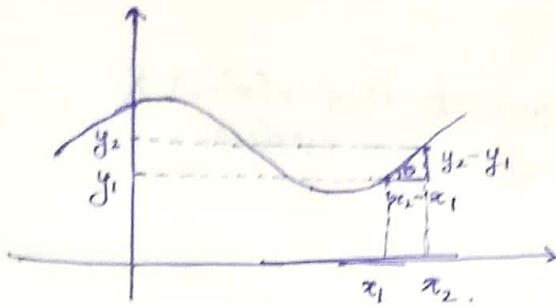


# Solving Optimization Problems

## Differentiation



$\frac{dy}{dx}$  — diff<sup>n</sup> of  $y$  wrt  $x$ .  
— intuitively means rate of change of  $y$  wrt  $x$ .

$$\frac{y_2 - y_1}{x_2 - x_1} = \frac{\Delta y}{\Delta x} = \tan \theta$$

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

Geometrically,  $\left. \frac{dy}{dx} \right|_{x=x_1}$  = slope of tangent to  $f(x)$  at  $x=x_1$ .

~~\*\*\*~~  $y = f(x)$

"All of deep-learning is just chain rule" — Jeff Dean.

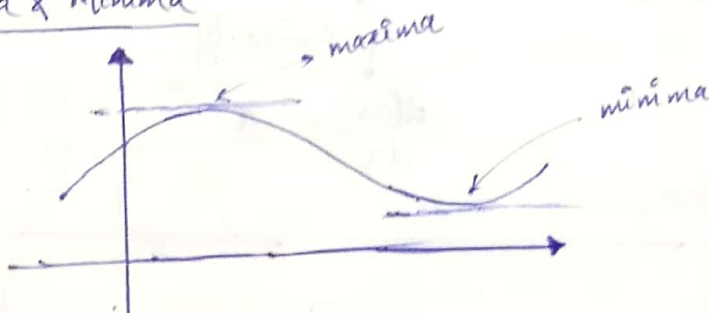
## Online derivative calculators

↳ derivative-calculator.net.

symbolab.com — step by step breakdown.

wolframalpha.com — plots.

## Maxima & Minima



At maxima & minima,  
slope = 0.

A function may/mayn't have maxima & minima, or even may have more than one maxima & minima.

Ex  $f(x) = x^2 - 3x + 2$ .

$$\frac{df}{dx} = \text{slope} = 0 \Rightarrow 2x - 3 = 0$$

$$\Rightarrow x = 3/2 \text{ — whether maxima or minima?}$$

Random check  $\rightarrow f(3/2) < f(0)$ . So, minima.

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

$$\frac{df}{dx} = \frac{a \cdot \exp(ax)}{1 + \exp(ax)} = 0 \quad \left\} \text{ solving this not so trivial}$$

We'll learn some computational techniques to find minima & maxima  
↳ Gradient Descent

Vector differentiation: grad

$$x: \text{vector} \quad x = \langle x_1, x_2, \dots, x_d \rangle$$

$$y = a^T x \quad a = \langle a_1, a_2, \dots, a_d \rangle$$

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

$$\frac{df}{dx} = \nabla_x f \rightarrow \text{grad or del.}$$

$$y = a^T x = \sum_{i=1}^d a_i x_i = a_1 x_1 + a_2 x_2 + \dots + a_d x_d$$

$$\nabla_x f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_d \end{bmatrix} = a$$

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

$\frac{\partial f}{\partial x_i}$  = partial diff.

$$\text{so, } \nabla_x (a^T x) = a$$

similarity

$$\frac{d(a^T x)}{dx} = a$$

Gradient Descent Algorithm

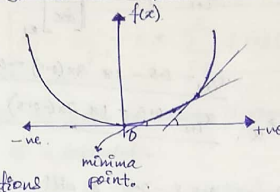
↳ Iterative algorithm

say the optimizat<sup>n</sup> eq<sup>n</sup> is  $x^* = \text{argmin}_x f(x)$

Initially, start off with some initial guess. ( $x_0$ ).

Iteratively update  $x$  value & after say  $k$  iterations, it'll converge towards the optimal value.

core geometric idea of gradient descent



Left side of minima  $\rightarrow$  slope (-)ve.  
Right side of minima  $\rightarrow$  slope (+)ve.  
At minima  $\rightarrow$  slope is 0.

Observations

- \* slope changes its sign from (-)ve to (+)ve at minima.
- \* slope gradually increases as we move from left to right.
- \* slope gradually decreases as we move from right to left.

gradient descent

① pick an initial pt  $x_0$  at random.

②  $x_{i+1} = x_i - r \cdot \left[ \frac{df}{dx} \right]_{x_i}$  causes  $x$  to move towards minima irrespective of which side of minima it is currently.

step-size (+)ve.

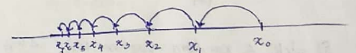
if  $|x_{k+1} - x_k|$  is very small, then we terminate the iterations

& declare  $x^* = x_k$

As the pt moves more & more closer towards minima, the rate of convergence decreases.

$$x_{i+1} = x_i - r \left[ \frac{df}{dx} \right]_{x_i} \rightarrow \text{update function}$$

$$\left[ \frac{df}{dx} \right]_{x_0} > \left[ \frac{df}{dx} \right]_{x_1} > \left[ \frac{df}{dx} \right]_{x_2}$$

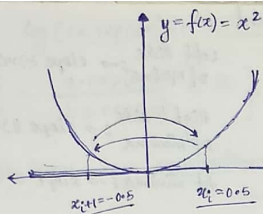


Learning Rate ( $r$ ).

there is a problem of oscillation if we keep the learning rate constant.

let  $y = f(x) = x^2$ ,  $x_0 = 0.5$  &  $r = 1$ .

$$\text{so, } \frac{df}{dx} = 2x$$



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

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

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

So, we basically keep on oscillating between 0.5 & -0.5, & will never converge to minima.

Remedy for oscillation → change  $r$  with each iteration.

one method → reduce  $r$  with each iterat<sup>n</sup>.

$$r = h(i) \quad i = \text{iterations} \quad i \uparrow, r \downarrow$$

In deep learning we'll cover more efficient ways of handling  $r$ .

### Gradient descent for Linear Regression

Optimized function

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

{for simplicity, leaving out  $w_0$  & regularization term}

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

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

① pick a random vector  $w_0 = \langle w_1, w_2, \dots, w_d \rangle$ .

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

$$w_2 = w_1 - r \times \sum_{i=1}^n (-2x_i)(y_i - w_1^T x_i)$$

$$w_0, w_1, w_2, \dots, w_k, w_{k+1}, \dots$$

stop when  $w_{k+1} - w_k$  is a vector of very small values.

But there is a small issue with this gradient descent algo.

$$w_j = w_{j-1} - r \times \sum_{i=1}^n (-2x_i)(y_i - w_j^T x_i) \rightarrow \text{update function}$$

For the update to occur, it requires to go through the whole dataset of  $n$  datapoints.

If  $n$  is large ( $n \approx 1$  million), each converging step will take a significant amount of time.

In order to handle this problem we have stochastic gradient descent.

### Stochastic Gradient Descent (SGD)

↳ the most important optimization algorithm.

$$b1D: w_{j+1} = w_j - r \times \sum_{i=1}^n (-2x_i)(y_i - w_j^T x_i)$$

$$sgd: w_{j+1} = w_j - r \times \sum_{i=1}^k (-2x_i)(y_i - w_j^T x_i), \quad 1 \leq k \leq n$$

So, SGD basically picks a random set of  $k$ -pts to update the weight vector.

$k = \#$  of random pts that we pick at each iteration for updating. (also called batch size)

$$w_{old}^* = w_{sgd}^* \leftarrow \text{proven}$$

b1D: 100 iterations to converge. ( $n=1M$ )

SGD:  $k=1000$ , takes  $>100$  iterations to converge. (500).

SGD:  $k=100$ , (1000) iterations to converge.

$k$ : batch size in SGD.

often times  $k=1 \rightarrow$  sgd.

$n \ll k \rightarrow$  batch SGD.



## Constrained Optimization

General form of a constrained optimization is as follows -

$$\begin{aligned} \max_x f(x) & \quad \text{Objective function} \\ \text{s.t. } g(x) &= c \quad \text{equality constraint} \\ & \quad \& \quad h(x) \geq d \quad \text{inequality constraint.} \end{aligned}$$

## Lagrangian Multipliers

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

$\lambda$  &  $\mu$  are Lagrangian Multipliers.

$$\lambda, \mu \geq 0$$

The value of  $x$  found by solving  $\rightarrow \frac{\partial L}{\partial x} = 0, \frac{\partial L}{\partial \lambda} = 0, \frac{\partial L}{\partial \mu} = 0$  is same as solving the original constrained optimization.

Revisiting the PCA problems we have

$$\begin{aligned} \max_u \frac{1}{n} \sum_{i=1}^n (u^T x_i)^2 & \rightarrow \frac{1}{n} \sum_{i=1}^n (u^T x_i)(u^T x_i) \\ \text{s.t. } u^T u &= 1 \end{aligned}$$

$$\text{So, } \max_u u^T S u, \text{ s.t. } u^T u = 1$$

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

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

$$\Rightarrow 2(Su - \lambda u) = 0$$

$$\Rightarrow Su = \lambda u \quad \text{def'n of eigen value \& eigen vector.}$$

$u_i$  - eigen vector  
 $\lambda_i$  - eigen value of  $S$

$$= \frac{1}{n} \sum_{i=1}^n (u^T x_i)(x_i^T u)$$

$$= \frac{1}{n} u^T \left[ \sum_{i=1}^n (x_i x_i^T) \right] \cdot u$$

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

$\downarrow$   
co-variance matrix

## Logistic Regression Revisited

$$w^* = \arg \min_w (\text{logistic loss}) + (\lambda w^T w) \quad \text{regularizer}$$

Regularizer is nothing but Lagrangian multiplier.

We can write logistic regression as constrained optimization problem as follows  $\rightarrow$

$$w^* = \arg \min_w (\text{logistic loss}) \quad \text{s.t. } w^T w = 1$$

$\downarrow$   
objective function

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

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

$\downarrow$   
same as L.R. with regularization

\* So, regularization can be thought as imposing an equality constraint.

## Why L1-regularization creates sparsity

$$\frac{L_2}{\min_w \text{loss} + \lambda \|w\|_2^2}$$

only difference in both is  $\|w\|_2^2$  &  $\|w\|_1$   
So will focus on that.

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

considering only 1 component for simplicity

$$\min_{w_1} (w_1^2) \quad \downarrow \quad L_2(w_1)$$

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

$$\text{So, } (w_1)_{j+1} = (w_1)_j - r(2 * w_{1,j})$$

$$r = 0.01, \quad w_{1,j} = 0.05$$

$$2 * r * w_{1,j} = 0.001 \quad \text{small}$$

Becomes further smaller on convergence.  
So  $L_2$  req doesn't change value after few iterations.

$$\frac{L_1}{\min_w \text{loss} + \lambda \|w\|_1}$$

$$\min_{w_1, w_2, \dots} (|w_1| + |w_2| + \dots + |w_d|) \quad \downarrow \quad L_1(w_1)$$

$$\min_{w_1} (|w_1|) \quad \downarrow \quad L_1(w_1)$$

$$\frac{\partial L_1}{\partial w_1} = \begin{cases} 1, & w_1 > 0 \\ -1, & w_1 < 0 \end{cases}$$

say  $w_1$  is (+)ve.

$$\text{Then, } (w_1)_{j+1} = (w_1)_j - (r) * 1$$

$$r * 1 = 0.01 \quad \text{small}$$

So converges faster than  $L_2$ , hence more likely to converge towards 0.

\* broadient is constant - major advantage.