

35-1

Solving Optimization Problems

Differentiation

We have seen optimization problems like:

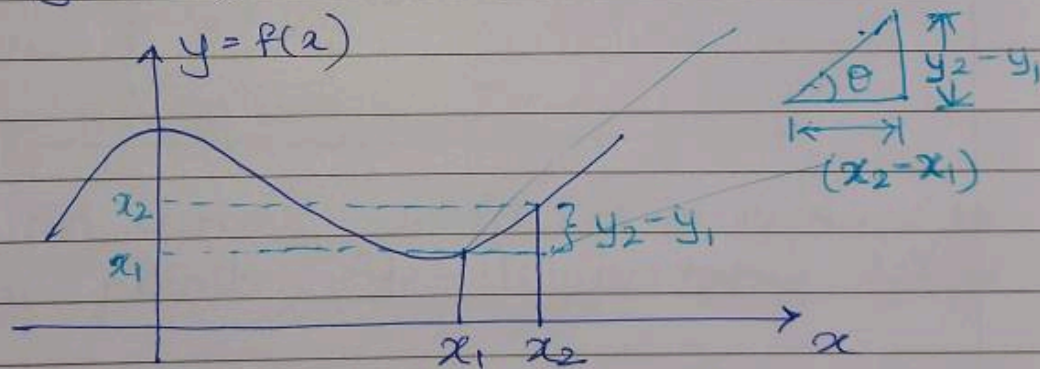
- PCA
- Logistic Regression
- Linear u

Concepts used to solve these problems

- Differentiation
- Maxima & Minima

Single Variable differentiation

$$y = f(x)$$



$$\frac{dy}{dx} = \frac{df}{dx} = y' = f'$$

= differentiation of y w.r.t. x

Intuitively it means

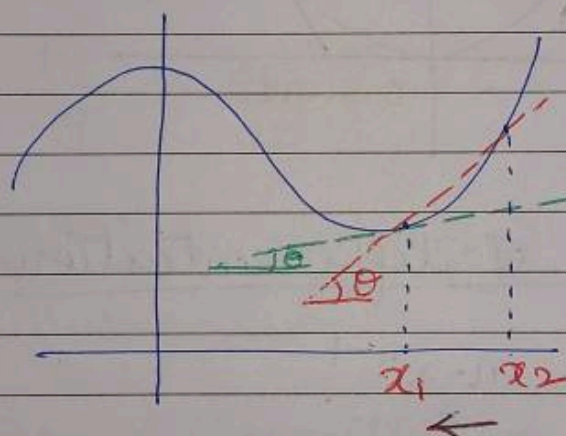
- How much does y change as x changes
- Rate of change of y as x changes

From the graph, how much does 'y' changes as 'x' changes.

$$= \frac{y_2 - y_1}{x_2 - x_1} = \frac{\Delta y}{\Delta x} \quad \left. \begin{array}{l} \text{Rate of change of } y \\ \text{in vicinity of } x. \end{array} \right\} = \tan \theta$$

Mathematically,

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



As x_2 moves
closer to x_1 ;

Tangent of $f(x)$
at x_1

As $x_2 \rightarrow x_1$; $\Delta x \rightarrow 0$

As $x_2 \rightarrow x_1$, the secant becomes tangent

Tangent is the hypotenuse that we obtain as $\Delta x \rightarrow 0$

$$\tan \theta = \frac{dy}{dx} ; \text{ as } \Delta x \rightarrow 0 ; \frac{\Delta y}{\Delta x} = \frac{dy}{dx}$$

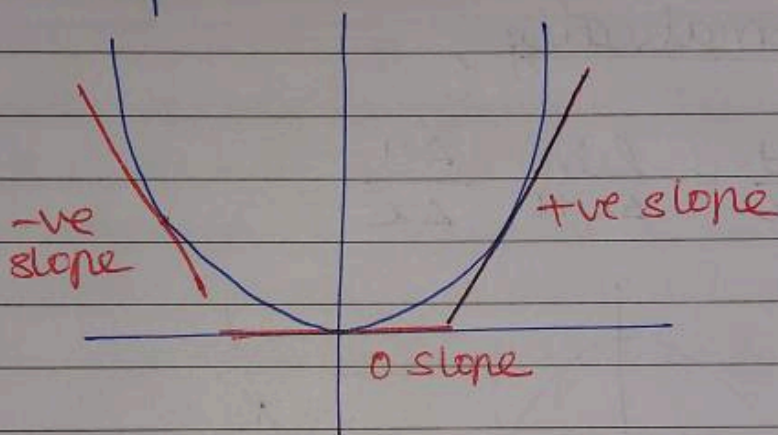
$$\frac{dy}{dx} = \text{Slope of tangent to } f(x)$$

$$\left[\frac{dy}{dx} \right]_{x_1} = \text{slope of the tangent to } f(x) \text{ at } x = x_1$$

slope is +ve; $0 < \theta < 90^\circ$

slope is -ve; $90 < \theta < 180^\circ$

Sum it up:



Basic Rules of differentiation

$$\frac{d}{dx} (x^n) = n \cdot x^{n-1}$$

$$\frac{d}{dx} (c) = 0$$

$$\frac{d}{dx} (\log x) = \frac{1}{x}$$

Chain Rule :

$$\frac{d}{dx} [f(g(x))] = \frac{df}{dg} \cdot \frac{dg}{dx}$$

$$f(g(x)) = (a-bx)^2$$

$$\text{where } g(x) = (a-bx)$$

$$f(x) = x^2$$

$$\therefore \frac{d}{dx} f(g(x)) = \frac{df}{dg} \cdot \frac{dg}{dx}$$

$$\text{Now } \frac{dg}{dx} = \frac{d}{dx} (a-bx) = -b.$$

$$\text{Now } \frac{df}{dg} ; \text{ Now take } g(x) = z$$

$$\therefore f(g(x)) = z^2$$

$$\therefore \frac{df}{dg} = \frac{dz^2}{dz} = 2z = 2(a-bx)$$

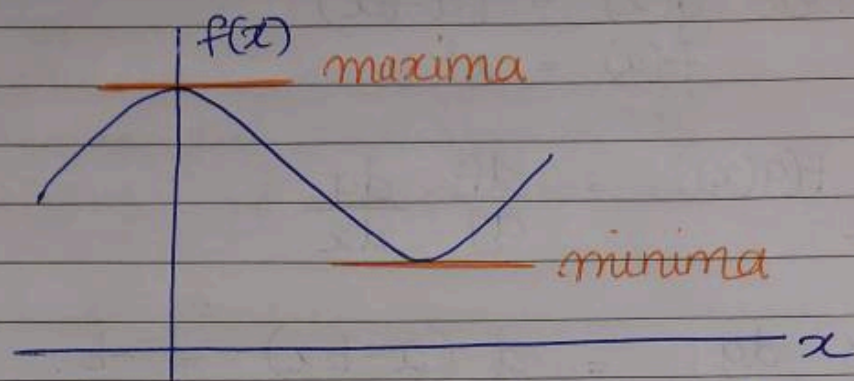
$$\therefore \frac{d}{dx} f(g(x)) = \underline{-2b(a-bx)}$$

35.2 Online differentiation tool

Visit derivative-calculator.net

symbolab.com

wolframalpha.com \rightarrow plots the derivative

35-3Minima and Maxima

At minima and maxima ; slope = 0

Eg: ① $f(x) = x^2 - 3x + 2$

I know that slope = 0 at min & max

$$\therefore \frac{df}{dx} = 0$$

$$\therefore 2x - 3 = 0$$

$$x = 3/2$$

Now is this maxima/minima?

If we evaluate $x = 1.5$

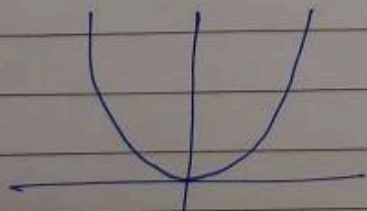
$$\therefore f(1.5) = -0.25$$

Now, we check $f(1) = 0$

$$\text{Now } f(1.5) < f(1)$$

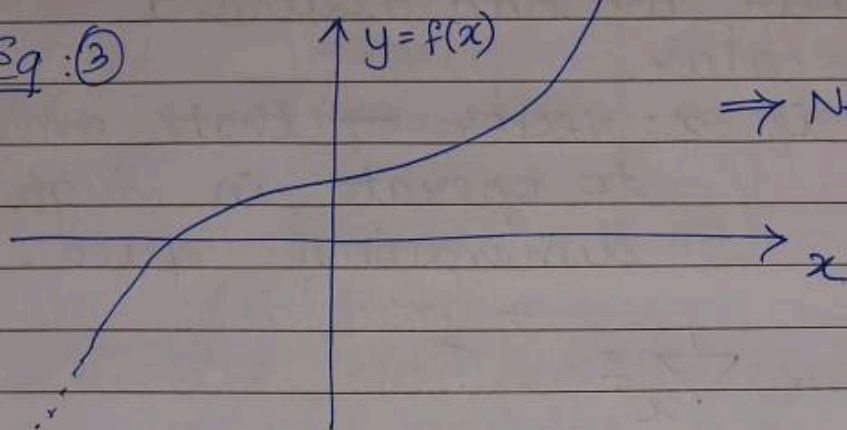
Hence 1.5 cannot be a maxima as there is some other point greater than that.

② $y = f(x) = x^2$ has no maxima;
it has minima at $x = 0$



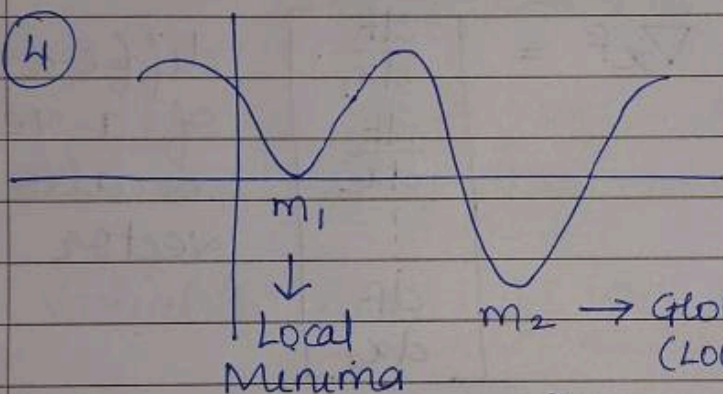
We may/may not have maxima or minima

Eg: (3)



\Rightarrow Neither Maxima
Nor Minima

(4)



m_1 and m_2 are
two minima

(It is co-r.t locality)
in its neighbourhood

$m_2 \rightarrow$ Global Minima
(Lowest amongst all
minima)

(5) $f(x) = \log(1 + \exp(ax)) \Rightarrow$ similar to
Logistic Loss

$$\text{Now } \frac{d}{dx} = \frac{a \exp(ax)}{1 + \exp(ax)} = 0$$

Solving this is not trivial
We should use Gradient descent

35.4 Vector differentiation

Till now, we had assumed
 x : scalar.

what if x : vector \rightarrow that's only way
 to operate in high
 dimensional space.

$$\frac{df}{dx} = \nabla_x f \quad \xrightarrow{\text{grad/del}}$$

Vector

written as

$$\frac{df}{dx_1} = \frac{\partial f}{\partial x_1}$$

since its

partial differentiation

$$\nabla_x f = \begin{bmatrix} \frac{df}{dx_1} \\ \frac{df}{dx_2} \\ \vdots \\ \frac{df}{dx_d} \end{bmatrix} \quad \begin{array}{l} \text{differentia}^n \\ \text{of a vector} \\ \text{is also a} \\ \text{vector} \end{array}$$

$$\therefore \nabla_x f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix} \in \mathbb{R}^d$$

$$f(x) = y = a^T x = \sum_{i=1}^d a_i x_i$$

$$= a_1 x_1 + a_2 x_2 + \dots + a_n x_n$$

$$\frac{\partial f}{\partial x_1} = a_1 ; \frac{\partial f}{\partial x_2} = a_2 ;$$

$$\therefore \nabla_x f = \begin{bmatrix} \partial f / \partial x_1 \\ \partial f / \partial x_2 \\ \vdots \\ \partial f / \partial x_d \end{bmatrix} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_d \end{bmatrix} = a \Rightarrow \text{Vector}$$

② Consider our logistic loss

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

We know that $\langle x_i, y_i \rangle \rightarrow \text{Constants}$ because they are obtained from D_{Train}

Variable here is w

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

chain Rule

For simplicity, we considered $w^T w = w^2$

$$\therefore \frac{d}{dw}(w^2) = 2w \rightarrow \text{Remembers } w \text{ is a vector}$$

We will solve it using Gradient Descent

35.5

Gradient descent algo

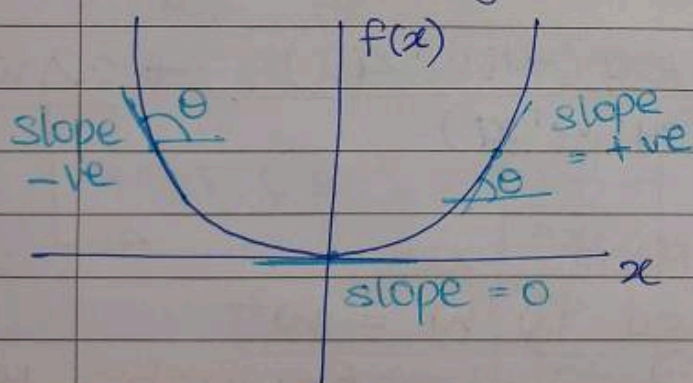
When solving these $\frac{df}{dx} = 0$; $\nabla_x f = 0$ equations are not easy then we prefer Gradient descent algo.

Its an iterative algo.

- ① We randomly select $x_0 \rightarrow$ first guess x^*
- ② That's iteration 1 \rightarrow select x_1
- ③ That's iteration 2 \rightarrow select x_2
- ④ Then iteration $k \rightarrow$ select x_k

Minima/
Maxima

we want to move closer towards x^*

Understanding Geometrically

$$x^* = \operatorname{argmin}[f(x)]$$

$$x^* = \operatorname{argmax}[-f(x)]$$

To find Minima;

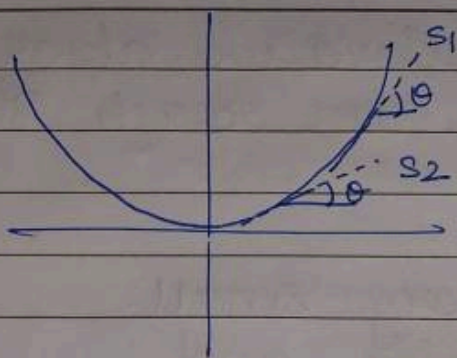
one side :- slope :- +ve

other side :- slope :- -ve

At minima : slope :- 0

- ① Slope changes its sign from positive to Negative at Minima

(2)



slope(s_1) > slope(s_2)
slope reduces as
we move towards
minima

So, to summarize:

① Pick an initial point = x_0 at random

② Let $x_1 = x_0 - \alpha \left[\frac{df}{dx} \right]_{x_0}$

step size = 1 (Typically)

Consider the slope at $s_1 : \rightarrow +ve$

$$x_1 = x_0 - 1 (+ve \text{ value})$$

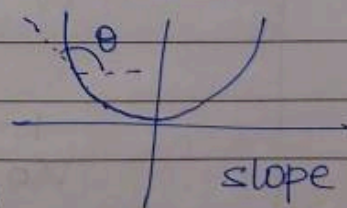
So x_1 is reducing and moving closer to 0.

If my x_0 was on other side

$$\therefore x_1 = x_0 - \alpha \left[\frac{df}{dx} \right]$$

$$x_1 = x_0 - \alpha (-ve \text{ value})$$

Hence x_1 will move close to 0.



slope will
have -ve
value

③ Let $x_2 = x_1 - \alpha \left[\frac{df}{dx} \right]_{x_1}$

I will get x_2 closer to x^*

This is a simple iterative algo.
Going on forward, we reach x_k
such that

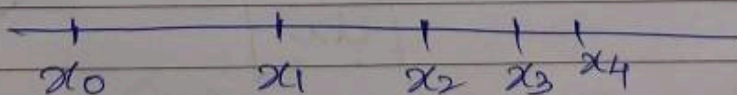
$(x_{k+1} - x_k)$ is very small
then terminate algo and
finalize $x_k = x^*$

If x is a vector then we can
replace $\left[\frac{df}{dx} \right]$ with ∇_x

~~Plz~~ Please Note :

$$\left(\frac{df}{dx} \right)_{x_0} \geq \left(\frac{df}{dx} \right)_{x_1} \geq \left(\frac{df}{dx} \right)_{x_2} \dots$$

In first iteration, we would
make a big jump from $x_0 \rightarrow x_1$.
The difference would slowly reduce



difference reduces gradually

35.6 Learning Rate

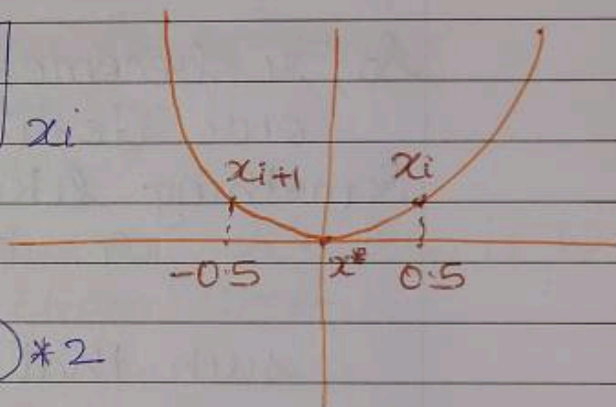
$$x_i = x_{i-1} - \eta * \left(\frac{df}{dx} \right)_{x_{i-1}} \Rightarrow \text{Update Equation}$$

↓
Step Size

Let's say $f(x) = x^2$
and x_i at a particular point = 0.5.

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

Now $\frac{df}{dx} = 2x$



$$\therefore x_{i+1} = 0.5 - (1) * (0.5) * 2$$

$$\underline{x_{i+1} = -0.5}$$

Here we went to other side rather than reaching x^* .
While jumping from x_i to x_{i+1} ,
I jumped over x^*

$$x_{i+2} = x_{i+1} - \eta \left[\frac{df}{dx} \right]_{x_{i+1}}$$

$$= (-0.5) - (1) [2 * (-0.5)]$$

$$\underline{x_{i+2} = 0.5}$$

So, x_{i+2} came back to 0.5.

My x_{i+3} would be again -0.5,

So, I would keep oscillating between 0.5 & -0.5

This is because ' α ' has been kept constant. Due to this, I would never converge towards x^* .

Remedy for this problem:
(oscillation)

- \Rightarrow change $\boxed{\alpha}$ with each iteration
- \Rightarrow one technique: reduce α with each iteration

So, α becomes a function of our iteration.

Something like

$$\alpha = h(i)$$

\uparrow iteration

such that as ~~$x \uparrow$~~ $i \uparrow$; $\alpha \downarrow$

Then, we will eventually converge to right solution.

When we learn Deep Learning, we would know how to modify α more effectively.

35.7 Gradient descent for Linear Regression

Linear Regression:

$$w^* = \underset{w}{\operatorname{argmin}} \sum_{i=1}^m (y_i - \underbrace{w^T x_i}_{\substack{\uparrow \\ \text{not using} \\ \text{intercept} \\ w_0}})^2 \quad \underbrace{\hspace{10em}}_{\substack{\downarrow \\ \text{not using} \\ \text{Regularization} \\ \text{term.}}}$$

for simplicity

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

Here (x_i, y_i) are constants; variable is w
 $\langle x_i, y_i \rangle$ we get from D_{train}

We have to solve

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

① Pick a random vector $w_0 = \langle \dots \dots \rangle$

② $w_1 = w_0 - \underset{\substack{\uparrow \\ \text{will change slightly}}}{\eta} * \sum_{i=1}^m (-2x_i)(y_i - \underline{w_0^T x_i})$

③ $w_2 = w_1 - \underset{\substack{\downarrow}}{\eta} * \sum_{i=1}^m (-2x_i)(y_i - \underline{w_1^T x_i})$

I will compute $w_0, w_1, w_2, \dots, w_k$ as long as w doesn't change much

If $(w_{k+1} - w_k) = \text{very small value}$
then we stop there

There is one problem with this:

$$w_i = w_{i-1} - \alpha \left[\sum_{i=1}^n (-2x_i)(y_i - w_i^T x_i) \right]$$

Every time I go from w_{j-1} to w_j ,
I have to compute summation over
all n ; In M.L. problems, n
may traverse in Millions.

Computing the summation requires
us to go through whole dataset.

It takes a lot of time if n
is large (which is not uncommon in ML)

Hack around it: Gradient descent
is converted to

Stochastic Gradient descent

85.8 Stochastic Gradient Descent (SGD)

It says instead of using all ' n '
points, use only k points ($k \leq n$)
We can pick a random set of k points

If we do iteratⁿ with k points, the
solution we get with Gradient descent
is same as that with Stochastic G.D.

$$w_{G.D.}^* = w_{SGD}^*$$

Stochastic = Probabilistic = Randomized.
(Randomly select K)

One observation

Gradient Descent : $n = 1M$; require 100 Iteration to converge to x^*

Stochastic G-D ; $K = 1000$; require 500 Iteration
~~~~~ ;  $K = 10$ ; require 1000 ~~~~~

Iterations increase as  $K \downarrow$

$K$  = Random points we pick at each iteration

Remember:  $K$  should be diff for every iteration.

$K$  is often called as Batch Size since we pick batch of random points

Often times  $K=1 \Rightarrow$  called SGD

When  $1 < K \ll n \Rightarrow$  ~~~~ Batch SGD  
with Batch Size  $K$



### 35.9 Constrained Optimization and PCA

Formulation of PCA: 
$$\max_{U} \frac{1}{n} \sum_{i=1}^n (U^T x_i)^2$$
  
such that  $U^T U = I$

This is similar to  $\max_x f(x)$   
such that  $g(x) = c$

The function which we want to maximize is called as Objective function; any condition on it is called as Constraint.

PCA: 
$$\max_{U} \frac{1}{n} \sum_{i=1}^n (U^T x_i)^2$$
 objective func<sup>n</sup>  
such that  $U^T U = I$  constraint

The optimization problem with constraint are called Constrained optimization prob

### General constraint optimization

$\max_x f(x) \rightarrow$  Objective func<sup>n</sup>

such that  $g(x) = c \rightarrow$  Equality constraint  
and  $h(x) \geq d \rightarrow$  Inequality constraint

We are not putting  $\leq$  constraint because we can easily convert it to  $\geq$  constraint  
i.e.  $K(x) \leq e \Rightarrow -K(x) \geq -e$

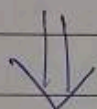


How to solve this problem?

We will modify the problem using Lagrangian Multiplier

So,

$$\rightarrow \underset{x}{\text{Max}} f(x) \text{ such that } g(x) = c \text{ and } h(x) \geq d$$



$$\mathcal{L}(x, \lambda, \mu) = f(x) - \lambda [g(x) - c] - \mu [d - h(x)]$$

$\lambda$  &  $\mu$  :- Lagrangian Multipliers  
and  $\lambda \geq 0$  ;  $\mu \geq 0$

To optimize this; we have a method

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

By solving these three equations, we get value of  $x = \tilde{x}$ , this value is same as  $\underset{\uparrow}{\text{Max}} x = \underline{x^*}$

Proof of this is beyond scope of syllabus



PCA :

$$\left. \begin{array}{l} \text{Max}_U \frac{1}{n} \sum_{i=1}^n (v^T x_i)^2 \\ \text{such that } v^T v = 1 \end{array} \right\} \begin{array}{l} \text{same as} \\ \downarrow \\ \text{Max}_U v^T S v \\ \text{such that } v^T v = 1 \end{array}$$

and

$$S = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$$

Covariance  
Matrix of  $X$ 

PCA prob

$$\therefore \text{PCA prob: } \text{Max}_U v^T S v$$

$$\text{such that } v^T v = 1$$

Please note that  $S$  is constant  
because  $S = P(x_i) = \text{Covariance Matrix of } X$

$$\mathcal{L}(u, \lambda) = v^T S v - \lambda(v^T v - 1)$$

$$\frac{\partial \mathcal{L}}{\partial v} = 0 \Rightarrow \frac{\partial}{\partial v} [v^T S v - \lambda v^T v - \lambda] = 0$$

$$\therefore S v - \lambda v = 0$$

$$\therefore S v = \lambda v$$

We have seen this equation

$$\begin{array}{ccc} S v = \lambda v & \rightarrow & \text{definition of} \\ \uparrow & \uparrow & \text{Eigen value and} \\ \text{Covariance} & \text{Vector} & \text{Eigen Vector} \\ \text{variance} & & \end{array}$$

$v$  = Eigen Vector of  $S$

$\lambda$  = Eigen value of  $S$



When we learnt PCA, we wrote the optimization problem, we had said the solution to this would be to take top Eigen Values and Eigen Vectors of Covariance Matrix  $S$ .

The best  $v$  we get is the one with highest value of  $\lambda$  (Eigen Value)

### 35.10 Logistic Regression Revisited

$$w^* = \underset{w}{\operatorname{argmin}} \left[ (\text{logistic loss}) + \lambda w^T w \right]$$

We would try to give an interpreta<sup>n</sup> that  $\lambda w^T w$  is similar to Lagrange Multiplier.

We have to find  $w^*$  which minimizes our logistic loss such that  $w \perp \lambda$ .

So, we rewrite the optimiza<sup>n</sup> problem:

$$w^* = \underset{w}{\operatorname{argmin}} (\text{logistic loss}) \quad \rightarrow \text{objective function}$$

such that  $w^T w = 1 \quad \rightarrow \text{equality constraint}$

This is constrained optimiza<sup>n</sup> problem.

We solve this through Lagrange Multiplier



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

can be also written  
as  $+ \lambda (w^T w - 1)$

$$= (\text{logistic loss}) - \lambda + \lambda w^T w$$

If we see, these two terms are same as our optimization problem

Regulariza<sup>n</sup> can be thought of as an equality constraint as a solution through Lagrangian Multipliers.

35.11 Why does L1 Regularization create sparsity?

We know that

Logistic Regression  $\rightarrow$  L1 regularization creates sparsity in  $w$  as compared to L2 Regularization

Sparsity :  $w = \langle w_1, w_2, \dots, w_d \rangle$   
more of the  $w_i = 0$



L2 Regulariza<sup>n</sup>

$$\min_w [\text{loss} + \lambda \|w\|_2^2]$$

L1 Regulariza<sup>n</sup>

$$\min [\text{loss} + \lambda \|w\|_1]$$

Since the loss is present  
on both the terms,

we can ignore it for now

Similarly  $\lambda$  is a <sup>hyper</sup>parameter; but  
its a fixed value; hence it can  
be ignored as well.

Hence its equivalent to considering  
L1 & L2 Reg

$$\text{Minimize } [w_1^2 + w_2^2 + \dots + w_d^2]$$

$$\text{Minimize } [|w_1| + |w_2| + \dots + |w_d|]$$

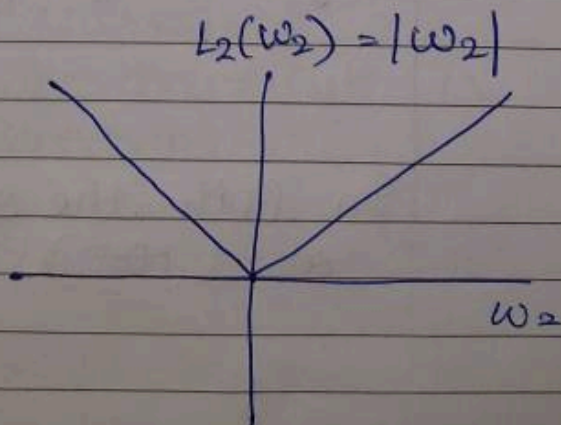
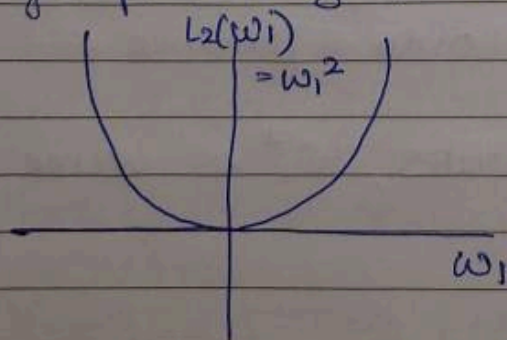
So we see that  $w_1$   $\Rightarrow$  same here.  
occurs only in 1 term

So lets see from  $w_1$   
perspective

From  $w_1$  perspective It would be  
 $\min(w_1^2)$

$\min(|w_1|)$

Graphically seeing





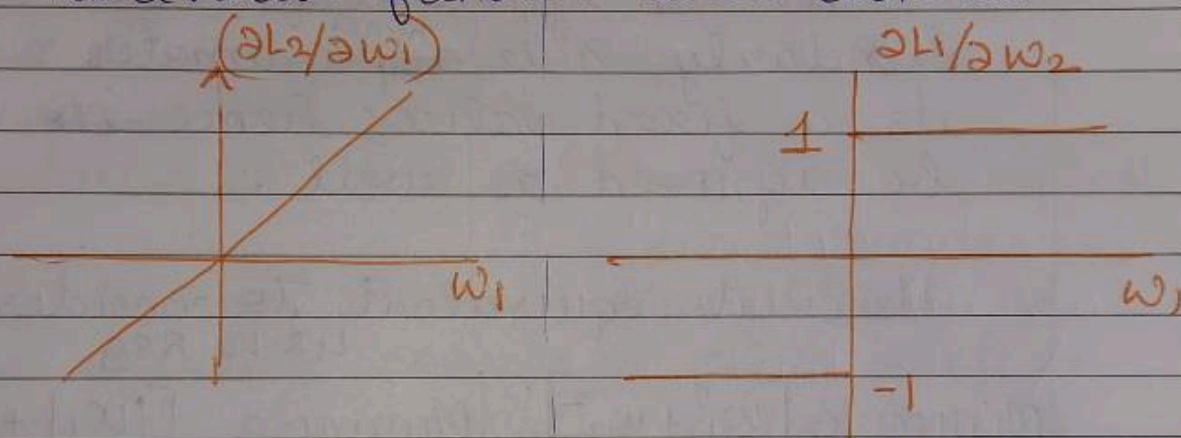
$$L_2(w_1) = w_1^2$$

$$L_1(w_1) = |w_1| = \begin{cases} w_1 & \text{if } w_1 > 0 \\ -w_1 & \text{if } w_1 < 0 \end{cases}$$

$$\frac{\partial L_2}{\partial w_1} = 2 * w_1$$

$$\frac{\partial L_1}{\partial w_1} = \begin{cases} +1 & ; \text{ if } w_1 > 0 \\ -1 & ; \text{ if } w_1 < 0 \end{cases}$$

Derivative function will look like



As a part of Gradient descent

$$(w_1)_{j+1} = (w_1)_j - \eta \left[ \frac{\partial L_2}{\partial w_1} \right]_{(w_1)_j} \quad (w_1)_{j+1} = (w_1)_j - \eta \left[ \frac{\partial L_1}{\partial w_1} \right]_{(w_1)_j}$$

Let  $(w_1)_j$  be a +ve value.

$$(w_1)_{j+1} = (w_1)_j - \eta (2 * w_j) \quad (w_1)_{j+1} = (w_1)_j - \eta (1)$$

→ only difference ←

In both the scenarios,  $w_1^*$  is same, i.e. 0 at  $(0,0)$



L2L1

Assume that at  $i^{\text{th}}$  iterations

$$(w_1)_j = 0.05 \text{ (very small)}$$

Let's say  $\alpha$  is a small value (~~say  $\alpha = 1$~~ )

Due to this

$$(2 * w_j) \rightarrow \text{small}$$

If  $\alpha$  is very small then

$$\alpha * 2 * w_j \rightarrow \text{very small}$$

$$\text{Let's say } \alpha = 0.01 \text{ \& } w_{ij} = 0.5$$

The subtraction term would be

$$\alpha (2 * w_j) = 0.01 \times 2 \times 0.5 \quad \alpha = 0.01$$

$$= \underline{0.001}$$

$$w_{1(j+1)} = 0.05 - 0.001 \quad w_{1(j+1)} = 0.05 - 0.01$$

$\downarrow$   
 $w_{ij}$

Here as the iteration increases

the derivative changes The derivative is  
and  $w_{1(j+1)}$  comes constant  
closer to  $w^*$

derivative becomes smaller

Speed of convergence decreases with every  $\alpha \times 10^m$

Speed of convergence is same everytime



L2

Since the derivative ~~is~~ becomes smaller with every iteration, L2 Regularizer doesn't change the value of  $w_i$  much, from one iteration to other

L1

Since the derivative is constant, L1 Regularizer continues to constantly reduce the value of  $w_i$  towards 0.

Hence L1 creates sparsity due to constant slope/gradient

chance is Less

← Hence, the chance that  $w_i = 0$  at end of  $K$  iterations are more

L1 will remove more of useless feature by making their weight to 0