# Gradient Descent
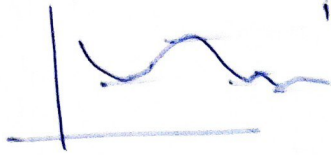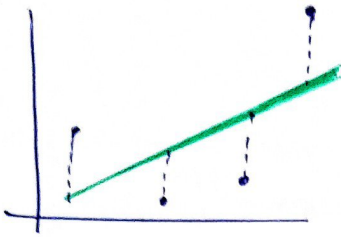
if we provide differtiable fun' it gives the minimum.

Ordinal least square

$$\hat{y} = mX_i + b$$

$$J = \sum_{0=1}^{1} (Y_i - \hat{Y}_i)^2$$

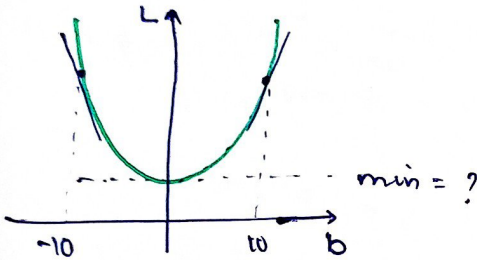$$J = \sum_{j=1}^{n} (m_i - mX_i - b)^2$$

$m = 78.35$

$$J(m,b) = \sum (Y_i - \overset{78.35}{m}X_i - b)^2$$

now only dependent upon $J(b)$

Step-1 select a Random b
let ∄ b = 0

How algo know where go up↑ down↓

How to find slope of $Y = x^2$

$$\frac{dy}{dx}\Big|_{x=5} = 2x = 10$$

Rel" of L, f b
$L \to (b)^2$

min = ?

$-10$　　$10$　$b$

$$\boxed{b_{new} = b_{old} - slope}$$

if slope(+ve) ⤳ b↓
　　-ve ⤳ b↑

Due to slope it change very fast so

$b = -10$
let slope = -50

$b_{new} = +0 - (-50)$
$= +0$

since ↯→

$b_{new} = 10 - (50$
$= -40$ ←↯

$$b_{new} = b_{old} - \eta\, slope$$

$\eta$ - learning Rate
0.01

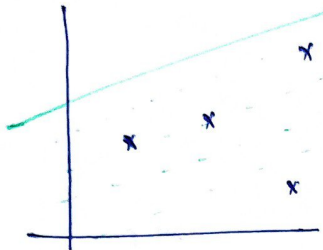eg- $b_{new} = -10 - (0.01 \times -50)$
$= -10 + 0.5$
$= -9.5$

**(?) when to stop**

$b_{new} - b_{old} \sim 0$

1. diff $> 0.0001$
2. Iteration (epochs)

## Mathematical Formulation

$m = 78.35$ (Known)

$(m, b)$

**Step-1** Start with Random $b$

for i in epochs:

$$b_{new} = b_{old} - \eta \, slope_{b=old}$$

$$J = \sum_{i=1}^{n} (Y_i - \hat{Y}_i)^2$$

$$\frac{d}{db} J = \frac{d}{db} \sum_{i=1}^{n} (Y_i - mx_i - b)^2$$

$$= 2 \sum_{i=1}^{n} (Y_i - mx_i - b)(-1)$$
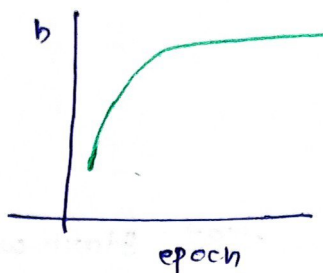
$$slope = -2 \sum_{i=1}^{n} (Y_i - mx_i - b)$$
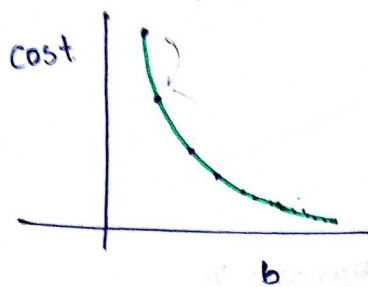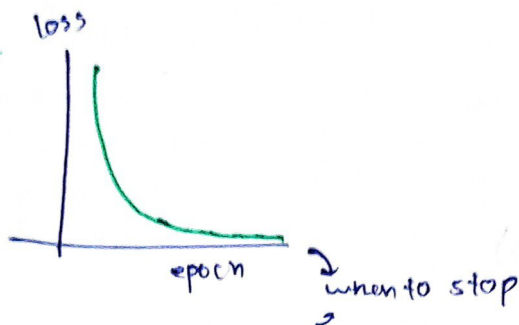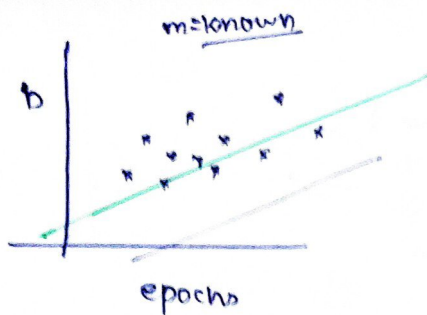
$$slope \Big|_{\substack{b=0 \\ m=78.35}} = -2 \sum_{i=1}^{n} (Y_i - 78.35X_i - 0)$$

we have 4 point now slove this ↑
is4 & we got slope
then we got $b_{new}$.

$$b_{new} = b_{old} - \underbrace{\eta \, slope_{b=old}}_{step-size}$$

m=known · loss

cost · b · epochs · epoch · when to stop · epoch

## Adding 'm'

S1: Random m & b     $L_r = 0.01$
                     copochs = 100

S2: for i in epochs:

$$b = b - \eta \cdot b\text{-slope}$$

$$m = m - \eta \cdot m\text{-slope}$$
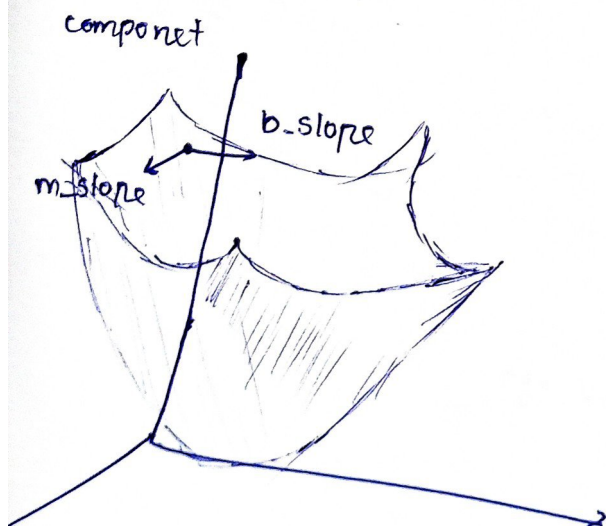
## Cost fun^c



loss · b · m

215 pv

$$J = \sum (Y_i - \hat{Y})^2$$

$$J_{(m,b)} = \sum (Y_i - mX_i - b)^2$$

$$\frac{\partial L}{\partial b} = -2\sum (Y_i - mX_i - b) = b\text{-slope}$$

$$\frac{\partial L}{\partial m} = 2\sum (Y_i - mX_i - b)(X_i) = m\text{-slope}$$

componet



b-slope · m-slope

, Standardized Data



cost

epochs

b

epochs

m

epoch

## Effect of loss function

$$J = \sum (Y_i - \hat{Y}_i)^2$$

↳ convex func | non covex func



cost

planato

touch two point | local | Global
w/o intercept | min | min
ay other.

Type

## Gradient Descent



| Batch GD | Stochastic GD | Mini-Batch GD |
|---|---|---|
| Total | one data | some data |

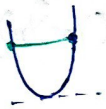| cgra | iq | lpa | |
|---|---|---|---|
| $X_1$ | $X_2$ | $Y$ | |
| 8.1 | 98 | 3.2 | $Y_1$ |
| 9.2 | 97 | 5.2 | $Y_2$ |

$$\hat{Y} = \beta_0 + \beta_1 X_1 + \beta_2 X_2$$

lpa  cgra  iq

1. Random values

$$\beta_0 = 0 \quad , \beta_1, \beta_2 = 1$$

2. epochs = 100 , lr = 0.01

$$\beta_0 = \beta_0 - \eta \, b_0\_slope$$

$$\beta_1 = \beta_1 - \eta \, b_1\_slope$$

$$\beta_2 = \beta_2 - \eta \, b_2\_slope$$

Our loss func in

d $J\langle \beta_0, \beta_1, \beta_2 \rangle$

$X_{i1}$  $X_{i2}$



| $X_1$ | $X_2$ | $Y$ |
|---|---|---|
| $X_{11}$ | $X_{12}$ | $Y_1$ |
| $X_{21}$ | $X_{22}$ | $Y_2$ |

row → 2
col → 2+1

$$\hat{Y} = \beta_0 + \beta_1 X_1 + \beta_2 X_2$$

$$\hat{Y_1} = \beta_0 + \beta_1 X_{11} + \beta_2 X_{12}$$

$$\hat{Y_2} = \beta_0 + \beta_1 X_{21} + \beta_2 X_{22}$$

$$J = \frac{1}{n} \sum_{i=1}^{n} (Y_i - \hat{Y_i})^2$$

$$= \frac{1}{2} \left[ (Y_1 - \hat{Y_1})^2 + (Y_2 - \hat{Y_2})^2 \right]$$

$$= \frac{1}{2} \left[ (Y_1 - \beta_0 - \beta_1 X_{11} - \beta_2 X_{12})^2 + (Y_2 - \beta_0 - \beta_1 X_{21} - \beta_2 X_{22})^2 \right]$$

$$\frac{\partial J}{\partial \beta_0} = \frac{1}{2} \left[ 2(Y_1 - \hat{Y_1})(-1) + 2(Y_2 - \hat{Y_2})(-1) \right]$$

$$= \frac{-2}{2} \left[ (Y_1 - \hat{Y_1}) + (Y_2 - \hat{Y_2}) \right]$$

let we have n rows ---

$$= \frac{-2}{n} \left[ (Y_1 - \hat{Y_1}) + (Y_2 - \hat{Y_2}) + \cdots + (Y_n - \hat{Y_n}) \right]$$

$$\frac{\partial J}{\partial \beta_0} = \frac{-2}{n} \sum_{i=1}^{n} (Y_i - \hat{Y_i})$$

$$\frac{\partial J}{\partial \beta_{①}} = \frac{1}{2} \left[ 2(Y_1 - \hat{Y_1})(-X_{11}) + 2(Y_2 - \hat{Y_2})(-X_{21}) \right]$$

$$= \frac{-2}{n} \left[ (Y_1 - \hat{Y_1})(+X_{11}) + (Y_2 - \hat{Y_2})(+X_{21}) + \cdots + (Y_n - \hat{Y_n})(X_{n1}) \right]$$

$$= \frac{-2}{n} \sum_{i=1}^{n} (Y_i - \hat{Y_i})(X_{i①})$$

$$\frac{\partial J}{\partial \beta_{②}} = \frac{-2}{n} \sum_{i=1}^{n} (Y_i - \hat{Y_i})(X_{i②})$$

$$\frac{\partial J}{\partial \beta_{(m)}} = \frac{-2}{n} \sum_{i=1}^{n} (Y_i - \hat{Y_i})(X_{i(m)})$$

$$\frac{\partial J}{\partial \beta_0} = \frac{-2}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)$$

Mean ↗ (over $\frac{1}{n}\sum$)

x_train ↙   prediction ↘

| $X_1$ | $X_2$ | $X_3$ | Y |
|---|---|---|---|
| $\hat{y}_1$ — | — | — | — |
| $\hat{y}_2$ — | — | — | — |

$= -2 \times \frac{1}{n} \sum K_i$

$=$ it is a scalar

$\hat{y} = np.dot(\text{xtrain\_coeff}) + \beta_0$

$\hat{y} = np.dot(X\_train, coeff) + \beta_0$

$= (353,10)\ (10,1) + \beta_0$

$= (353,1) + \beta_0$

$\therefore y\_pred$

**OR**

$\hat{y}_1 = \beta_0 + \beta_1 X_{11} + \beta_2 X_{12} + \beta_3 X_{13}$

$\hat{y}_2 = \beta_0 + \beta_1 X_{21} + \beta_2 X_{22} + \beta_3 X_{23}$

$\hat{y} = \beta_0 + [X_{11}\ X_{12}\ X_{13}] \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix}$

$\hat{y} = np.dot(coeff-, X\_train) + \beta_0$

| $X_1$ | $X_2$ | Y | $\hat{y}$ |
|---|---|---|---|
| 1 | 2 | 5 | 6 |
| 3 | 4 | 7 | 8 |

$$\frac{\partial J}{\partial \beta_1} \cdots \frac{\partial J}{\partial \beta_m} = \text{all dot} \begin{bmatrix} (y_i - \hat{y}) & \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \\ (1 \times 2) & (2 \times 2) \end{bmatrix} \times \frac{-2}{n}$$

$$= \left[ (y_i - \hat{y}_i)\ X\_train \right] \times \frac{-2}{n}$$

↙ no of row

$$\frac{\partial J}{\partial \beta_1} = \frac{-2}{n} \sum_{i=1}^{2} (y_i - \hat{y}_i)(X_{11})$$

$=$

$y - \hat{y} = [5\ 7][6\ 8]$

$= [-1\ -1]$

$(y-\hat{y})(X_n) = [-1\ -1] \begin{bmatrix} 1 \\ 3 \end{bmatrix}_{2 \times 1}$    $_{1 \times 2}$

$= [(-1 \wedge1 -3) + (-1 -3)]$

$= [(-4) + (-4)]$

$= 8$

$\frac{-2}{2}(8) = +8$

$\frac{\partial J}{\partial \beta_2} = \frac{-2}{2} \left[ [-1\ -1] \begin{bmatrix} 2 \\ 4 \end{bmatrix} \right]$

$= -1 \times [(-2-4) + (-2-4)]$

$= -8$

# The problem w/ Batch GD

ex.

| | |
|---|---|
| $\eta \to 100$ | $10^5$ |
| col $\to 5$ | $10^2$ |
| epoch $\to 50$ | $10^3$ |

for 1 col $\to 1000$
5 col $\to 6000$

Total = $6000 \times 50$ , $10^{10}$

1. Slow v. big data
2. Hardware

## Stochastic GD

single row one update
less epochs

$$\frac{\partial L}{\partial \beta_0} = \frac{-2}{n} \sum_{i=}^{n} (Y_i - \hat{Y})$$

$\hookrightarrow$ we calculate
for single row
so $n=1$

$$= -2(Y_i - \hat{Y}_i)$$

## Time Comparsion

e=100

batch $\quad\searrow$ Stochastic

100 update $\qquad$ $100 \times n$ updates

if the epoch is same then batch GD is faster.

due to this SGD don't required that much epoch so eve at the end it fast.
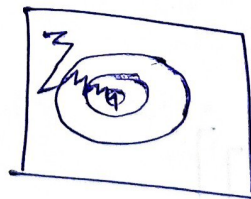
## learning Schedules

Due to randomness at end also it don't show optimal point so we vary w/ the epoch.

t0, t1 = 5, 50
def learning_rate (t):
$\qquad$ return t0/(t + t1)
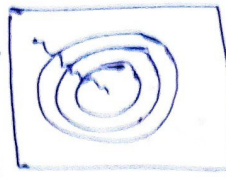
for i in range (epochs):
$\qquad$ for j in range (x.shape[0])
$\qquad$ $\alpha$ = learning_rate(i × x.shape)

when to use $\qquad$ big Data $\qquad$ Non convex

# Mini Batch GD

SGD is not optimal sol$^n$.

n rows = 1000

    batches = 100

    means = $\dfrac{1000}{100}$ = 10 updates/epoch

n rows
batch_size
$\rightarrow$ n    1/ epoch   BGD
$\rightarrow$ 1    n/epoch   SGD