

Gradient descent

$w^{k+1} = w^k - \underbrace{\text{update}}_{\text{d} \cdot \text{error}} \rightarrow \text{Gradient}$

Essential of statistics learning

by Hastie (Maths)

Introduction of statistics learning (R codes)

C.M. Bishop - for beginners

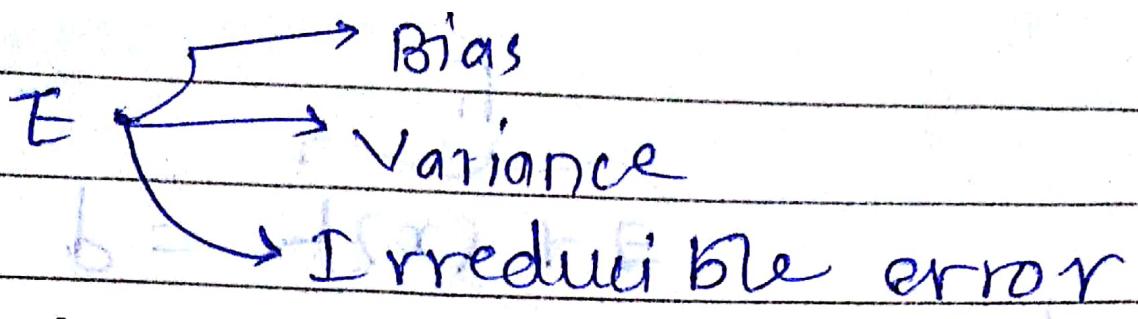
3) Mitchell - D.T., Gini, entropy.

function estimation / Approximation

$$p = m + c$$

$$P \rightarrow f(x) + E$$

y



~~Bias~~

b = variance

The moment you decide that model will grow as i thought

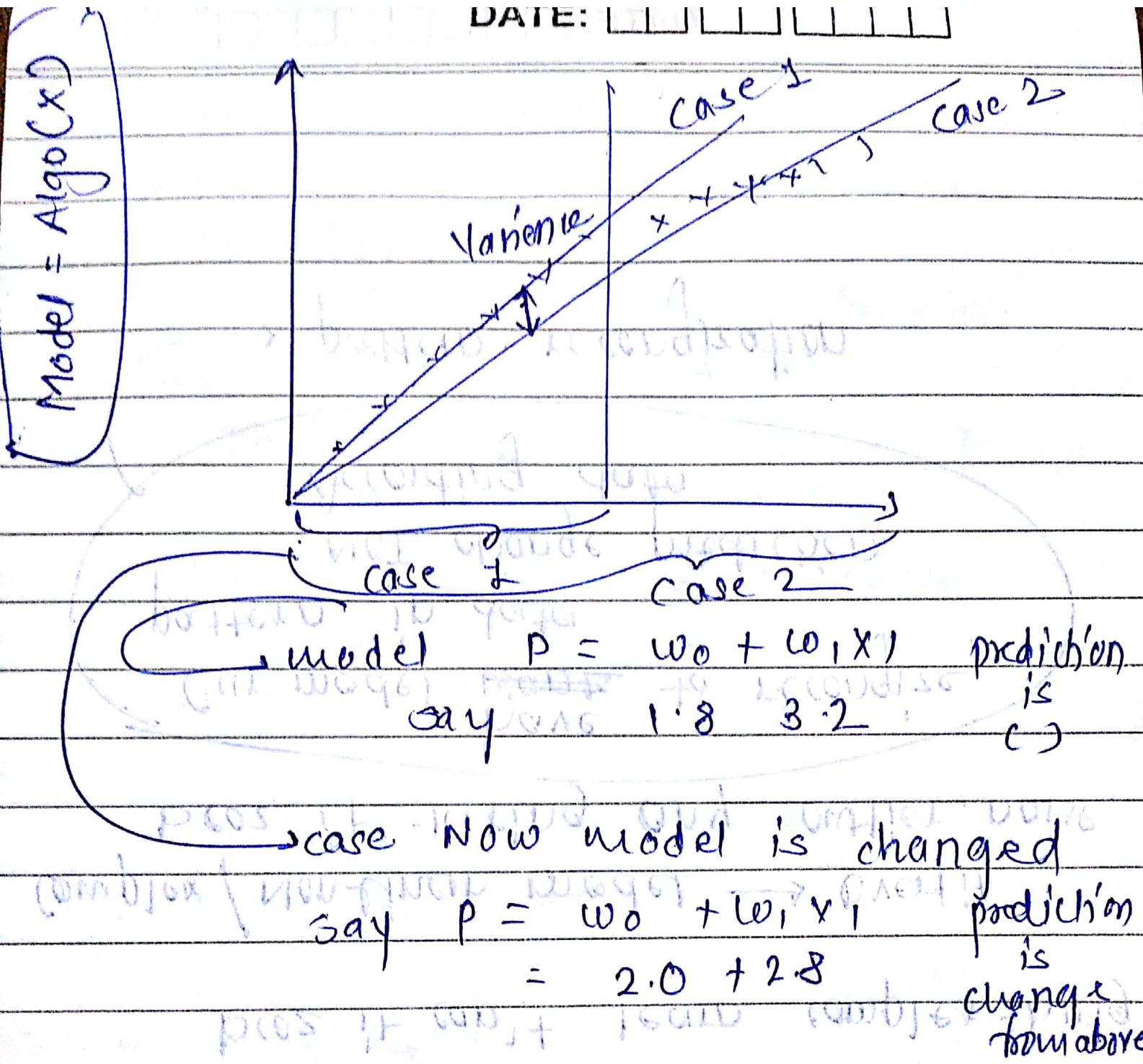
error due is introduce due to Assumption or naive Assumption or false assumption called it Bias error

The moment when decide model then condition may be total -

Variance:

⑤ who model is sensitive to data

The moment we introduce



so value of w_0 & w_1 is change due to you're train the model when you take new data points.

- So value of w_0 & w_1 is varied
- old data prediction is change when new data came in

high

~~bias~~

low

~~high~~

variance

highly Non-linear model → high variance

simple model → underfit

bcz it can't learn complex thing

Complex / Non-linear model → Overfit

bcz it taking any outlier, noise

Our model ~~wants~~ have to recognize pattern in data

NOT change prediction

According data

→ pattern recognition.

Irreducible error:

- can't reduce error whatever you try
- electronic, EMF noise, thermal noise, Eddy current

Instrument noise due to process data Acquisition

- Noise are occur due to process from process which getting us data.

Loss / cost function :-

Mean square loss

$$L = \frac{1}{N} \sum (Y_i - P_i)^2 \quad \text{M.S.E}$$

$$L = \frac{1}{2N}$$

$$\text{RMSF} = L = \sqrt{\frac{1}{N} \sum (Y_i - P_i)^2}$$

Why we take 2nd order? in MGE

$$\text{Loss } f^n \rightarrow \underset{w_i}{\operatorname{argmin}} (L(\hat{x}_i | y_i, w_i))$$

optimization
problem
formulation

I want to find
~~to~~ for which
argument my weight
is minimum

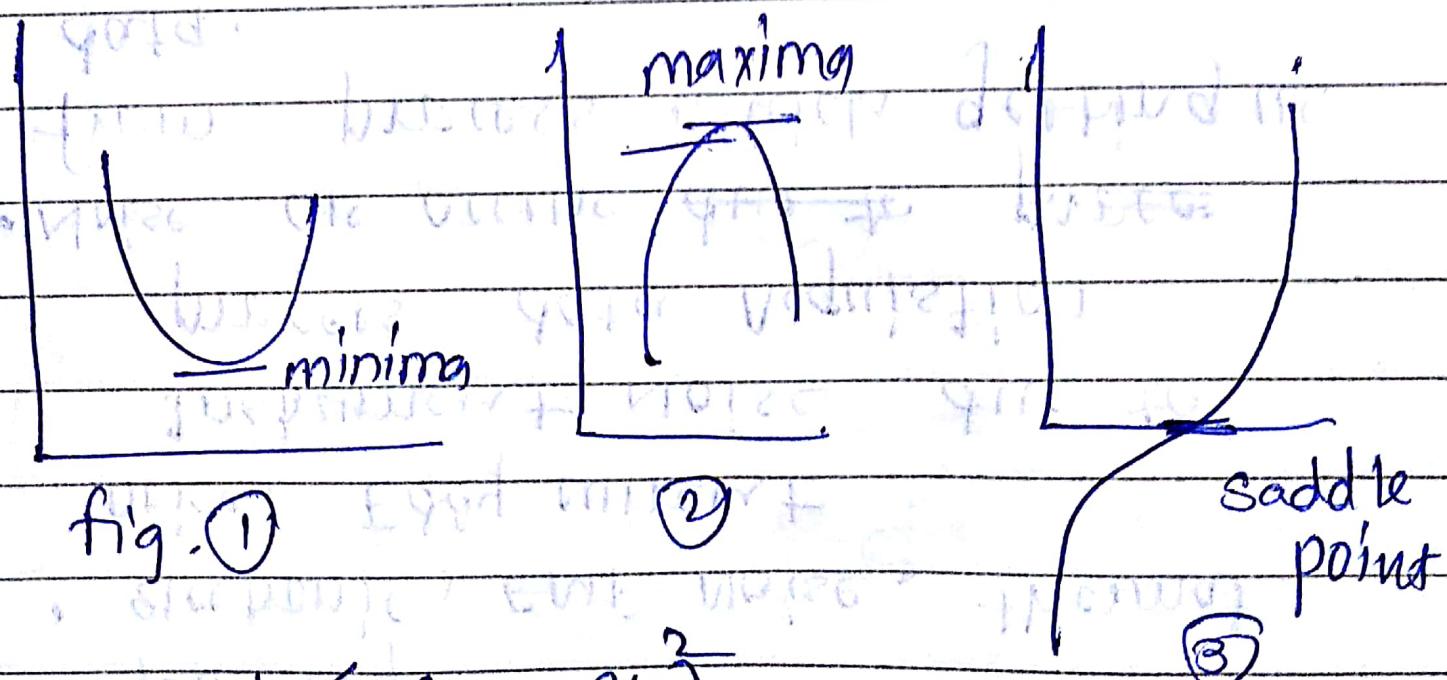
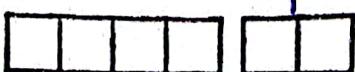


fig. ①

$$L = \frac{1}{N} \sum (y_i - p_i)^2$$

never gives saddle point
if need to ~~point~~ that saddle



DATE need to go high
order

$$\frac{\partial L}{\partial w_i} = 0$$

$$\frac{\partial^2 L}{\partial w_i^2}$$

= +ve (minimum).

$$\frac{\partial^2 L}{\partial w_i^2} = \text{constant} \rightarrow \text{minimum}$$

- using loss function we eliminate combination needed to decide (A) & (B)

$$\text{error} = y_i - p_i \quad \text{we}$$

let say me loss.

$$b_i = (y_i - p_i)^2 \rightarrow \text{this is for single}$$

$$\text{error} = y_i - \{w_0 + w_1 x_{i1}\} \rightarrow \text{instance for } i^{\text{th}}$$

$$x_{i1} = z_i$$

\downarrow
for i^{th} instant of x_1 feature.

$$L_i = [y_i - \{w_0 + w_1 z_i\}]^2$$

$$= y_i^2 - 2y_i(w_0 + w_1 z_i) + (w_0 + w_1 z_i)^2$$

$$\downarrow$$
$$w_0^2 + 2w_0 w_1 z_i + (w_1 z_i)^2$$

$$\frac{\partial L}{\partial w_0} = -2y_i + 2w_0 + 2w_1 a_i$$

$$\frac{\partial L}{\partial w_0} = -2[y_i - p_i] \xrightarrow{(1)} \boxed{\frac{\partial L}{\partial w_0} = 2}$$

$$\frac{\partial L}{\partial w_1} = -2y_i a_i + 2w_0 + 2w_1 a_i^2$$

$$= -2[y_i - p_i] a_i \xrightarrow{\text{sign}}$$

$$\frac{\partial L}{\partial w_1} = -2[y_i - p_i] a_i \xrightarrow{(2)}$$

$$\boxed{\frac{\partial L}{\partial w_1} = -2a_i^2}$$

ordinary
least square

(Analytic soln)

(Ordinary least square)

SGD.

computer don't does OLS part.

The equation which have statics
solution / deterministic solution

$$\text{like } y = 2x^2 - 3x + 1$$

$$\text{soln } y = 0$$

So this is NOT optimization

E we use optimization when we don't have clear solution about expression.

- for unconstraint optimization we use derivative & double derivative

Necessary condition

sufficient condition

vanishing

Next derivak
should be

derivative

+ve or

is equal

if 2nd derivak
+ve

to zero

if 2nd derivak
then go

for 3rd derivak

if 3rd derivak

then go

for 4th derivak

if 4th derivak

then go

for 5th derivak

if 5th derivak

then go

for constraint optimization! E. constraint
for equality constraint.

Lagrangian's equation

$$L = f(x) + \lambda_1(g_1(x)) + \lambda_2(g_2(x))$$

Now you have $\frac{\partial L}{\partial x}$ $\frac{\partial L}{\partial \lambda_1}$ $\frac{\partial L}{\partial \lambda_2}$

so Constraint optimization is difficult than unconstraint.

SVM = Maximal Marginal Classification

which requires optimization

so we are here

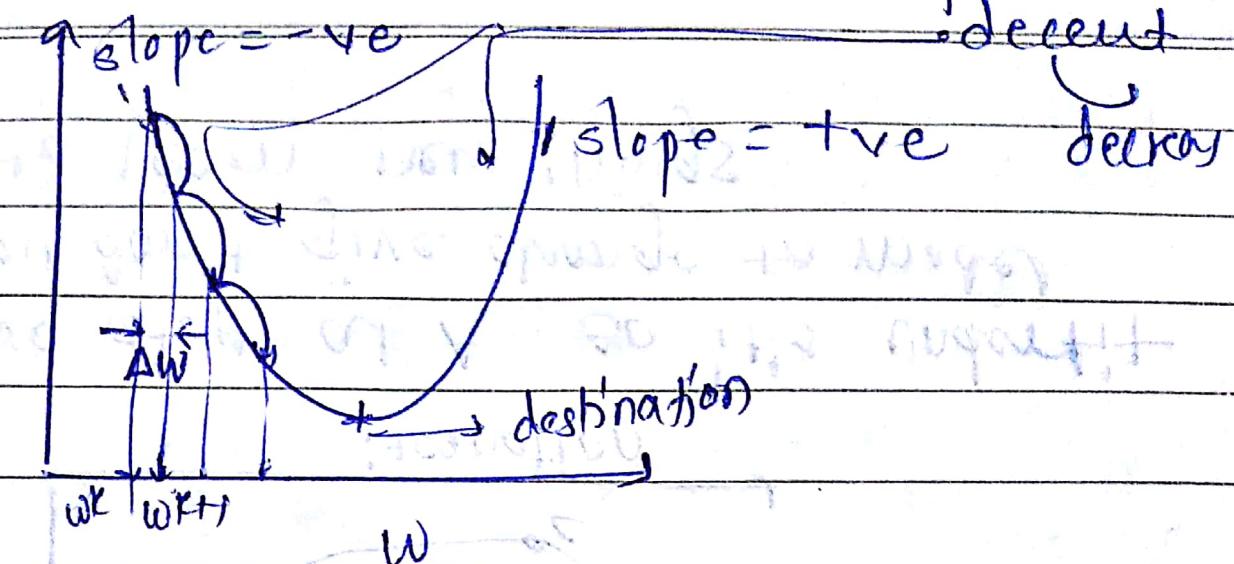
$$\frac{\partial L}{\partial w_i} = -2[y_i - p_i] * g_i \quad \text{features}$$

gradient $w_i^{k+1} = w_i - \Delta w$

descent

$$\Delta w = -\alpha \left(\frac{\partial L}{\partial w} \right)$$

DATE: [] gradient descent



- how to Algorithm know do i have to go forward (-ve slope) or backward (+ve slope) is known by gradient.
- and at what step / speed is α

$$\frac{dL}{dw}$$

$$\Delta w = \alpha \left(\frac{dL}{dw} \right)$$

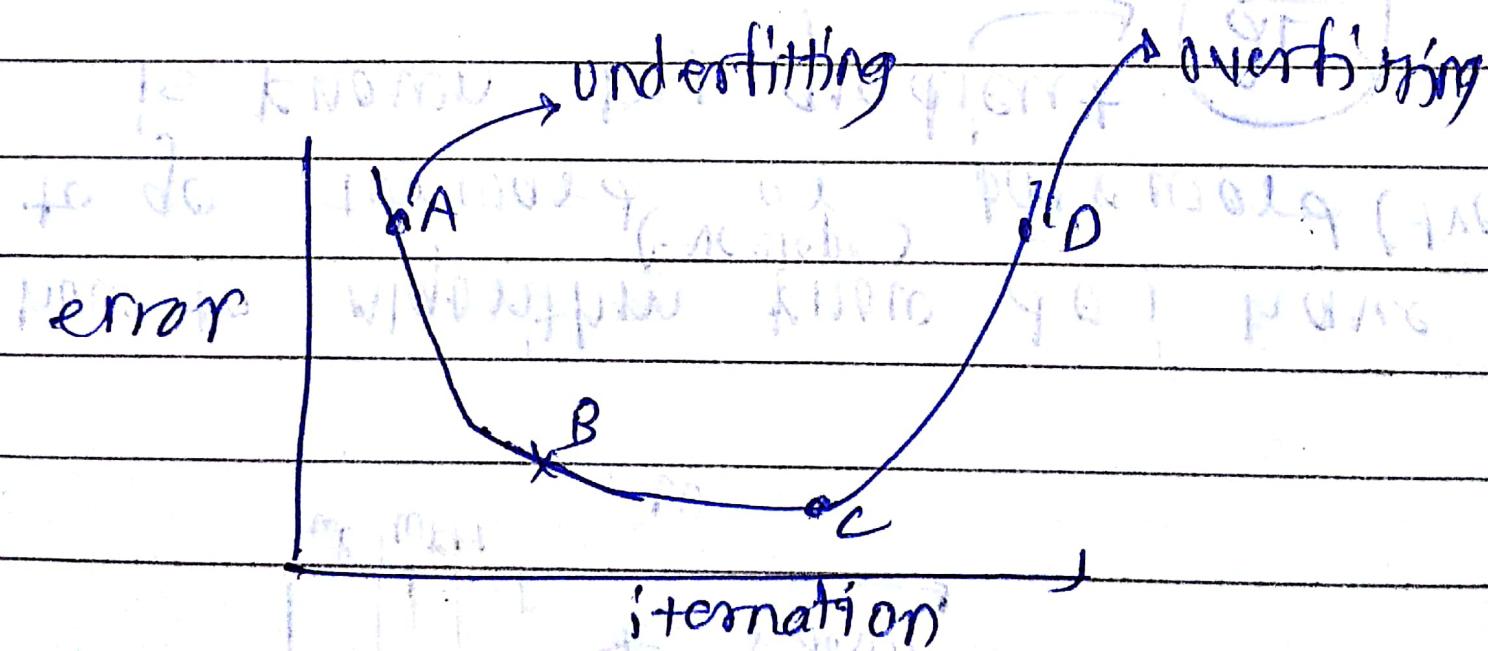
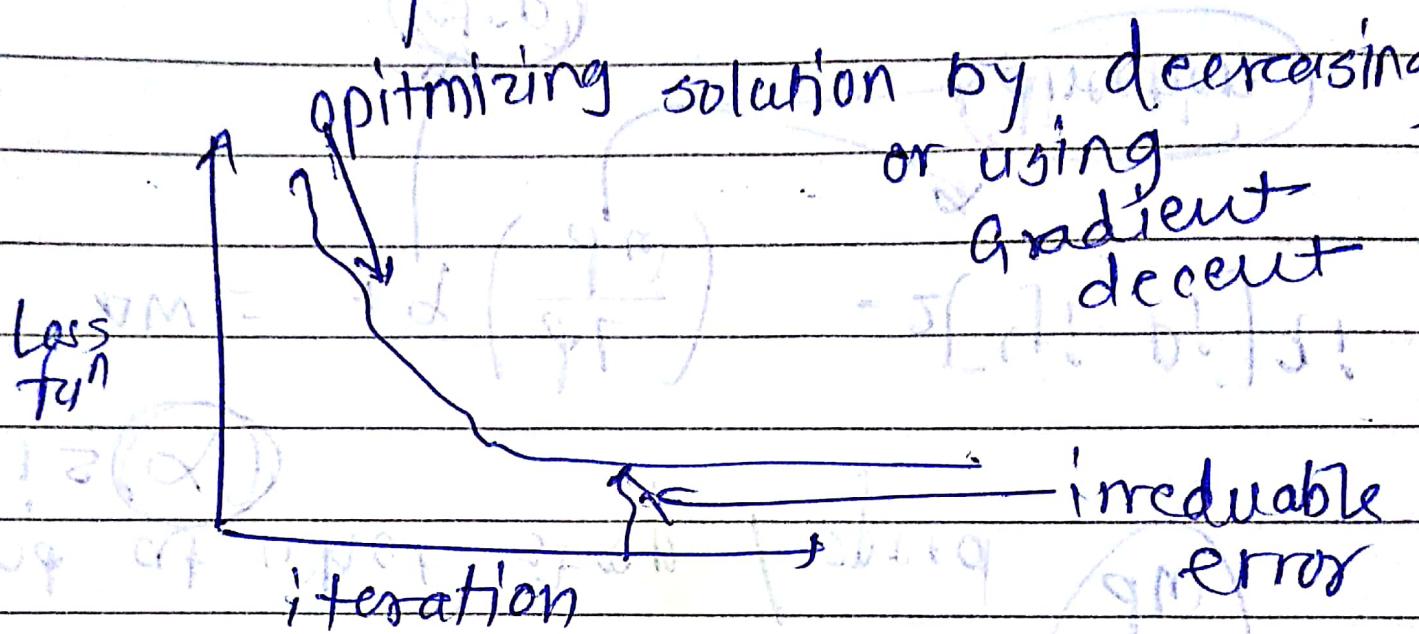
$$-2[y_i - p_i] \gamma_i$$

step size

direction

for minimization $w_i^{k+1} = w_i^k - \Delta w$ for Maximization $w_i^{k+1} = w_i^k + \Delta w$

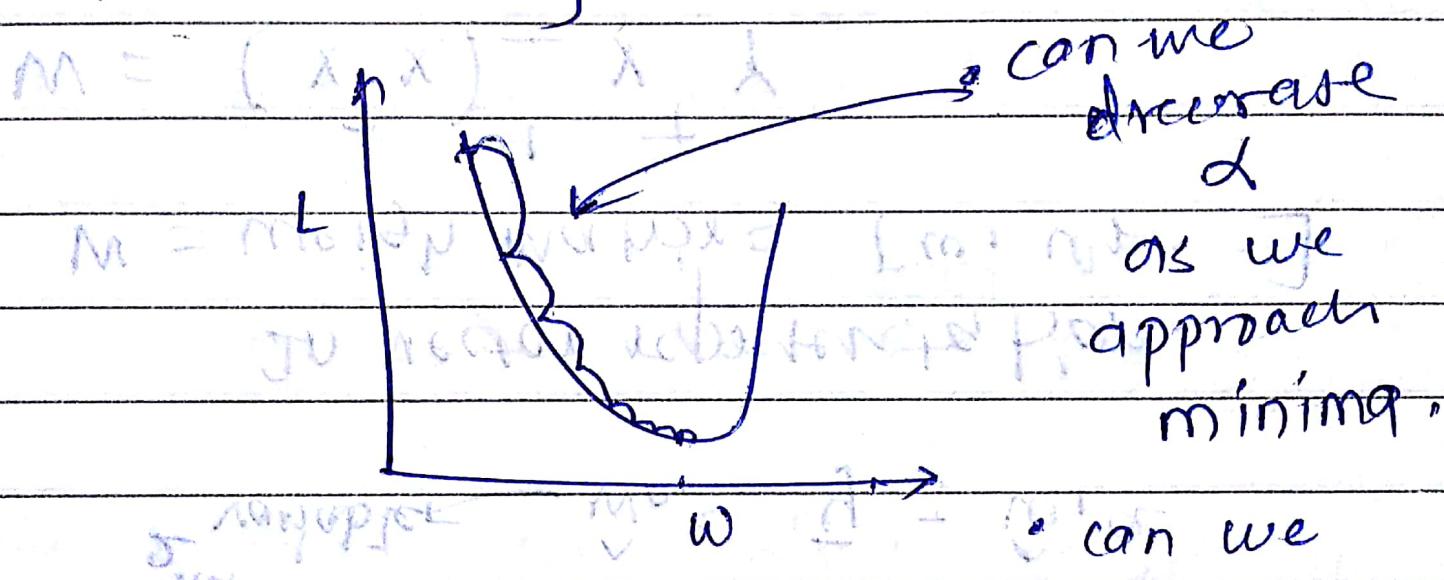
This proof the problem is optimization and this use gradient decent.



- if we stop at A so it's underfit you don't give chance to model to learn new things

- we get at D due to overfitting, complex model, ~~hang~~ more iteration

α = learning rate.



- but gradient descent ← make already doing that

because of vanishing slope
as slope goes with constant α
we never get to minima.

- Until Now we just consider only one instant. \rightarrow in dataset.

Now take for all dataset
 \sum term

for $L = \frac{1}{N} \sum (y_i - p_i)^2$

\rightarrow square error \rightarrow w_1 \rightarrow \hat{y}

\rightarrow specific w1 depends on x^2 values

st ↴

specific equation eg. $\hat{y} = a_1 + a_2 + a_3$

• perf. requirement present \leftarrow 3 vars

2nd variable $w_0 \in \hat{y} = \hat{w}_1 a_1$

In vector representation

W = weight Matrix = $[w_0, w_1, \dots]$

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

α = learned weights

$$Y = \begin{bmatrix} w_0 + w_1 x_1 \\ w_0 + w_1 x_2 \\ \vdots \\ w_0 + w_1 x_n \end{bmatrix} = X^T W$$

$$X = \begin{bmatrix} w_0 & w_1 \\ 1 & 1 \\ 1 & 3 \\ 1 & 5 \end{bmatrix}$$

$$X^T X = 0$$

$$J = \frac{1}{N} (y - p)^T (y - p)$$

similar to

$$L = \frac{1}{N} \sum (y_i - p_i)^2$$

logistic regression: gives linear funⁿ

$$\log \left[\frac{p(x)}{1-p(x)} \right] = w_0 + w_1 x$$

$$p(x) = \frac{e^{w_0 + w_1 x}}{1 + e^{w_0 + w_1 x}}$$

even.

why it called regression if it gives classification?

Cuz we working on number

we working b/w 0 to 1

- continuous value

- and we checking value is > 0.5 or < 0.5

Homework

prove that $\log \left[\frac{p(x)}{1-p(x)} \right] = w_0 + w_1 x$

given $p(x) = \frac{e^{w_0 + w_1 x}}{1 + e^{w_0 + w_1 x}}$