

Day 2

June 16, 2024

Linear Regression:

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)}))^2 \quad (\text{MSE Loss})$$

Gradient Descent (Generalized)

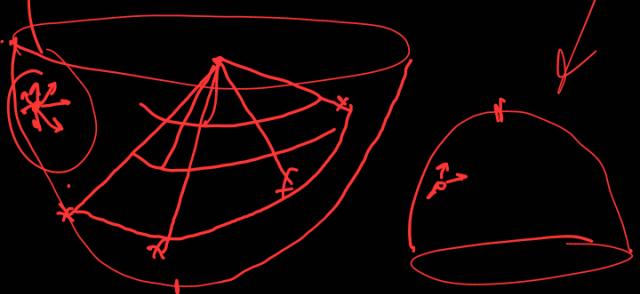
$$\theta_j^{(t+1)} \leftarrow \theta_j^{(t)} - \eta \cdot \nabla_{\theta_j} J(\theta)$$

$$\nabla_{\theta_j} J(\theta) = -\frac{1}{m} \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

Multiple Directions

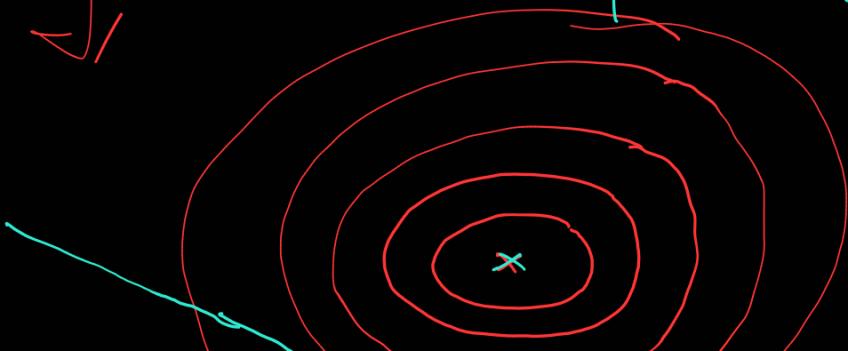
Contour Plots:

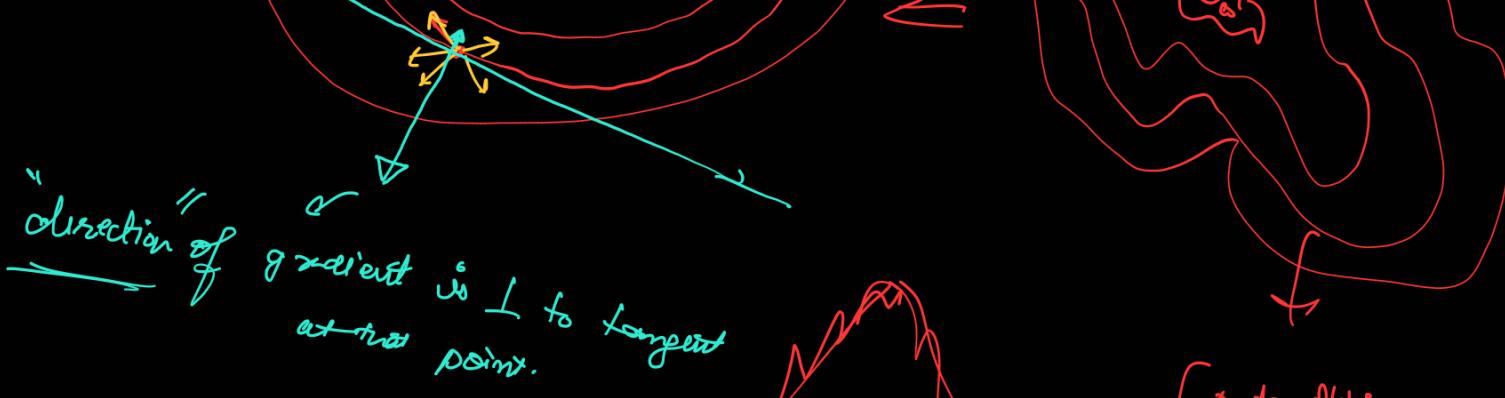
$\nabla J(\theta)$ → In all those directions, If I move opposite to direction of Gradient at that point, then I will be on the fastest route to minima / maxima.



Contour line: Joining points of same distance from a ref point.

Depending on the curve.





Contour plots



Paper ← Katori^o

$X \mapsto$ Gaussian Surface



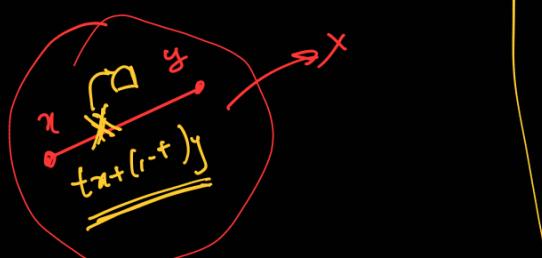
Convex Sets:

Any set $X \subseteq \mathbb{R}^d$ is convex if,

$t\mathbf{x} + (1-t)\mathbf{y} \in X, \mathbf{x}, \mathbf{y} \in X, t \in [0, 1]$

Relate to Section
formula (Interior)

$$\int \int \int \dots \int$$

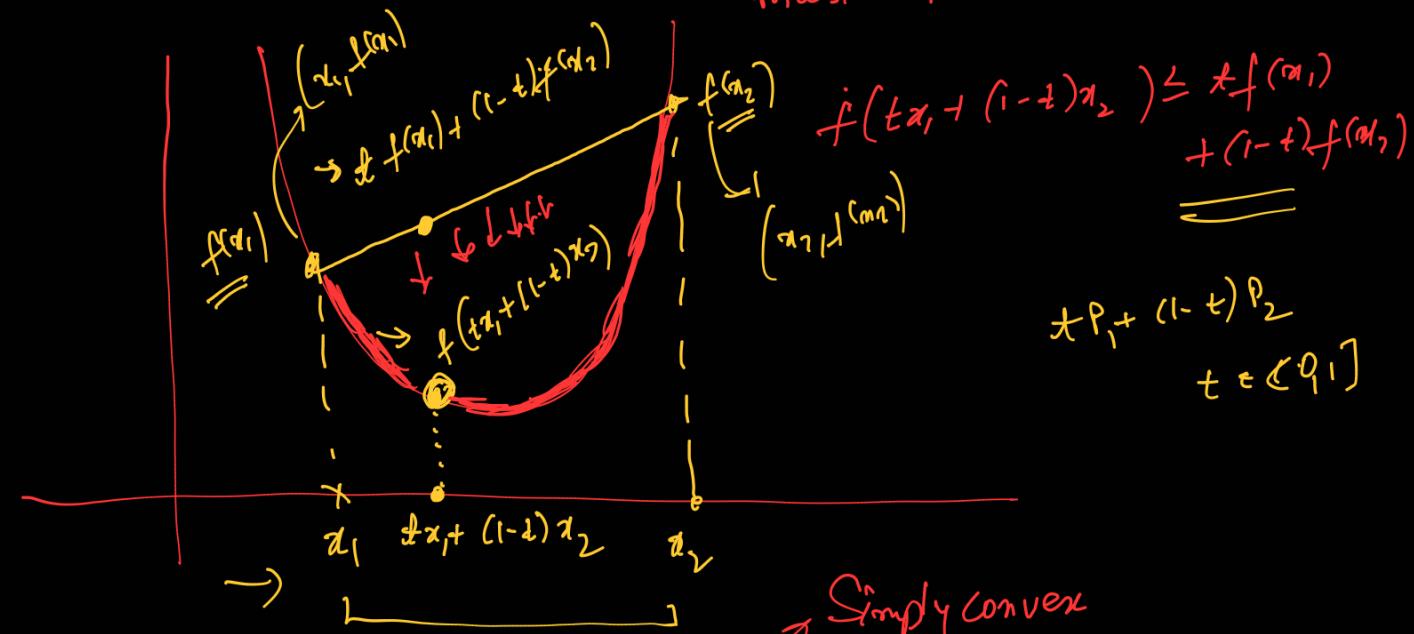


Convex functions: A function f is convex if: \rightarrow

$$f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2)$$

$\forall x_1, x_2 \in \text{Domain}(f) \quad \forall t \in [0, 1]$

must also be a convex set.



Simply convex

$$f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2)$$

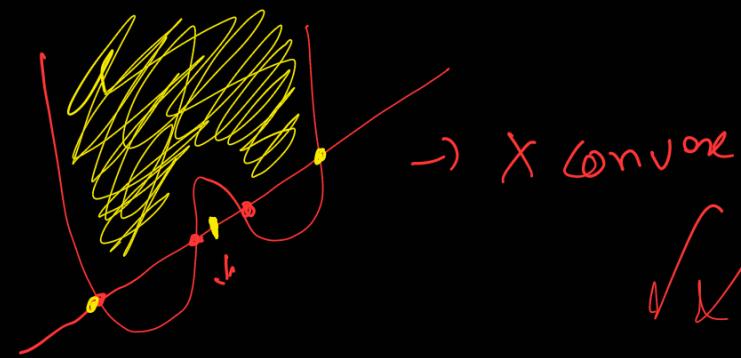
↓

Strictly convex \Rightarrow

$$f(tx_1 + (1-t)x_2) < tf(x_1) + (1-t)f(x_2)$$

\rightarrow Convex functions: \rightarrow local minima = global minima / maxima / maxima





✓ Set
✓ Function

$$l(\theta_1, \theta_2) = \theta_1^2 + \theta_2^2 + 2\theta_1\theta_2$$

$$\nabla_{\theta_1} l = \begin{bmatrix} - & 2\theta_1 + 2\theta_2 \\ - & 2\theta_2 + 2\theta_1 \end{bmatrix} \begin{array}{l} \frac{\partial}{\partial \theta_1} \\ \frac{\partial}{\partial \theta_2} \end{array} \left(\theta_1^2 + \theta_2^2 + 2\theta_1\theta_2 \right)$$

$$\nabla_{\theta_2} l = \begin{bmatrix} - & 2\theta_1 + 2\theta_2 \\ - & 2\theta_2 + 2\theta_1 \end{bmatrix} \begin{array}{l} \frac{\partial}{\partial \theta_1} \\ \frac{\partial}{\partial \theta_2} \end{array} \left(\theta_1^2 + \theta_2^2 + 2\theta_1\theta_2 \right)$$

$$\frac{\partial l(\theta)}{\partial \theta} = \begin{bmatrix} \frac{\partial l(\theta_1, \theta_2)}{\partial \theta_1} \\ \frac{\partial l(\theta_1, \theta_2)}{\partial \theta_2} \end{bmatrix}$$

Whenever I have a multi-variable function, its derivative will be a vector i.e. take partial derivative w.r.t every variable in the function.

$$l(\theta) = \theta_1^2 + \theta_2^2 + 2\theta_1\theta_2$$

$$\frac{\partial l(\theta)}{\partial \theta} = \begin{bmatrix} \frac{\partial l(\theta_1, \theta_2)}{\partial \theta_1} \\ \frac{\partial l(\theta_1, \theta_2)}{\partial \theta_2} \end{bmatrix} \in \begin{bmatrix} 2\theta_1 + 2\theta_2 \\ 2\theta_1 + 2\theta_2 \end{bmatrix} \xleftarrow[\quad]{\quad} \nabla l(\theta)$$

$$\nabla_{\theta} \left(\frac{\partial \ell(\theta)}{\partial \theta} \right) = \frac{\partial \ell(\theta)}{\partial \theta^2} \rightarrow \frac{\partial}{\partial \theta} \begin{pmatrix} \frac{\partial \ell_1}{\partial \theta_1} \\ \frac{\partial \ell_2}{\partial \theta_2} \end{pmatrix}$$

$$\nabla_{\theta} \left(X_{m \times n} \right)$$

\downarrow

$$D_{m \times n \times k \times l}$$

$$\frac{d}{dx} x^2 = 2x$$

$$\frac{d}{dx_1, dx_2} x^2 = \frac{dx}{dx_1} \frac{dx}{dx_2}$$



$$\nabla_{\theta} \left(2 \ell^2 \right) = \begin{pmatrix} \frac{\partial \ell_1}{\partial m_{11}} & \frac{\partial \ell_1}{\partial m_{12}} \\ \frac{\partial \ell_1}{\partial m_{21}} & \frac{\partial \ell_1}{\partial m_{22}} \\ \hline \frac{\partial \ell_2}{\partial m_{11}} & \frac{\partial \ell_2}{\partial m_{12}} \\ \frac{\partial \ell_2}{\partial m_{21}} & \frac{\partial \ell_2}{\partial m_{22}} \end{pmatrix}$$

matrix
↓
vector

$$\begin{pmatrix} \Delta_{11} & \Delta_{12} \\ \Delta_{21} & \Delta_{22} \end{pmatrix} \quad f = \theta_1^2 + \theta_2^2 + 2\theta_1\theta_2$$

$$\begin{matrix} \nabla f \\ \frac{\partial f}{\partial \theta_1} \end{matrix} \rightarrow \begin{pmatrix} \frac{\partial f}{\partial \theta_1} & \frac{\partial f}{\partial \theta_2} \end{pmatrix}$$

$$\rightarrow \begin{pmatrix} \frac{\partial f}{\partial \theta_1} \\ \frac{\partial f}{\partial \theta_2} \end{pmatrix} \quad \checkmark$$

$$\begin{pmatrix} (1x2) & (1x2) \\ (1x2) & (2x1) \end{pmatrix}$$

$$\begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} \xrightarrow{\text{J}} \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} \quad \nabla f = \begin{bmatrix} 2\theta_1 + 2\theta_2 \\ 2\theta_1 + 2\theta_2 \end{bmatrix}$$

$$\nabla f = \begin{bmatrix} 2\theta_1 + 2\theta_2 & 2\theta_1 + 2\theta_2 \end{bmatrix}$$

$$\begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} \xleftarrow{\text{J}} \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix}$$

$$\begin{bmatrix} e_1 \\ e_2 \\ e_3 \end{bmatrix} \quad \begin{bmatrix} e_1 \\ e_2 \\ e_3 \end{bmatrix} \quad \begin{bmatrix} e_1 & e_2 \\ e_3 & e_4 \end{bmatrix}$$

$$J(\theta) = \sum_{i=1}^m \left(y^{(i)} - h_{\theta}(x^{(i)}) \right)^2 \quad (\text{Loss})$$

$$\theta_j^{(t+1)} \leftarrow \theta_j^{(t)} - \eta \cdot \frac{\partial J(\theta)}{\partial \theta_j} \quad \text{Iterative}$$

Convergence Criteria: Σ can define myself.

T have 2 options: i) if $J(\theta)$ monitoring the loss value

I have captured (either I monitor a parameter of my interest.

(II) either I monitor a parameter of my interest.

(III) either I monitor $\text{grad}(\mathcal{J})$

ml model \rightarrow what we are trying to learn.

$$1. \left| \mathcal{J}(\theta^{t+1}) - \mathcal{J}(\theta^t) \right| < \delta$$

$$2. \left| \theta^{t+1} - \theta^t \right| < \epsilon$$

$$3. \left\| \nabla_{\theta} \mathcal{J}(\theta) \right\|_2 \leq \delta = 0$$

28th

29th

$\theta: 10$

t

$$(\theta_g)(0.008) \rightarrow (0.007954) \rightarrow 0.0005$$

↳ Vector

gradient (vec)

$$\left(\left[\quad \right] - \left[\quad \right] \right)$$

$$\left\| \left[\quad \right] \right\|_N \leq \text{threshold}$$

$$(2 \quad 3)$$

$$L2: \sqrt{2^2 + 3^2} = \sqrt{a+g} = \sqrt{13} \leq 4$$

$$y = mx + c$$

↓ ↓ ↓ ↓

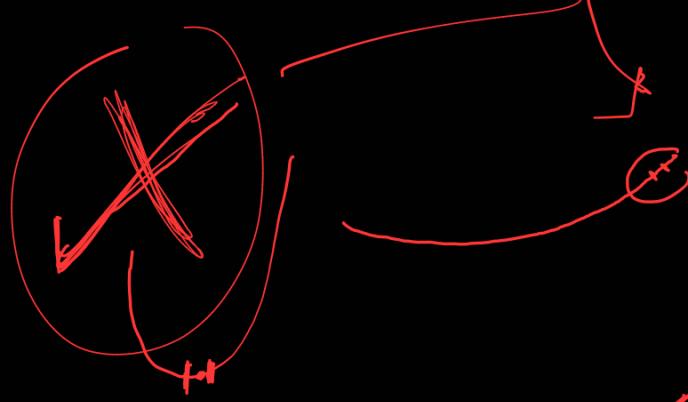
Param 1 Param 2 Param 1 Param 2

→ trying to learn params of st. line.

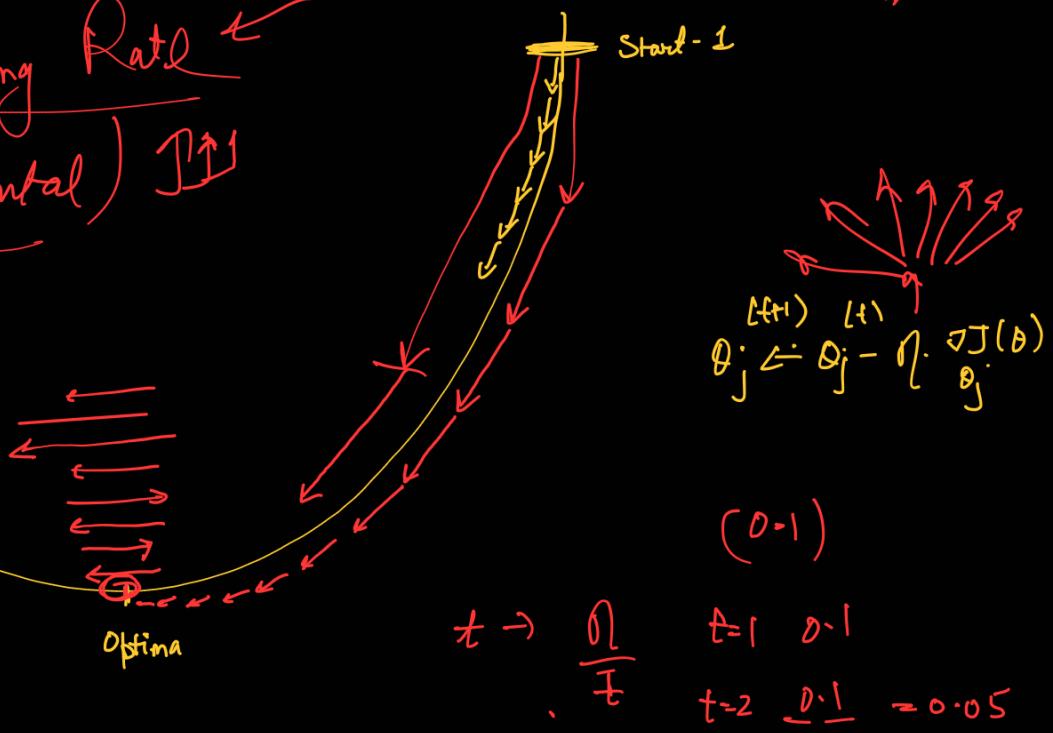
Good obs!

Hyperparameters: We give them.

$$4. \quad \left| \frac{\nabla J(\theta) - \nabla J(\theta^{(t)})}{\|\theta^{(t+1)} - \theta^{(t)}\|} \right| < \epsilon$$



Dynamic Learning Rate
(Experimental) ↑



$$\eta_t \propto \frac{1}{t} \quad \text{or} \quad \eta_t \propto \frac{1}{\sqrt{t}} \quad \not\propto \frac{1}{t^2}$$

$$\begin{aligned} t &\rightarrow \frac{1}{t} & t=1 & \eta_1 = 1 \\ && t=2 & \eta_2 = 0.5 \end{aligned}$$

$$\text{Initial: } \eta_t = \frac{\eta_0}{t} \quad n_{t+1} = n_t$$

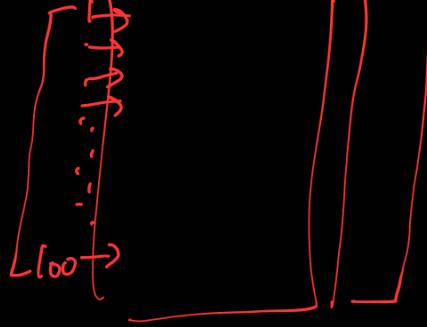
$$(\eta_{t+1} | \mathcal{L}(n_t))$$

Stochastic Gradient Descent (SGD)

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2 \quad (\text{MSE})$$

→ L m # entire training data.

I add up all losses, then take average of m losses.



\rightarrow [Advantage] \rightarrow I have seen my entire data.

$$\rightarrow \hat{\theta}_j^{(t+1)} \rightarrow \hat{\theta}_j^{(t)} - \eta \cdot \left[\frac{1}{m} \sum_{i=1}^m (y_i - \hat{\theta}_j^{(t)} x_j^{(i)}) x_j^{(i)} \right]$$

I have my 1 simple update after seeing entire data.

BATCH Gradient Descent

I can guarantee that my loss will converge (reduce to acceptable criteria) with each passing iteration.

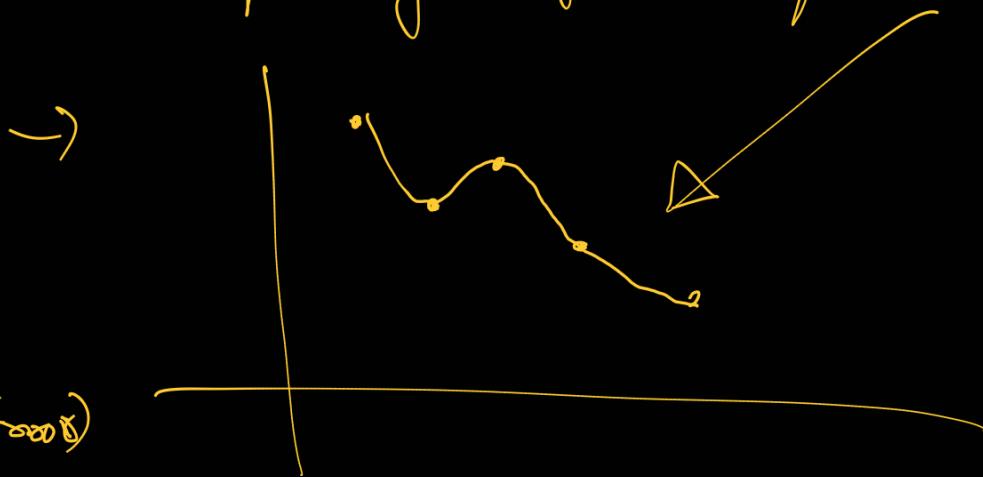
100 data points \rightarrow 100 \rightarrow $1000 \times 10 \rightarrow \dots 10^6, 10^7$
 millions, billions

\hookrightarrow Other extreme \rightarrow 1 data point

\downarrow (update) \rightarrow error has step $\frac{1}{\epsilon}$ count



I am plotting magnitude of loss values (Loss Curve)



$J_1 \rightarrow 0.9$
 $J_2 \rightarrow 0.7$
 $J_3 \rightarrow 0.8$
 $J_4 \rightarrow 0.6$
 $J_5 \rightarrow 0.5$

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

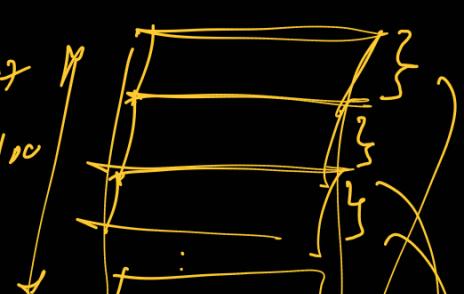
We are plotting this \rightarrow Value of $J(\theta)$ at $(x^{(i)}, y^{(i)})$

Mini-Batch Gradient Descent

$b_s = 10$

$$J^{(b)}(\theta) = \frac{1}{2b} \sum_{i=1}^b (y^{(i)} - h_{\theta}(x^{(i)}))^2$$

↳ b is my batch-size



$$\theta_j^{(t+1)} = \theta_j^{(t)} - \eta \cdot \frac{1}{b} \sum_{i=1}^b (y^{(i)} - h_{\theta}(x^{(i)})) q_j^{(i)}$$

Update

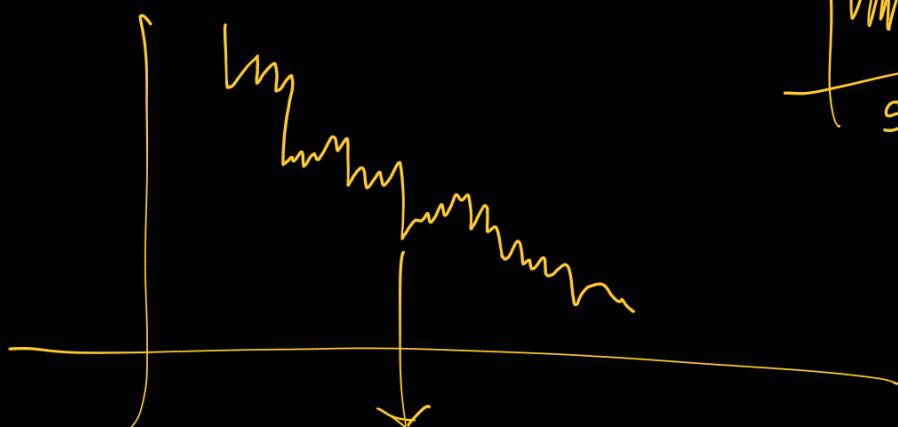
$\frac{1}{\delta} = \text{coupling}$

(δ)

on examples

6 batches

no of batch: $\left(\frac{m}{6}\right)$



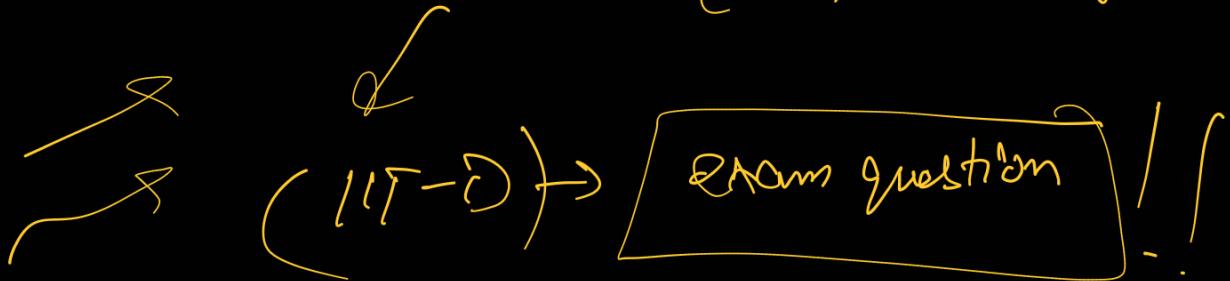
SG



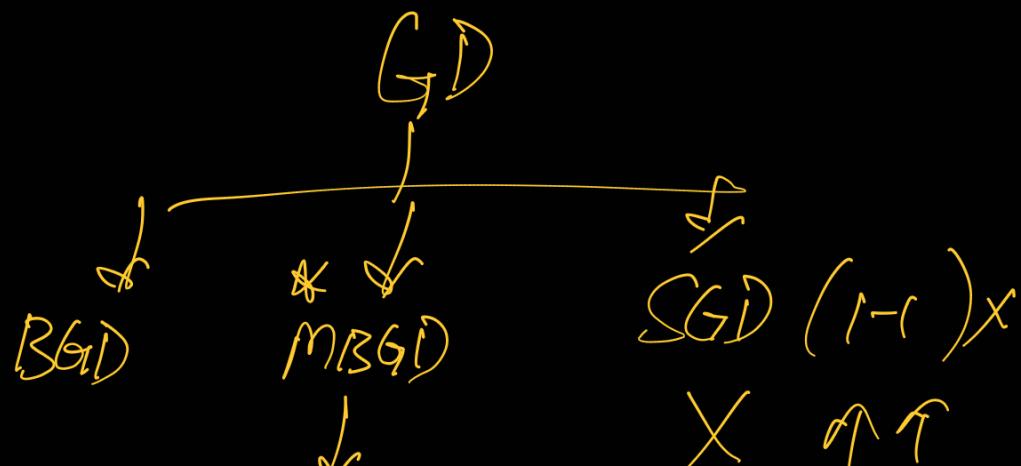
This I can prove
Mathematically.



In expectation, SGD (MBGD) would converge in the long run to BGD.



$$E \left[J^{(b)}(D) \right] \rightarrow \left\{ \text{Potential Exam Question} \right\}$$



$\rightarrow (SGD)$

Analytical Solution to Linear Regression

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)}))^2, \quad h_{\theta}(x^{(i)}) = \theta^T x^{(i)}$$

feature matrix:

$$X = \begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix}_{m \times n}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}_{(n+1) \times 1}$$

$$Y = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}_{m \times 1}$$

+ x_0

$$X = \begin{bmatrix} 1 & \vdots & 1 \end{bmatrix}_{(n+1) \times m}$$

$$\begin{array}{ccc} X & \theta & Y \\ (m \times (n+1)) & ((n+1) \times 1) & (m \times 1) \end{array}$$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (f_i - h_{\theta}(x^{(i)}))^2$$

$$J(\theta) = \frac{1}{2m} \|Y - X\theta\|_2^2$$

$$\frac{\|X\|_2}{X^T X}$$

$$= \frac{1}{2m} (Y - X\theta)^T (Y - X\theta)$$

$$(m \times 1) \quad ((n+1) \times 1)$$

$$\begin{matrix} (\mathbf{m} \times \mathbf{n}) \\ (\mathbf{m} \times \mathbf{l}) - (\mathbf{m} \times \mathbf{l}) \\ ((\mathbf{m} \times \mathbf{l}))^T \\ (\mathbf{l} \times \mathbf{m}) \quad (\mathbf{m} \times \mathbf{l}) \end{matrix}$$

$\equiv 1 \rightarrow \text{scalar}$

$$J(\theta) \leftarrow \frac{1}{2m} (y - x\theta)^T (y - x\theta)$$

$$(AB)^T = B^T A^T$$

$$J(\theta) = \frac{1}{2m} \left[(y^T - \theta^T x^T) (y - x\theta) \right]$$

$$J(\theta) = \frac{1}{2m} \left[y^T (y - x\theta) - \theta^T x^T (y - x\theta) \right]$$

$$= \frac{1}{2m} \left[y^T y - \underbrace{y^T x \theta}_{(1 \times m)(m \times (n+1))} - \underbrace{\theta^T x^T y}_{(n+1 \times 1)} + \theta^T x^T x \theta \right]$$

$$= \frac{1}{2m} \left[y^T y - \underbrace{y^T x \theta}_{(1 \times m)} - \underbrace{(y^T x \theta)^T}_{(m \times (n+1))} + \theta^T x^T x \theta \right]$$

$$= \frac{1}{2m} \left[y^T y - \underbrace{y^T x \theta}_{(1 \times m)} - \underbrace{\theta^T x^T x \theta}_{(1 \times 1)} \right]$$

$$\begin{aligned} \nabla_{\theta}(\theta^T A \theta) &= 2A\theta \\ \nabla_{\theta}(a^T \theta) &= a \\ \nabla_{\theta}(a\theta^T) &= a \end{aligned}$$

$$J(\theta) = \frac{1}{2m} \left[y^T y - 2y^T x \theta + \theta^T x^T x \theta \right]$$

$$J(\theta) = \frac{1}{2m} \left[y^T y - 2(x^T y)\theta + \theta^T (x^T x)\theta \right]$$

$$J(\theta) = \frac{1}{2m} \left[0 - 2x^T y + 2x^T x \theta \right] = 0$$

$$X^T X \theta = X^T Y$$

We are not sure
whether this inv.
exists.

$\rightarrow p_{inv}$

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

Normal Equations Method

Gradient Descent

Based Methods

N.Eqns

Inv calculation
is exp. consuming.

$$\mathcal{O}(n^2 m + n^3)$$

$$m = 1000$$

$$m = 1 \text{ Million}$$

$$\vec{v} = 2\hat{i} + 3\hat{j} + 4\hat{k}$$

$$\|\vec{v}\| = \sqrt{2^2 + 3^2 + 4^2} = \sqrt{29} \text{ norm}$$

$$|2| + |3| + |4| : L1$$

P-norm (v_1, \dots, v_n)

$\left(v_1^P + v_2^P + \dots + v_n^P \right)^{1/P} \rightarrow \|v\|_P$

$$\sqrt{2^2 + 3^2 + 4^2} = \sqrt{29} = \sqrt{S} = \sqrt{\sum v_i^2}$$

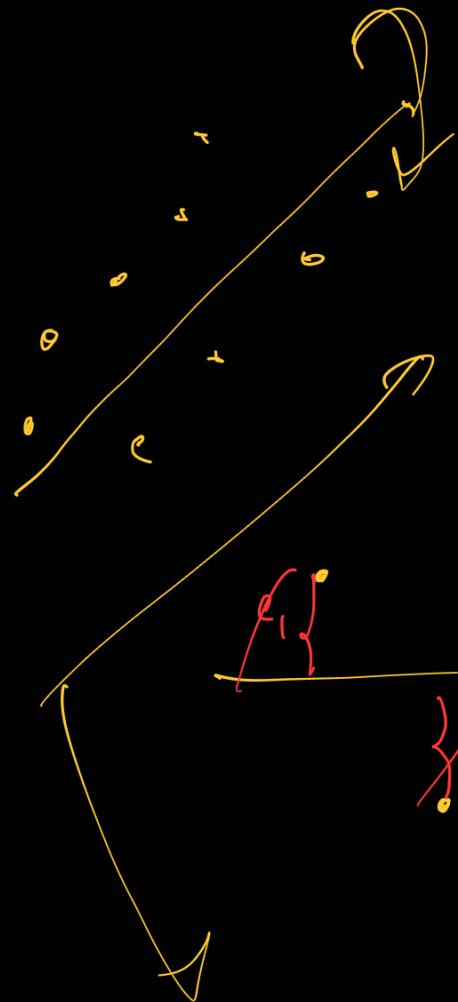
$B^T A - J \equiv$

$\sum x_i$



$$\text{Error}_{(1)} = (y_p - y_a)^2$$

$$\boxed{\text{Error}_{(2)} = (y_p - y_a)^2}$$



$$e_1 = e_2$$

↓

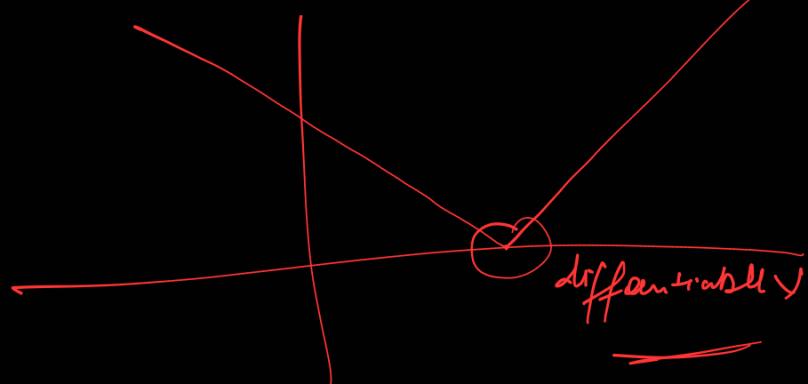
Sum

When we take sum
then 2 errors can cancel
each other.

$$\oplus (y_p - y_a) \cdot x$$

$$\backslash (+ (+ ($$

LL norm



Magnitude
Cancellation

Dif ✓

Squared
error

MCX

Dif ✓

