

Data Science and Artificial Intelligence

Machine Learning



Regression

Lecture No. 04

By- SIDDHARTH SABHARWAL SIR

GATE WALLAH

Recap of Previous Lecture



Topic

Representation of data

Topic

2D data

Topic

Final expression of B.

Topic

Topic

Topics to be Covered



Topic

R^2 Factor

Topic

MSE

Topic

Gradient descent.

Topic

Topic





How the data is represented in matrix format

N points, D dimension

$$X_{\text{matrix}} = \begin{bmatrix} 1 & x_1^1 & x_1^2 & \dots & x_1^D \\ 1 & & & & \\ 1 & & & & \\ \vdots & & & & \end{bmatrix} \begin{array}{l} \leftarrow 1^{\text{st}} \text{ data} \\ \leftarrow 2^{\text{nd}} \text{ data} \\ \vdots \\ \end{array}$$

NX(D+1)

Each Row Show a data point



$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_D \end{bmatrix}_{D+1 \times 1}$$

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}_{N \times 1}$$



The final expression of Beta matrix..

$$\text{So } \beta = (X^T X)^{-1} (X^T Y)$$

If we have any data point

$$(1 \quad x_i^1 \quad x_i^2 \quad x_i^3 \quad \dots \quad x_i^D)$$

\Rightarrow So Predicted value $\Rightarrow (\beta_0 + \beta_1 x_i^1 + \beta_2 x_i^2 + \dots + \beta_D x_i^D)$

$$\Rightarrow \begin{pmatrix} X & \beta \end{pmatrix} \Rightarrow \begin{bmatrix} 1 & x_1^1 & x_1^2 & x_1^3 & \dots & x_1^D \\ 1 & x_2^1 & x_2^2 & x_2^3 & \dots & x_2^D \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \vdots \\ \beta_D \end{bmatrix} \Rightarrow \begin{bmatrix} \beta_0 + \beta_1 x_1^1 + \beta_2 x_1^2 + \dots + \beta_D x_1^D \\ \beta_0 + \beta_1 x_2^1 + \beta_2 x_2^2 + \dots + \beta_D x_2^D \\ \vdots \end{bmatrix}$$

$X\beta \Rightarrow$ we get $\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \\ \hat{y}_4 \\ \vdots \\ \hat{y}_N \end{bmatrix}$

So $Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$ \leftarrow actual values given in data.

$X\beta \Rightarrow \hat{Y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_N \end{bmatrix}$

\leftarrow Training data

• loss $f(x, \eta) \Rightarrow \sum_{i=1}^N (y_i - \hat{y}_i)^2$

\leftarrow Calculated on training data.



Linear Regression



How to represent the Loss function in the matrix format

Example $M = \begin{bmatrix} a \\ b \\ c \end{bmatrix}$

Now I want $(a^2 + b^2 + c^2)$

$$M^T M = \begin{bmatrix} a & b & c \end{bmatrix} \begin{bmatrix} a \\ b \\ c \end{bmatrix} = (a^2 + b^2 + c^2)$$

So if $X \Rightarrow$ Training data
 $Y \Rightarrow$ actual Y values of tr. data
 $\beta \Rightarrow$ linear reg. Parameters.

$$(X\beta) \Rightarrow \hat{Y}$$



Linear Regression



How to represent the Loss function in the matrix format

$$\text{So } Y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_N \end{bmatrix}$$

$$\hat{Y} = X\beta = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \\ \vdots \\ \hat{y}_N \end{bmatrix}$$

$$\text{So, } (Y - \hat{Y}) = \begin{bmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \\ \vdots \\ y_N - \hat{y}_N \end{bmatrix}$$

$$\begin{aligned} \text{Since loss fn} &\Rightarrow \sum_{i=1}^N (y_i - \hat{y}_i)^2 \Rightarrow (Y - \hat{Y})^T (Y - \hat{Y}) \\ &\Rightarrow (Y - X\beta)^T (Y - X\beta) \end{aligned}$$



Linear Regression



How to represent the Loss function in the matrix format

done

$$\mathcal{L} = (Y - X\beta)^T (Y - X\beta)$$



Linear Regression



How to represent the derivative of L by Beta in matrix format

↓
Recall the 2D data

$$\hat{y} = \beta_0 + \beta_1 x_i^1 + \beta_2 x_i^2$$

$$L = \sum_{i=1}^N \left(y_i - (\beta_0 + \beta_1 x_i^1 + \beta_2 x_i^2) \right)^2$$

$$\frac{\partial L}{\partial \beta_0} = -2 \left[\sum_{i=1}^N y_i - \beta_0 \sum 1 - \beta_1 \sum x_i^1 - \beta_2 \sum x_i^2 \right]$$

$$\frac{\partial L}{\partial \beta_1} = -2 \left[\sum_{i=1}^N x_i^1 y_i - \beta_0 \sum x_i^1 - \beta_1 \sum (x_i^1)^2 - \beta_2 \sum x_i^1 x_i^2 \right]$$

$$\frac{\partial L}{\partial \beta_2} = -2 \left[\sum_{i=1}^N x_i^2 y_i - \beta_0 \sum x_i^2 - \beta_1 \sum x_i^1 x_i^2 - \beta_2 \sum (x_i^2)^2 \right]$$



Linear Regression



How to represent the derivative of L by Beta in matrix format

$$\begin{bmatrix} \frac{\partial L}{\partial \beta_0} \\ \frac{\partial L}{\partial \beta_1} \\ \frac{\partial L}{\partial \beta_2} \end{bmatrix} = -2 \begin{bmatrix} \sum y_i - (\beta_0 \sum 1 + \beta_1 \sum x_{i1} + \beta_2 \sum x_{i2}) \\ \sum x_{i1} y_i - (\beta_0 \sum x_{i1} + \beta_1 \sum (x_{i1})^2 + \beta_2 \sum x_{i1} x_{i2}) \\ \sum x_{i2} y_i - (\beta_0 \sum x_{i2} + \beta_1 \sum (x_{i1}) x_{i2} + \beta_2 \sum (x_{i2})^2) \end{bmatrix}$$
$$= -2 \begin{bmatrix} \begin{bmatrix} X^T Y \\ \vdots \end{bmatrix} - \begin{bmatrix} (X^T X) \beta \end{bmatrix} \end{bmatrix}$$



Linear Regression



How to represent the derivative of L by Beta in matrix format

- So $L = (Y - \hat{Y})^T (Y - \hat{Y})$

- $\frac{\partial L}{\partial \beta} = \begin{bmatrix} \frac{\partial L}{\partial \beta_0} \\ \frac{\partial L}{\partial \beta_1} \\ \frac{\partial L}{\partial \beta_2} \\ \vdots \end{bmatrix} = -2 \begin{bmatrix} (X^T Y) - (X^T X) \beta \\ \dots \\ \dots \end{bmatrix}$



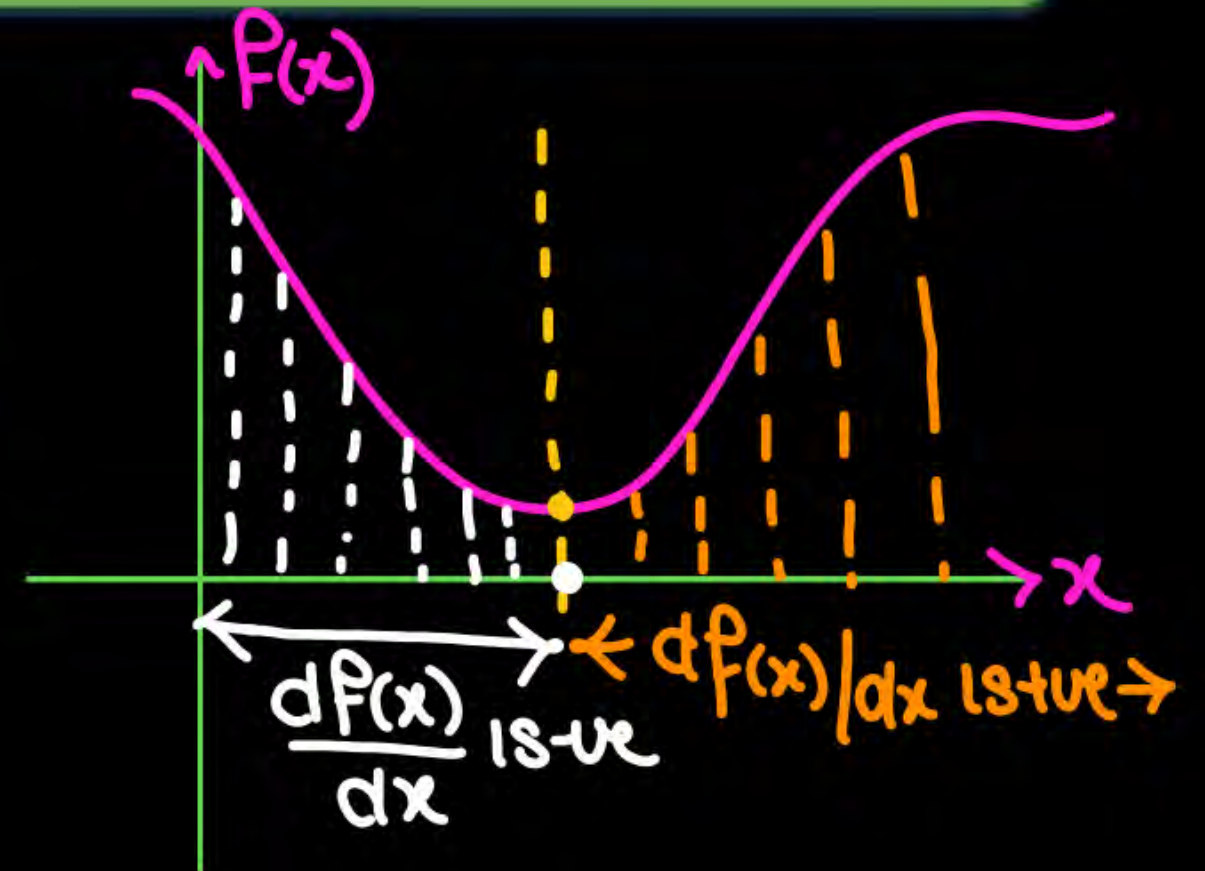
Linear Regression



What is gradient descent method

Concept \Rightarrow

- When any function is decaying then at that location derivative of $f(x) \Rightarrow -ve$
- When any function is Rising then at that location derivative of $f(x) \Rightarrow +ve$





Linear Regression



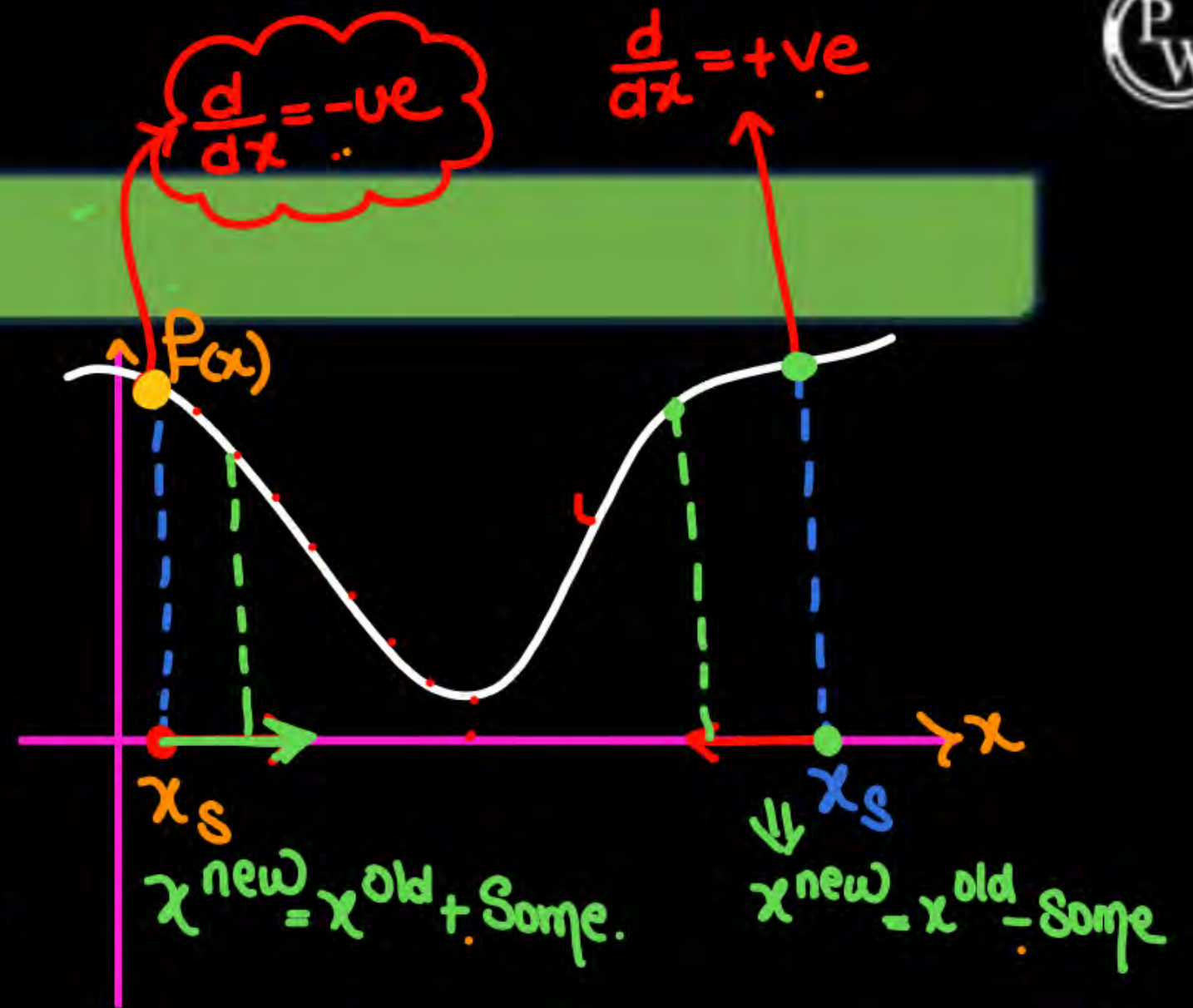
What is gradient descent method

* my objective is to reach minima location

we have a method
 $\frac{dP(x)}{dx} = 0$

So we use gradient descent \Rightarrow
 \Rightarrow Start with any Random value of x
 \Rightarrow then we move on to next value of x

$$x^{\text{new}} = x^{\text{old}} - \eta \left. \frac{dP(x)}{dx} \right|_{x^{\text{old}}}$$





Linear Regression



What is gradient descent method

So gradient descent method is as follows

So we have to find the min location of any $f(x)$

So Step 1. Start with any x value

Step 2. find new $x \Rightarrow x^{\text{new}} = \left(x^{\text{old}} - \eta \cdot \frac{df(x)}{dx} \Big|_{x^{\text{old}}} \right)$

\Rightarrow This is an iterative method
to reach to Result

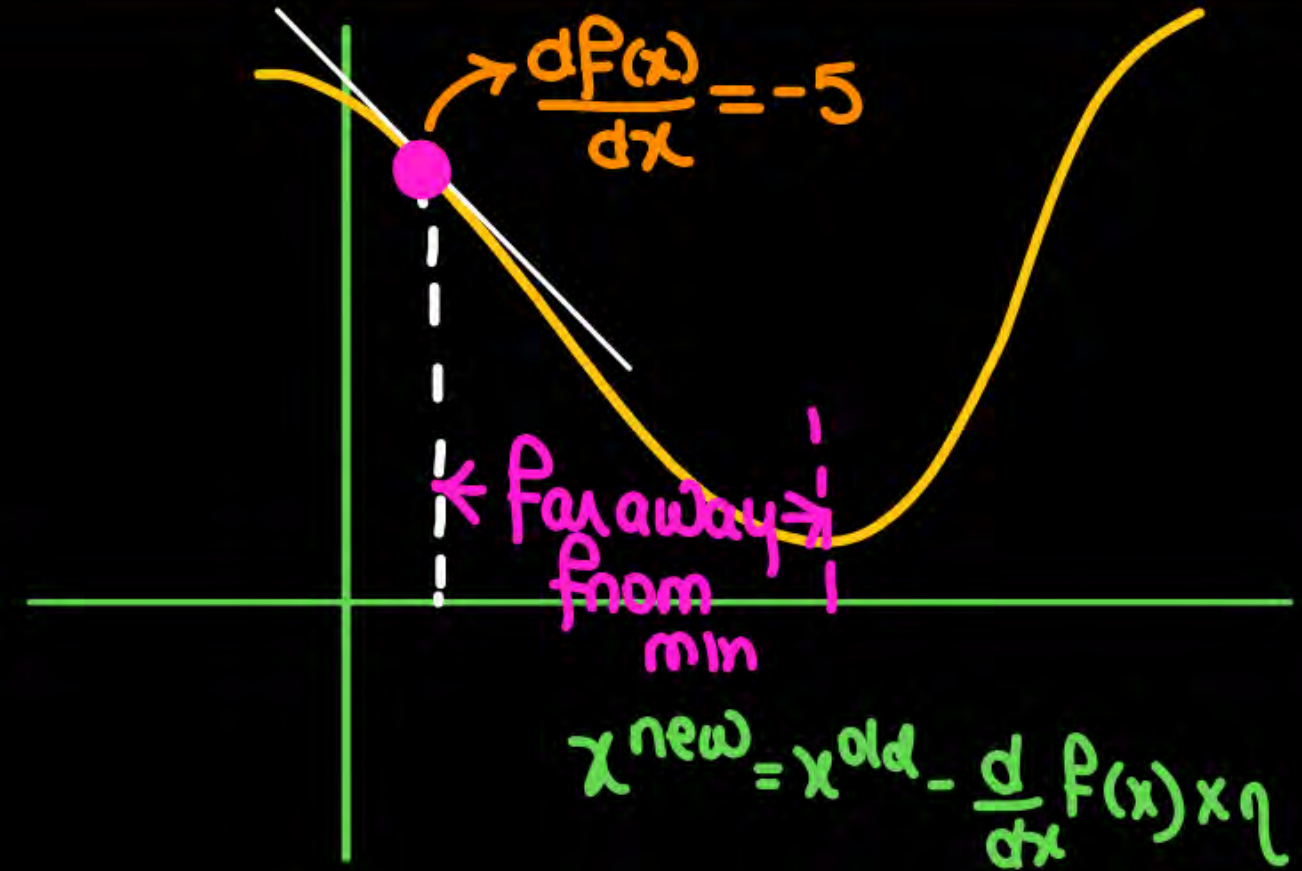
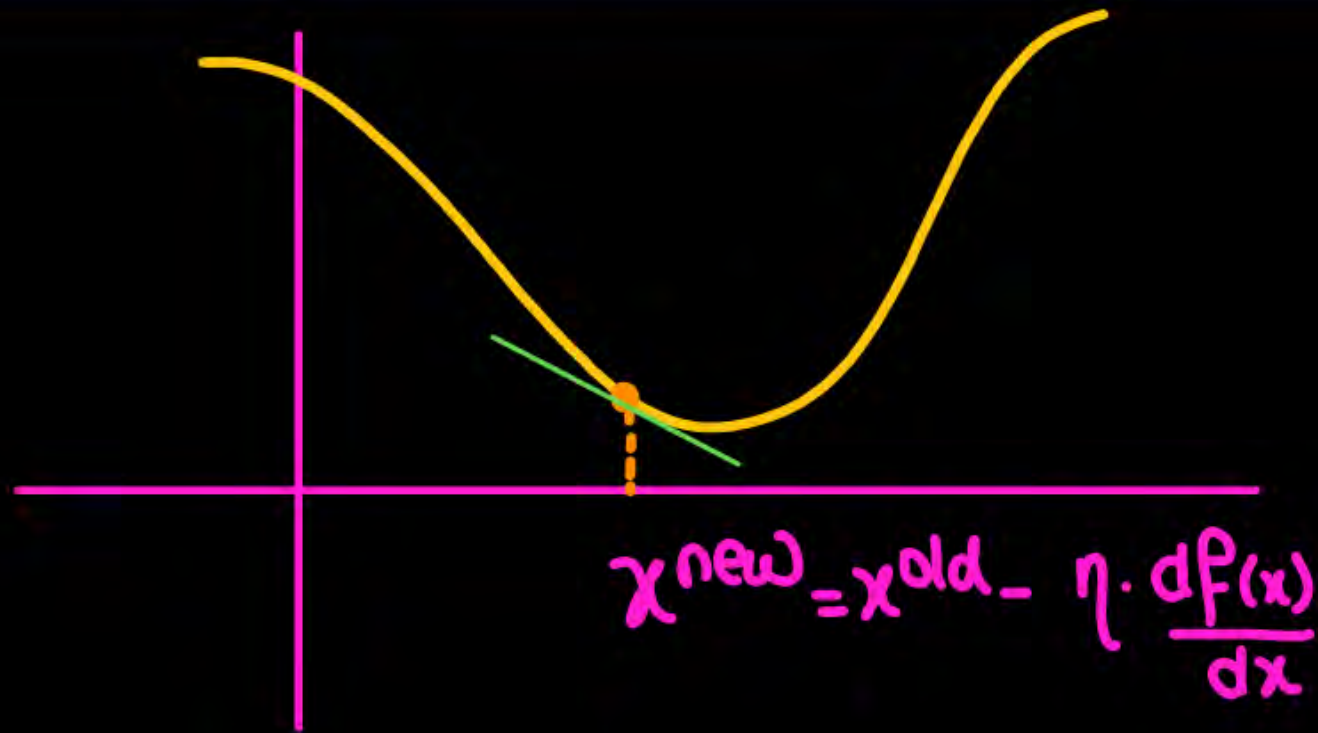
η : Learning Rate



Linear Regression



What is gradient descent method





Linear Regression



What is gradient descent method

So in linear regression \Rightarrow

So we have to minimize

$$L = \sum_{i=1}^N (y_i - \hat{y}_i)^2 = \sum_{i=1}^N \left[y_i - (\beta_0 + \beta_1 x_i^1 + \dots + \beta_D x_i^D) \right]^2$$

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_D \end{bmatrix}$$

So gradient descent method

So we start with Random β

$$\beta^{\text{new}} = \beta^{\text{old}} - \eta \left[\frac{\partial L}{\partial \beta} \right] \Big|_{\text{at } \beta^{\text{old}}}$$

matrix matrix matrix

$$\frac{\partial L}{\partial \beta} \Rightarrow -2 \left[X^T Y - (X^T X) \beta \right] \Big|_{\beta^{\text{old}}}$$

#Q. If $g(x, y) = x^2 + y^2 - 4x$, find the gradient vector $\nabla g(1, 2)$

example

Normal maths
function

$x=1$
 $y=2$

$$\text{So } \nabla g = \begin{bmatrix} \partial g / \partial x \\ \partial g / \partial y \end{bmatrix} = \begin{bmatrix} 2x - 4 \\ 2y \end{bmatrix} \quad @x=1, y=2$$

$$= \begin{bmatrix} -2 \\ 4 \end{bmatrix}$$

#Q. Let f be the function of two variables given by

$$f(x, y) = xy(x + y) = x^2y + xy^2$$

(a) Calculate the gradient vector ∇f and evaluate at the point $(1, 2)$.

H.W.

- Why Gradient descent \Rightarrow

$$\beta = (X^T X)^{-1} (X^T Y)$$

\rightarrow we know β solution for min. loss function

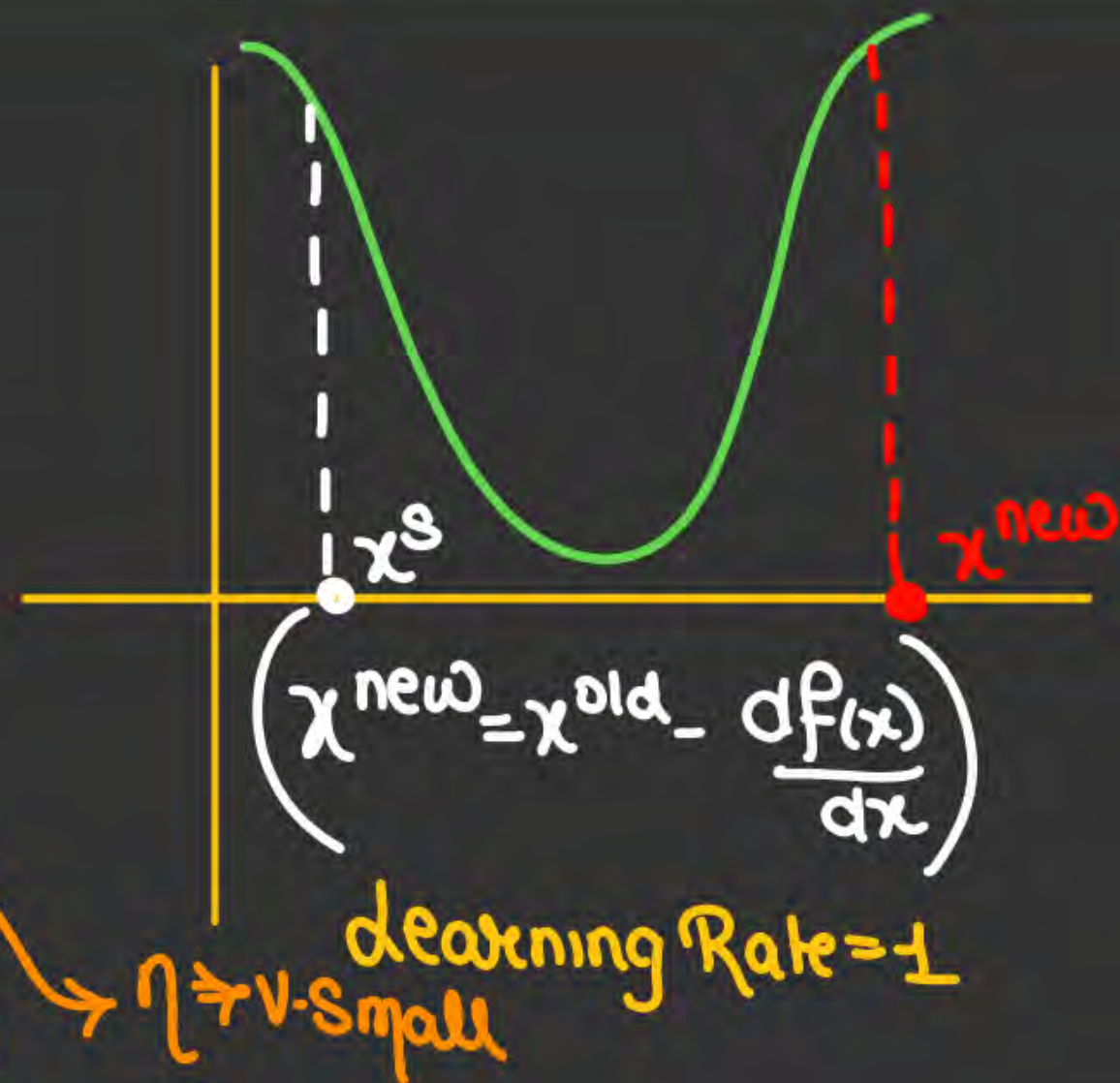
\rightarrow Problem in above eq. is that the calculation of $(X^T X)^{-1}$ is v.v. difficult bcoz data is huge.

\rightarrow that's why we use Gradient descent.

Significance of Learning Rate

→ Since we don't want that $\frac{df(x)}{dx}$ completely control the movement of x .

→ Thus we multiply $\frac{df(x)}{dx}$ by learning Rate so that movement from $x^{old} \rightarrow x^{new}$ is v. slow and we don't cross minima loc.



- if learning Rate is v. large \Rightarrow we will never reach min loc.

#Q. Let's consider regression in one dimension, so our inputs $x^{(i)}$ and outputs $y^{(i)}$ are in \mathbb{R} .

- (a) (4 points) Linny uses regular linear regression. Given the following dataset, (x, y) P.W.
- $$D = \{((1), 1), ((2), 2), ((3), 4), ((3), 2)\}$$

What value of θ and θ_0 optimize the mean squared error of hypotheses of the form $h(x; \theta, \theta_0) = \theta_1 x + \theta_0$?

find θ_1, θ_0

#Q. Suppose we have data about 5 people shown below.

Name	Level	Trials	Phase
Megda	1	10	1
Valerie	5	20	-1
Kumar	2	15	1
Octavia	6	30	1
Dorete	6	5	-1

Handwritten orange circle containing "P.W."

- (a) Suppose we want to model the level of each person, and use the following constant model:

$f_{\theta}(x) = \theta_1$. What is $\hat{\theta}_1$, the value that minimizes the average L_2 loss?

#Q. Consider a one-dimensional regression problem with training data $\{x_i, y_i\}$. We seek to fit a linear model with no bias term:

(a) Assume a squared loss $\frac{1}{2} \sum_{i=1}^N (y_i - \hat{y}_i)^2$ and solve for the optimal value of ω^* .

fp. ω

$\hat{y}_i = \omega x$
"find ω "

#Q. Consider the following 4 training examples:

X	Y
-1	0.0319
0	0.8692
1	1.9566
2	3.0343

$\hat{f}(x)$

We want to learn a function $f(x) = ax + b$ which is parametrized by (a, b) . Using squared error as the loss function, which of the following parameters would you use to model this function.

(a) (1, 1)

(b) (1, 2)

(c) (2, 1)

(d) (2, 2)

#Q. The linear regression model $y = a_0 + a_1x_1 + a_2x_2 + \dots + a_px_p$ is to be fitted to a set of N training data points having p attributes each. Let X be $N \times (p + 1)$ vectors of input values (augmented by 1's), Y be $N \times 1$ vector of target values, and θ be $(p + 1) \times 1$ vector of parameter values ($a_0, a_1, a_2, \dots, a_p$). If the sum squared error is minimized for obtaining the optimal regression model, which of the following equation holds?

(d)

(a) $X^T X = X Y$

(b) $X \theta = X^T Y$

(c) $X^T X \theta = Y$

(d) $X^T X \theta = X^T Y$



Linear Regression



Considering data of P Dimensions

R-squared in Regression Analysis in Machine Learning

$$R^2 = 1 - \frac{RSS}{TSS}$$

R^2 = coefficient of determination

RSS = sum of squares of residuals

TSS = total sum of squares



Linear Regression



Considering data of P Dimensions

R-squared in Regression Analysis in Machine Learning

$$RSS = \sum_{i=1}^n (y_i - f(x_i))^2$$

RSS = residual sum of squares

y_i = i^{th} value of the variable to be predicted

$f(x_i)$ = predicted value of y_i

n = upper limit of summation



Considering data of P Dimensions

R-squared in Regression Analysis in Machine Learning

$$TSS = \sum_{i=1}^n (y_i - \bar{y})^2$$

TSS = total sum of squares

n = number of observations

y_i = value in a sample

\bar{y} = mean value of a sample



Linear Regression



Considering data of P Dimensions

R-squared in Regression Analysis in Machine Learning

- ❖ The most important thing we do after making any model is evaluating the model.
- ❖ R-squared is a statistical measure that represents the goodness of fit of a regression model.
- ❖ The value of R-square lies between 0 to 1.
- ❖ Where we get R-square equals 1 when the model perfectly fits the data and there is no difference between the predicted value and actual value.
- ❖ However, we get R-square equals 0 when the model does not predict any variability in the model.



Considering data of P Dimensions

R-squared in Regression Analysis in Machine Learning

- ❖ R-Squared (R^2 or the coefficient of determination) is a statistical measure in a regression model that determines the proportion of variance in the dependent variable that can be explained by the independent variable.
- ❖ The most common interpretation of r-squared is how well the regression model explains observed data. For example, an r-squared of 60% reveals that 60% of the variability observed in the target variable is explained by the regression model.



Linear Regression



Considering data of P Dimensions

R-squared in Regression Analysis in Machine Learning

- ❖ The goodness of fit of regression models can be analyzed on the basis of the R-square method. The more the value of the r-square near 1, the better the model is.
- ❖ Note: The value of R-square can also be negative when the model fitted is worse than the average fitted model. .



Considering data of P Dimensions

Adjusted R - Squares

- ❖ Adjusted R-Squared is an updated version of R-squared which takes account of the number of independent variables while calculating R-squared.
- ❖ n is the total number of observations in the data
- ❖ k is the number of independent variables (predictors) in the regression model

$$Adjusted R^2 = 1 - \frac{(1-R^2) \cdot (n-1)}{n-k-1}$$



Linear Regression



Considering data of P Dimensions

Lets solve a question

Question 2: Given a simple linear regression model with an R-squared value of 0.64, what percentage of the variation in the dependent variable is explained by the predictor variable?



Linear Regression



Considering data of P Dimensions

Lets solve a question

Question 6: In a simple linear regression model, if the coefficient of determination (R-squared) is 0.81 and the total sum of squares (SST) is 400, what is the sum of squared errors (SSE)?

- a)76
- b)77
- c)54
- d)33



Linear Regression



What is Mean Square Error

4) Start with the initial guess of $[w_1, w_2] = [5, 5]$. Take the value of learning rate = 0.3. The value of w_1 after 2 iterations of gradient descent will be _____.

$$J(w) = w_1^2 + w_2^2 - 6w_1 + 8w_2 - 9$$

Solve . $\begin{pmatrix} w_1 \\ w_2 \end{pmatrix} = \begin{pmatrix} 5 \\ 5 \end{pmatrix}$

Step 1 \Rightarrow Start with initial value

Step 2 $\Rightarrow w^{\text{new}} = w^{\text{old}} - \alpha \left(\frac{\partial J}{\partial w} \right)_{w^{\text{old}}}$

$$\Rightarrow \frac{\partial J}{\partial w} = \begin{bmatrix} \frac{\partial J}{\partial w_1} \\ \frac{\partial J}{\partial w_2} \end{bmatrix} = \begin{bmatrix} 2w_1 - 6 \\ 2w_2 + 8 \end{bmatrix}$$

$$= \begin{bmatrix} 5 \\ 5 \end{bmatrix} - 0.3 \begin{pmatrix} 4 \\ 18 \end{pmatrix} \Rightarrow \begin{bmatrix} 3.8 \\ -0.4 \end{bmatrix}$$

1st iteration

2nd Iteration \Rightarrow

$$w_{\text{new}} = w_{\text{old}} - 0.3 \left(\frac{\partial J}{\partial w} \right) \Big|_{@w_{\text{old}}}$$

$$= \underbrace{\begin{bmatrix} 3.8 \\ -0.4 \end{bmatrix}} - 0.3 \begin{bmatrix} 1.6 \\ 7.2 \end{bmatrix}$$

$$= \begin{bmatrix} 3.32 \\ -2.56 \end{bmatrix}$$



Linear Regression



What is Mean Square Error

Consider the function $J(w) = w_1^2 + w_2^2 - 6w_1 + 8w_2 - 9$. Answer questions (1-6):

1) The theoretical value of $\min(J(w))$ is _____.

So find w_1, w_2 to $\min(J(w))$

$$\frac{\partial J}{\partial w_1} = 2w_1 - 6 = 0, w_1 = 3$$

$$\frac{\partial J}{\partial w_2} = 2w_2 + 8 = 0, w_2 = -4$$

Value of w_1, w_2

$$\text{min Value of } J(w) \Rightarrow -34$$



2 mins Summary



Topic

Topic

Topic

Topic

Topic

THANK - YOU