

## CH- Optimization problem solution

→ Solving optimization problems: differentiation.

We have seen optimization previously at 3 places.

① PCA

② Logistic Reg

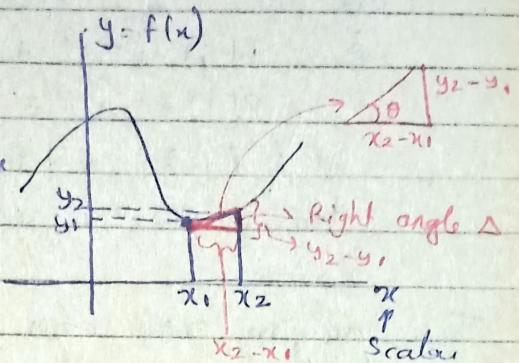
③ Linear Regression.

We use concepts like ① Differentiation ② Maxima - minima

→ Single variable differentiation.

so what does  $\frac{dy}{dx}$  intuitively mean.

$\frac{dy}{dx}$  means how much does  $y$  change as  $x$  changes or we can say rate of change of  $y$  w.r.t.  $x$ .

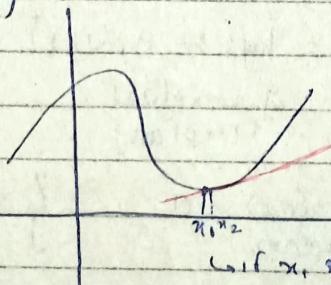


$$\frac{dy}{dx} \leftarrow \frac{\text{change in } y}{\text{change in } x} \Rightarrow \frac{\Delta y}{\Delta x} = \tan \theta$$

Now let's understand mathematically

$$\frac{dy}{dx} = \lim_{\Delta x \rightarrow 0} \frac{\Delta y}{\Delta x}$$

$\frac{dy}{dx}$  = tangent of slope of.



if  $x_1, x_2$  are very close  
hypotenuse becomes tangent.

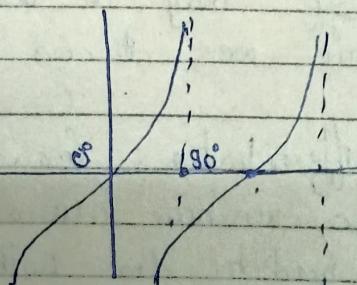
Let's understand Tan θ

$0 < \theta < 90^\circ \Rightarrow +ve$

~~Tan 0~~  $\theta = 0$

$\tan 90^\circ = \text{undefined}$

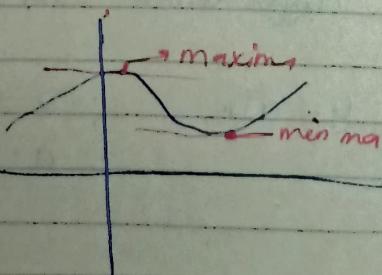
$90^\circ < \theta < 180^\circ \Rightarrow -ve$



\* Maxima & Minima

at both minima & maxima

slope = 0



Ex:-  $f(x) = \frac{3x^2}{2} + 3x + 2$  find minima & maxima  
 $\frac{df}{dx} = 2x + 3$

$\therefore 2x + 3 = 0 \quad \boxed{x = -\frac{3}{2}}$  → at  $x = -\frac{3}{2}$  slope = 0, so is the a maxima or minima.

One simple way is to compare ~~with~~ with  $x$  with 2 values one at  $\frac{3}{2}$  and one at some other value.

Ex:-  $f(\frac{3}{2}) = (\frac{3}{2})^2 + 3 \times \frac{3}{2} + 2$   
 $f(1) = 1 + 3 + 2$

If  $f(\frac{3}{2}) < f(1)$  then  $f(\frac{3}{2})$  is minima because there is some other value greater than  $f(\frac{3}{2})$ .

There may be some cases that func have no maxima or the func have no minima.

It is also possible that func have many local max and min.

Ex:-  $f(x) = \log(1 + \exp(ax))$

$$\frac{df}{dx} = \frac{a \cdot \exp(ax)}{1 + \exp(ax)}$$

$$\frac{a \cdot \exp(ax)}{1 + \exp(ax)} = 0 \quad \text{Solving this is not trivial.}$$

Here to calculate maxima or minima just trying to find  $g_f$  may be very hard sometimes.  
In such case we can use gradient descent

→ Vector differentiation : Grad.

Till now we have assumed that  $x$  is scalar. But in ML we need to deal with vectors. So let's see how to perform differentiation on vectors.

$y = a^T x$  ;  $x = \langle x_1, x_2, \dots, x_d \rangle$   
 $a = \langle a_1, a_2, \dots, a_d \rangle \rightarrow \text{constants.}$

$$\frac{df}{dx} = \nabla_x f \quad \begin{matrix} \text{grad or Del} \\ \text{vector} \end{matrix}$$

$$\nabla_x f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix} = \begin{bmatrix} \text{partial diff.} \\ \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix}$$

$$y = a^T x = a_{11}x_1 + a_{21}x_2 + a_{31}x_3 + \dots \quad \text{Partial}$$

$$\frac{\partial f}{\partial x_1} = a_1$$

$$\frac{\partial f}{\partial x_2} = a_2$$

$$\frac{\partial f}{\partial x_d} = a_d$$

$$\nabla_x(a^T x) = a \quad \text{vector}$$

$$L(w) = \sum_{i=1}^n \log \left( 1 + \exp(-y_i w^T x_i) \right) + \lambda w^T w$$

$$\nabla_w L = \frac{(-y_i x_i)}{1 + \exp(-y_i w^T x_i)} + 2\lambda w = 0$$

Now solving this is hard. we have to use gradient descent

Gradient Descent :- Now, let's learn gradient descent algorithm  
It is an iterative algorithm

$x_0 \leftarrow$  first guess of  $x$

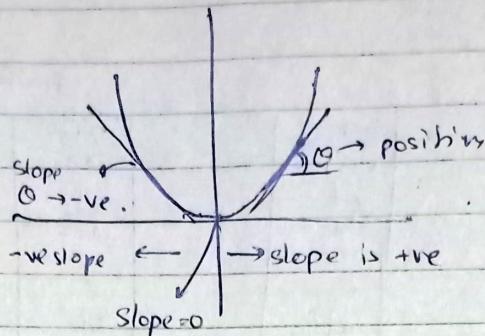
$x_1 \leftarrow$  iteration 1

$x_2 \leftarrow$  iteration 2

$x_n \leftarrow$   $n^{th}$  iteration

As our iteration number increases, we go closer to right  $x$ .

Slope changes its sign from +ve to -ve at minima.



Gradient descent Algo:-

① pick an initial pt  $x_0$  at random

②  $x_1 = x_0 - \gamma \left[ \frac{df}{dx} \right]_{x_0}$  → slope of the tangent :  
constant: stepsize.

Let  $\gamma = 1$  for now.

If  $\frac{df}{dx} > 0$  Then  $\gamma \cdot \frac{df}{dx} > 0$

$$x_1 < x_0$$

$$\therefore x_0 - \gamma \left[ \frac{df}{dx} \right] > x_0$$

∴ we move closer.

$$\text{If } \frac{df}{dx} \leq 0 \quad \therefore \quad x_1 = x_0 + \gamma \left[ \frac{df}{dx} \right]$$

$$x_1 > x_0$$

again we move closer. because this line never on other side

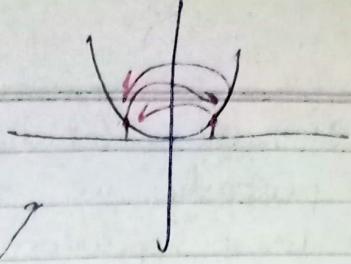
at  $i^{th}$  iteration

$$x_{i+1} = x_i + \gamma \left[ \frac{df}{dx} \right]_{x_i}$$

at any iteration  $K$

If  $(x_{K+1} - x_K)$  is very small then terminal  
at  $x^* = x_K$

$$x_{i+1} = x_i - \alpha * \left[ \frac{\partial f}{\partial x_i} \right]_{x_i}$$



$\alpha$  is called learning rate or step size.

Some time we get oscillation problem due to ~~gradient~~ constant step size. To remove that we can change  $\alpha$  at each iteration w.r.t. any of the following technique.

①  $\rightarrow$  reduce  $\alpha$  with each iteration.

② we can choose any func to change  $\alpha$ .

### Gradient descent for Linear Regression.

$$w^* = \underset{w_0}{\operatorname{argmin}} \sum_{i=1}^n (y_i - w^T x_i)^2$$

→ Not using  $w_0$  & Regular term for simplicity.

$$L = \sum_{i=1}^n (y_i - w^T x_i)^2$$

if  $y_i$ ,  $x_i$ , are part of training data  
hence they are ~~constant~~ vectors of  
constant, we do one optim.

$$\nabla_w L = \sum_{i=1}^n \{ 2(y_i - w^T x_i)(-x_i) \}$$

① pick a random vector  $w_0 = \langle \quad \rangle$

$$② w_1 = w_0 - \alpha \sum_{i=1}^n (-x_i)(y_i - w_0^T x_i)$$

$$③ w_2 = w_1 - \alpha \sum_{i=1}^n (-x_i)(y_i - w_1^T x_i)$$

$$(k) w_{k+1} = w_k - \alpha \sum_{i=1}^n (-x_i)(y_i - w_k^T x_i)$$

if ( $w_{k+1} - w_k$ ) ~~not~~ returns.

There is one problem here for every iteration we are computing the summation. If dataset is very large there is lot of time & computing power req.  
for each iteration we are visiting whole data.

To solve this problem we use technique called SGD

## ★ Stochastic Gradient descent (SGD)

Previously we saw

for linear regression

$$GD: w_{j+1} = w_j - \alpha \sum_{i=1}^n (-x_i) (y_i - w_j^T x_i)$$

computing sum on whole data set at each iteration.

$$SGD: w_{j+1} = w_j - \alpha \sum_{i=1}^{(K)} (-x_i) (y_i - w_j^T x_i)$$

$\hookrightarrow 1 \leq K \leq n$

Pick random set of  $K$ -points from total dataset of  $n$  points and apply GD on these  $K$ -points will give same result.

It is called stochastic because we are using random val.

If we pick smaller points i.e.  $K$  is very less then no. of iteration increases.

$K$  is also called batch-size.

and at each iteration we pick different random set of  $K$  points.

## ★ Constrained optimization:-

Till now we have seen unconstrained optimization problems like  $\max_n f(n)$  or  $\min_n f(n)$ . But PCA formulation is

$$PCA: \max_U \frac{1}{n} \sum_{i=1}^n (U^T x_i)^2 \quad \text{s.t. } U^T U = I$$

$\underbrace{\quad \quad \quad}_{\text{Objective func}}$        $\underbrace{\quad \quad \quad}_{\text{constraint.}}$

These types of problem are called constrained optimization problems.

General constraint func looks like:

$$\max_x f(x) \quad \text{s.t. } g(x) = c \rightarrow \text{equality constraint}$$

$h(x) \geq d \rightarrow \text{inequality constraint.}$

so how to find maxima/minima of such type of funcn.

$$\begin{array}{l} \max_x f(x) \\ \text{s.t. } g(x) = c \\ h(x) \geq d \end{array} \quad \left\{ \begin{array}{l} \xrightarrow{\text{modify it using}} \\ \text{Lagrangian multipliers} \end{array} \right. \quad \text{Lagrangian Multiplier.}$$

$$L(x, \lambda, \mu) = f(x) - \lambda g(x) - \mu h(x)$$

$\lambda, \mu$  are Lagrangian multipliers.  $\lambda \geq 0, \mu \geq 0$

$$\frac{\partial L}{\partial x} = 0; \quad \frac{\partial L}{\partial \lambda} = 0; \quad \frac{\partial L}{\partial \mu} = 0$$

Solve these equations: suppose we get  $\tilde{x}$  after solving this. Then  $\tilde{x}$  is the solution of the given optimization problem.

Now, let's revisit the PCA problem.

$$\text{PCA: } \max_u \frac{1}{n} \sum_{i=1}^n (u^T x_i)^2 \quad \text{s.t. } u^T u = 1$$

↓

$$\max_u u^T S u \quad \text{s.t. } u^T u = 1 \quad \times S = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$$

$S$  is a constant (because it contains covariance) covariance matrix

$$L(u, \lambda) = u^T S u - \lambda(u^T u - 1)$$

$$\frac{\partial L}{\partial u} = 0 \Rightarrow \frac{\partial}{\partial u} (u^T S u - \lambda u^T u - 1) = 0$$

$$= Su - \lambda u = 0$$

$Su = \lambda u$  // definition of eigen-value & eigen-

$u$  is eigen vector of  $S$

$\lambda$  is eigen value of  $S$

That is what we have done in PCA, we choosed the top eigen valued vectors of co-variance matrix. This is the proof.

Now, let's revisit Logistic Regression formulation in context of  
constrained optimization

LR:

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} (\text{logistic\_loss}) + \lambda \mathbf{w}^T \mathbf{w}$$

Actual problem is, to minimize the logistic loss s.t.  $\mathbf{w}$  is normal to hyperplane and  $\mathbf{w}$  is unit vector.

$$\therefore \mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} (\text{logistic\_loss}) \quad \text{s.t. } \mathbf{w}^T \mathbf{w} = 1$$

↑  
objective funcy      ↓  
constraint

Just use Lagrangian to solve this

$$\mathcal{L}(\cdot) = \text{logistic\_loss} - \lambda(1 - \mathbf{w}^T \mathbf{w})$$

$$\Rightarrow \text{logistic\_loss} - \lambda + \lambda \mathbf{w}^T \mathbf{w}$$

$\therefore$  regularization can be thought of as imposing an constraint

\* Why does  $L_1$ -regularization create sparsity.

while studying Logistic Regr we said that  $L_1$  reg creates sparsity in  $\mathbf{w}$  as compared to  $L_2$ -regularization. Now we will look at geometric intuition of it

$L_2$  regular

$$\min_{\mathbf{w}} \text{loss} + \lambda \|\mathbf{w}\|_2^2$$

$$\cancel{\|\mathbf{w}\|_2^2}$$

$$\min_{\mathbf{w}} \|\mathbf{w}\|_2^2$$

$$\min_{\mathbf{w}} (w_1^2 + w_2^2 + \dots + w_d^2)$$

Let choose only  $w_1$  for now

$$\min_{w_1} w_1^2 = L_2(w_1)$$

$L_1$  regular

$$\min_{\mathbf{w}} \text{loss} + \lambda \|\mathbf{w}\|_1$$

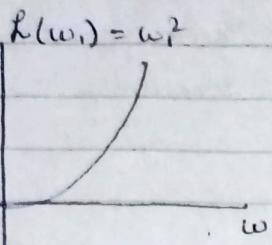
$$\cancel{\|\mathbf{w}\|_1}$$

$$\min_{\mathbf{w}} \|\mathbf{w}\|_1$$

$$\min_{\mathbf{w}} (|w_1| + |w_2| + \dots + |w_d|)$$

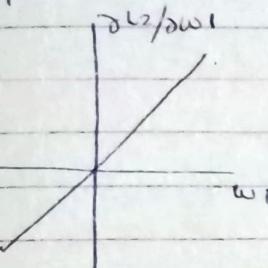
$$\min_{w_1} |w_1| = L_1(w_1)$$

12

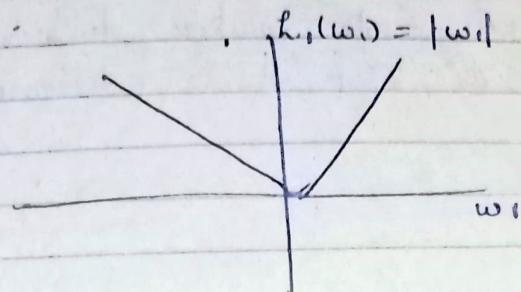


$$L_2(w_i) = w_i^2$$

$$\frac{\partial L_2}{\partial w_i} = 2w_i$$

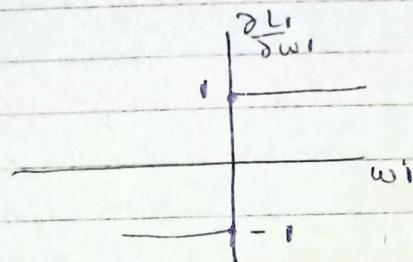


11



$$L_1(w_i) = |w_i|$$

$$\frac{\partial L_1}{\partial w_i} = \begin{cases} +1 & \text{if } w_i > 0 \\ -1 & \text{if } w_i < 0 \end{cases}$$



Let  $(w_i)_i$  is the value of  $w_i$  in  $i$ th iteration of Gradient Descent

Given:  $(w_i)_0, (w_i)_1, \dots, (w_i)_K$

$$(w_i)_{i+1}^* = (w_i)_i - \gamma \left[ \frac{\partial L_2}{\partial w_i} \right]_{w_i=i}$$

$$(w_i)_{i+1}^* = w_i - \gamma (2w_i)$$

If  $w_i$  is very small and  $\gamma$  is also very small  
 $\gamma \cdot (2w_i) \ll \gamma$

Hence in each iteration change in  $(w_i)_{i+1}^*$  in  $L_2$  is very less as compared to  $L_1$  where change is more.

Hence Rate of convergence is different in both.

All this happen due to  $\frac{\partial L_2}{\partial w_i}$  is a constant i.e. +1 or -1 whereas as  $\frac{\partial L_1}{\partial w_i}$  is not constant but reduces with increase in no. of iteration. Hence  $L_2$  converges very slow.

$L_2$  regul. doesn't change the value of  $w_i$  from one iteration to another.  
 $L_1$  " continues to constantly reduce  $w_i$  towards  $w_i^*$ .

$L_1$  creates sparsity due to its constant slope while slope of  $L_2$  decreases with iteration.