂ Score

In case of irrelevant column is added the R₂ score will either increase or remain constant, to avoid this we use Adj R₂ score

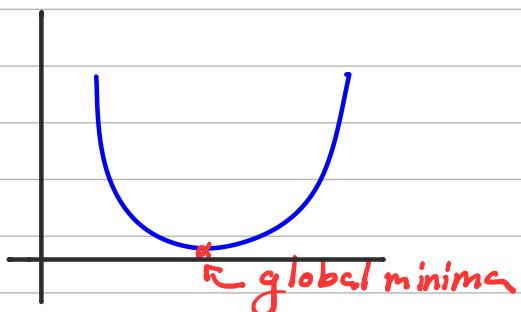
$$R^2_{adj} = 1 - \left[\frac{(1-R_2)(n-1)}{(n-1-k)} \right]$$

n → no of Rows
k → independent columns
R₂ → R₂ Score

when relatable column is added then R₂ increases

Mean Square Error

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$



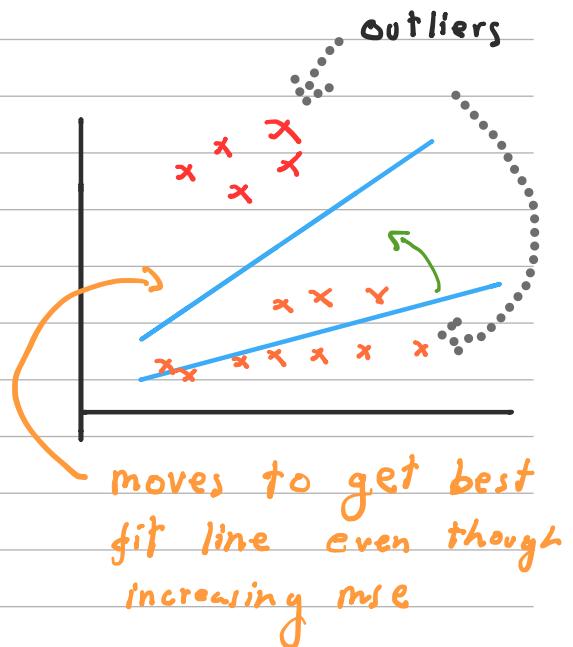
- ① mean Squared Error :

Advantage :

- ① Differentiable. (To find global minima)
- ② It has one global minima and local minima

Disadvantage :

- ① Non Robust to outliers
- ② It changes its unit



Mean absolute Error (MAE)

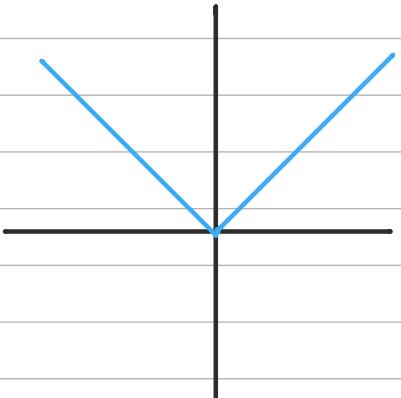
$$\text{cost function MAE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)$$

Advantage

1) It is robust to outliers

(The Best fit line if changed due to outliers might be minimal because its not squared)

2) It will be in same unit



Disadvantage

① Convergence usually takes more time

No gradient descent is formed in MAE

② optimisation is a complex process

RMSD

$$\text{RMSD} = \sqrt{\text{MSE}}$$

$$= \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}}$$

Advantage:

1) It will have same units

2) Differentiable

Note: The results by MAE & RMSD have different value and are not same.

Disadvantage

Not Robust to outliers

Polynomial Regression

$y = \beta_0 + \beta_1 x$ \rightarrow simple Linear Regression

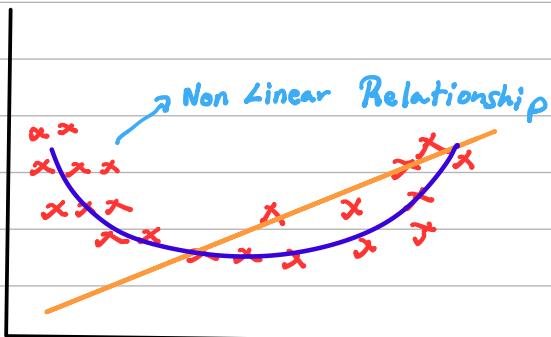
$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_n x_n$

\downarrow
multiple Linear
Regression

degree = 2 \rightarrow hyperparameter.

Simple Polynomial Regression : $y = \beta_0 + \beta_1 x + \beta_2 x^2$
(To extract non linear relation)

multiple Polynomial Regression : $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x_2 + \beta_4 x_2^2$



When we cannot use Linear Regression (it causes too much error) we use Polynomial Regression.

Polynomial Regression

\rightarrow Simple Polynomial Regression

\rightarrow Multiple Polynomial Regression

Simple Polynomial Regression

Polynomial degree = 0 $h_{\theta}(x) = \theta_0 x_1 = \theta_0 x_i \Rightarrow$ constant

Polynomial degree = 1

$h_{\theta}(x) = \theta_0 x_1 + \theta_1 x_i \Rightarrow$ Simple Linear Reg

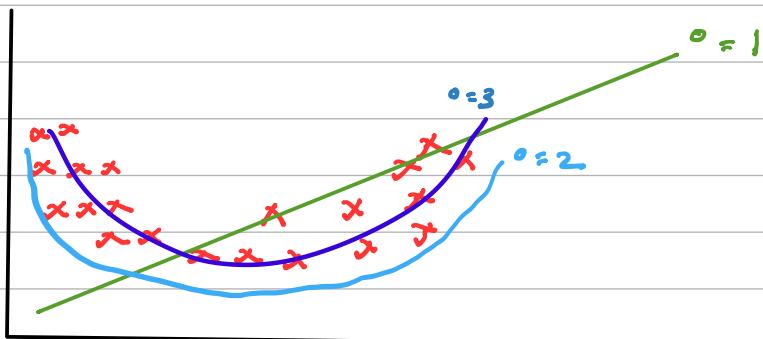
polynomial degree = 2

$$h_{\theta}(x) = \theta_0 x_i^0 + \theta_1 x_i^1 + \theta_2 x_i^2$$

:

polynomial degree = n

$$h_{\theta}(x) = \theta_0 x_i^0 + \theta_1 x_i^1 + \theta_2 x_i^2 + \dots + \theta_n x_i^n$$



multiple Polynomial Regression

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \quad \{ \text{multiple Linear Regression} \}$$

Polynomial degree = 2

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_1^2 + \theta_5 x_2^2 + \theta_6 x_3^2$$

Regularized Linear Regression Techniques.

① Ridge Regression → Reduces over fitting

Generalized Model

Low Bias

Aim to make generalized model

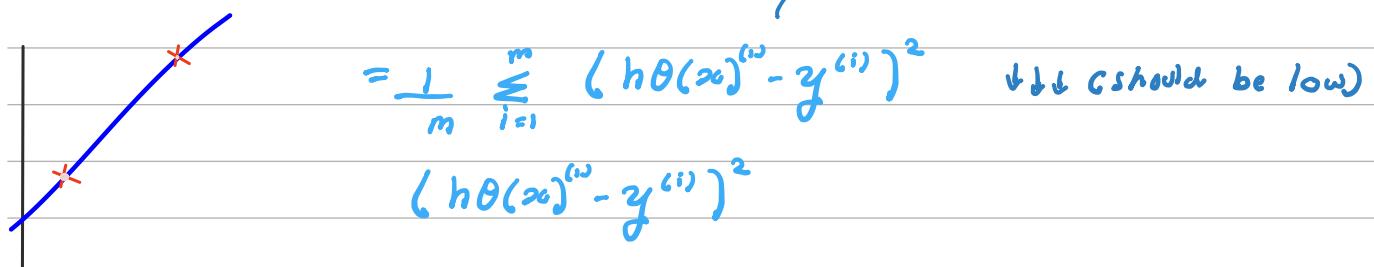
Low Variance

To prevent overfitting in Linear Regression we use Ridge regression

Linear Regression

cost function

→ Residual Error



= 0 ← If its zero then its overfitting

Ridge Regression (λ_2 Regularization)

$$\begin{aligned} &= (h\theta(x^{(i)}) - y^{(i)}) + \lambda (\text{slope})^2 \\ \lambda = 1 &= 0 + 1(2)^2 \\ &= 4 \end{aligned}$$

$\lambda \uparrow \uparrow$ Slope $\downarrow \downarrow$

Even though line passes through the points doesn't be 0 and doesn't cause overfitting

∴ Overfitting is removed

② Lasso Regression (λ , Regularization / norm) \rightarrow feature Selection

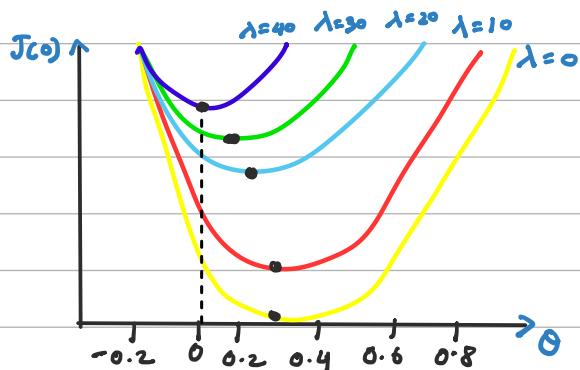
$$\text{Cost fn: } \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \boxed{\lambda + \sum_{i=1}^n |\text{slope}|}$$

$$(h_0(x)^{(i)} - y^{(i)})^2 + \lambda |\text{slope}|$$

1} Prevents overfitting

2} Feature Selection {feature Removal}

It uses feature Selection Process in which if the coefficient is less by using feature selection we eliminate the value by adjusting λ (hyperparameter) value.



Note

$\lambda \uparrow \uparrow \theta \downarrow \downarrow$

Till it becomes 0

③ Elastic Net Regression

Reducing Overfitting → Ridge
Feature Selection → lasso

$$\text{Cost function} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda_1 \sum_{i=1}^n (\text{slope})^2 + \lambda_2 \sum_{i=1}^n |\text{slope}|$$



Mean square Error + Reducing + Feature Selection
overfitting

λ_1, λ_2 = hyperparameter Tuning