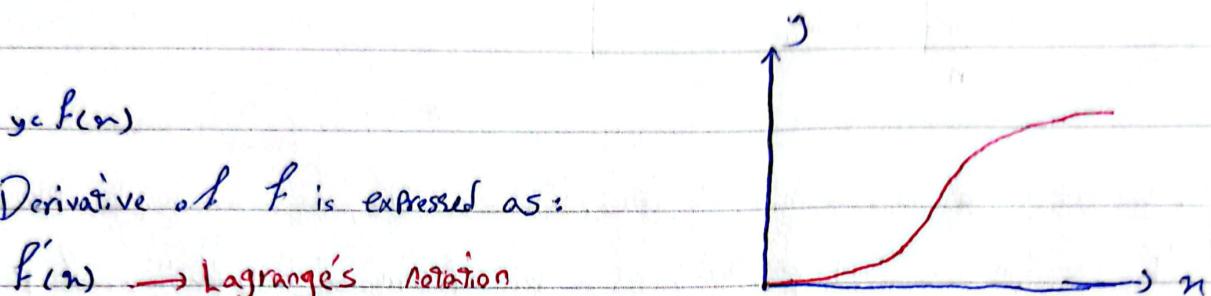


Derivatives: Instantaneous rate of change of a function

Instantaneous rate of change: A measure of how fast the relation between two variable is changing at any point \rightarrow Derivative

* Derivative of a func at a point is the slope of the tangent at that point

* Max & Min of a func occurs at one of the points where derivative is zero.

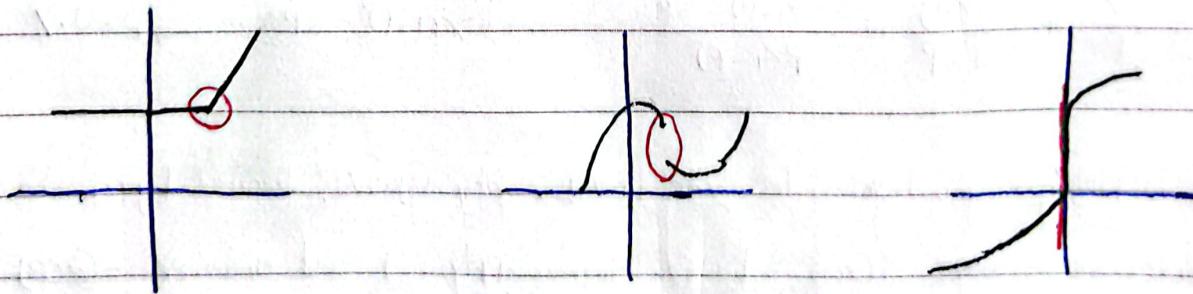


or $\frac{dy}{dx} = \frac{d}{dx} f(x)$ → Leibniz's notation

* If $g(x) = f^{-1}(x) \rightarrow g'(x) = \frac{1}{f'(x)}$ (derivative of inverse funcs)

* $f(x) = e^x \quad f'(y) = \log(y)$ (with e as base) $\frac{d}{dy} f^{-1}(y) = \frac{1}{y}$

Non-Differentiable functions



Corners / Cusps

jump Discontinuity

vertical tangents

* Derivatives are used for optimization (finding max or min value of a func) \rightarrow
Useful for finding the model that best fits your dataset

Ex: Find x that would minimize this func $(x-a)^2 + (x-b)^2$



$$\frac{d}{dx} [(x-a)^2 + (x-b)^2] \Rightarrow 2(x-a) + 2(x-b) = 0 \rightarrow 2x = a+b \rightarrow x = \frac{a+b}{2}$$

The square loss : Minimize $(x-a_1)^2 + (x-a_2)^2 + \dots + (x-a_n)^2$

$$\text{Sol: } x = \frac{a_1 + a_2 + \dots + a_n}{n}$$

dataset

Coin Toss : We are gonna toss a coin 10 times. If the results are seven heads followed by three tails, we win or we lose. We have a biased coin with probability of P for head and $(1-P)$ for tails. Find P

$$g(P) = P^7(1-P)^3 \rightarrow \text{maximize } g(P) \rightarrow \frac{dg}{dP} = 7P^6(1-P)^3 + 3P^7(1-P)^2(-1)$$

$$= P^6(1-P)^2(7-10P) = 0 \rightarrow P = 0 \Rightarrow \text{the coin would always land in tails}$$

$$P = 1 \Rightarrow \text{the coin would always land in heads}$$

$$P = 0.7 \checkmark$$

another way: if $g(P)$ is maximized $\rightarrow \log(g(P))$ is minimized

$$\log(g(P)) = \log(P^7(1-P)^3) = 7\log(P) + 3\log(1-P) = G(P)$$

$$\frac{dG(P)}{dP} = \frac{7}{P} + \frac{3(-1)}{1-P} = \frac{7(1-P) - 3P}{P(1-P)} = 7(1-P) - 3P = 0 \rightarrow P = 0.7$$

* $-G(P)$ is log loss - a useful loss func in ML. Why negative? because $\log(x)$ when $0 < x < 1$

is negative so we want $-G(P)$ to be pos so we apply $(-)$ and minimize $-G(P)$

why log?

1- Derivative of products is hard, derivative of sums is easy

2- Product of lots of tiny things is tiny

Functions of two variables: We have tangent plane in 2-d $f(x,y) = x^2 + y^2$

Partial Derivatives: Find Partial derivative of f with respect to x and y

$$\frac{\partial f}{\partial x} \rightarrow \text{Treat } y \text{ as a constant} \rightarrow \frac{\partial f}{\partial x} = f_x = 2x$$

$$\frac{\partial f}{\partial y} \rightarrow \text{Treat } x \text{ as a constant} \rightarrow \frac{\partial f}{\partial y} = f_y = 2y$$

Gradient: Collection of all the partial derivatives with respect to all the variables in

the function \rightarrow ^{notable} $\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix} = \begin{bmatrix} 2x \\ 2y \end{bmatrix}$

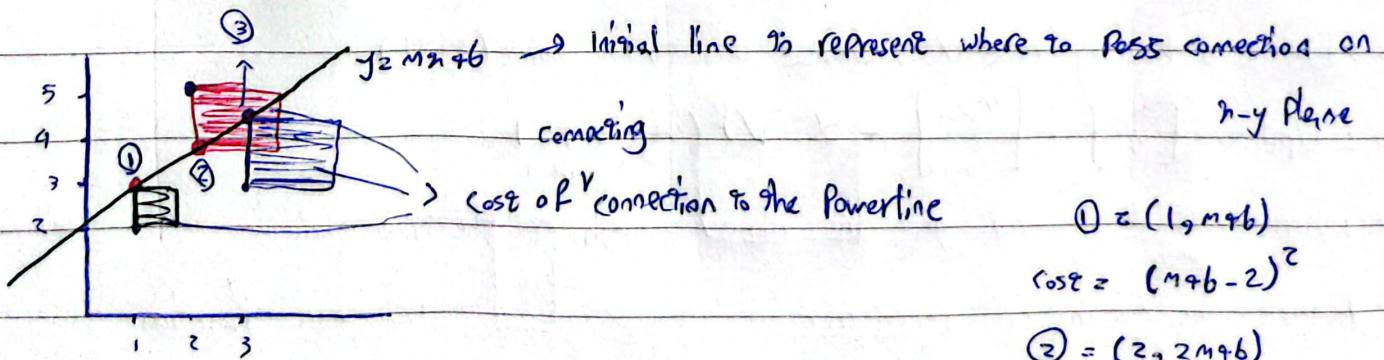
* In functions with two variables, the minimum point is where the tangent plane is parallel to the floor. It's the same with maximum.

If we are working with Partial derivatives, the min point (or max) is where both slopes are zero

$$\text{e.g.: } f(x,y) = x^2 + y^2 \rightarrow \text{Find min/max} \rightarrow \frac{\partial f}{\partial x} = 0 \rightarrow 2x = 0 \rightarrow x = 0$$

$$\frac{\partial f}{\partial y} = 0 \rightarrow 2y = 0 \rightarrow y = 0 \quad (0,0) \rightarrow \text{min}$$

Linear Regression: Analytical approach



Goal: Find m, b such that you minimize sum of square cost $③ (3, 3m+b)$

$$\text{cost} = (3m+b - 3)^2$$

Subject:

Date:

$$\text{cost func} = (m+b-2)^2 + (2m+b-5)^2 + (3m+b-3)^2$$

$$\rightarrow E(m,b) = 14m^2 + 3b^2 + 38 + 12mb - 42m - 20b$$

$$\frac{\partial E}{\partial m} = 28m + 12b - 42 = 0$$

$$\rightarrow m = 0.5 \quad b = \frac{7}{3} \quad E\left(\frac{1}{2}, \frac{7}{3}\right) \approx 4.167$$

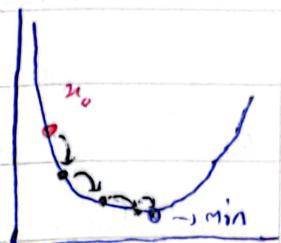
$$\frac{\partial E}{\partial b} = 6b + 12m - 20 = 0$$

Gradient Descent: Func: $f(w)$ goal: find min of func

a: Define a learning rate α , choose a starting point w_0

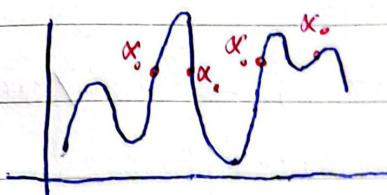
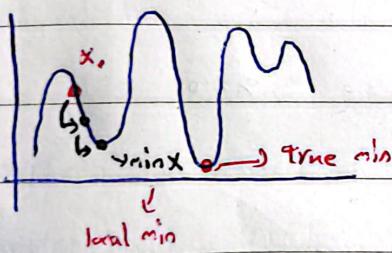
b: Update: $w_k = w_{k-1} - \alpha f'(w_{k-1})$

c: Repeat step 2 until you are close enough to the true min



Drawbacks: Find min!

Solution: Take multiple initial starting points



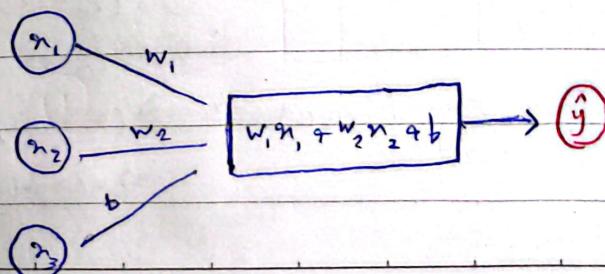
GD with two vars:

a: Define a learning rate α . choose a starting point (w_0, y_0)

b: Update $\begin{bmatrix} w_k \\ y_k \end{bmatrix} = \begin{bmatrix} w_{k-1} \\ y_{k-1} \end{bmatrix} - \alpha \nabla f(w_{k-1}, y_{k-1})$

c: Repeat step 2 until you are close enough to the true minimum

Regression with a Perceptron



Prediction Function

$$\hat{y} = w_1x_1 + w_2x_2 + b$$

$$\text{Loss Function} \rightarrow \text{error}$$

$$L(y, \hat{y}) = \frac{1}{2} (y - \hat{y})^2$$

Main Goal

Find w_1, w_2, b that gives the least error

Samane

Date:

Subject:

Using gradient descent

$$L(y, \hat{y}) = \frac{1}{2} (y - \hat{y})^2 \quad \hat{y} = w_1 n_1 + w_2 n_2 + b$$

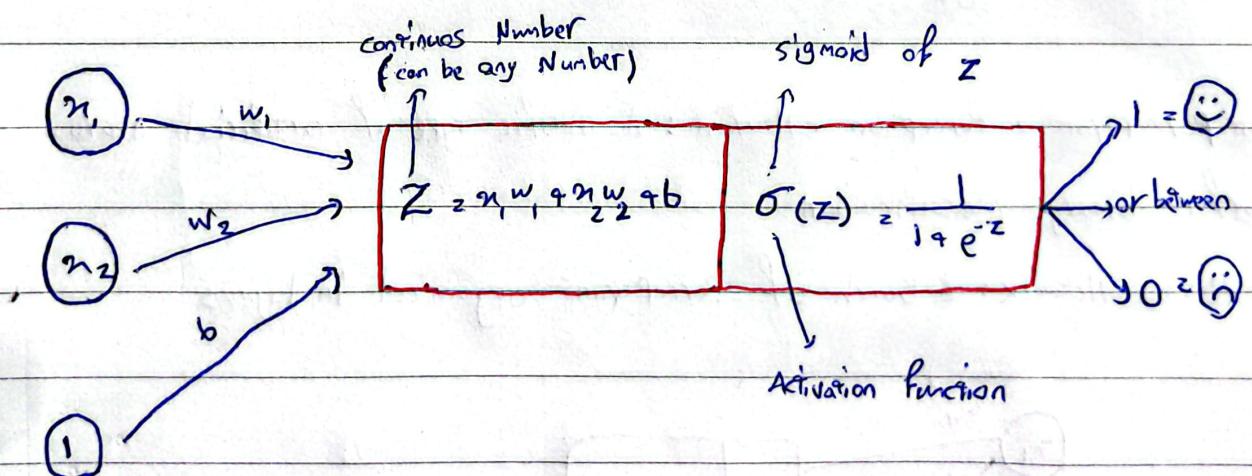
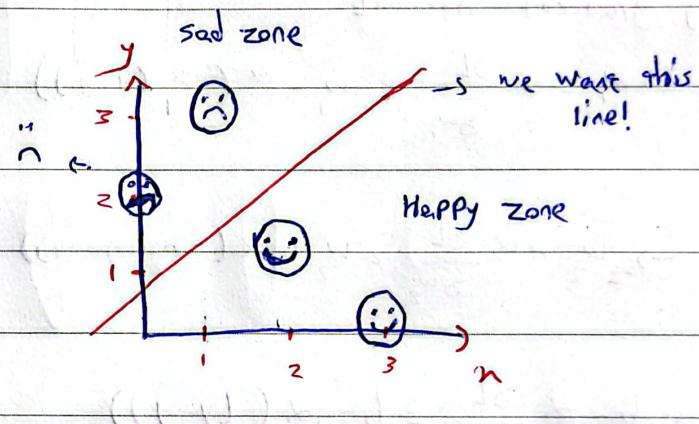
$$w_1 \rightarrow w_1 - \alpha \frac{\partial L}{\partial w_1} \rightarrow w_1 - \alpha (-n_1 (y - \hat{y})) \quad \frac{\partial L}{\partial w_1} = -(y - \hat{y}) n_1$$

$$w_2 \rightarrow w_2 - \alpha \frac{\partial L}{\partial w_2} \rightarrow w_2 - \alpha (-n_2 (y - \hat{y})) \quad \frac{\partial L}{\partial w_2} = -(y - \hat{y}) n_2$$

$$b \rightarrow b - \alpha \frac{\partial L}{\partial b} \rightarrow b - \alpha (-(y - \hat{y})) \quad \frac{\partial L}{\partial b} = -(y - \hat{y})$$

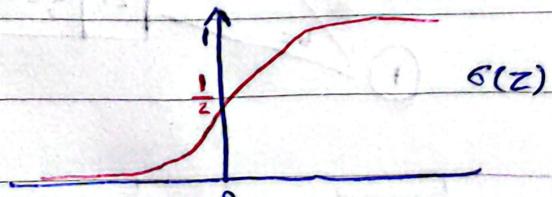
Classification with a Perceptron

sentence	n_1	n_2	y	Mood
nnn	3	0	0	:)
yy	0	2	1	:)
nyyy	1	3	1	:)
ynyn	2	1	0	:)

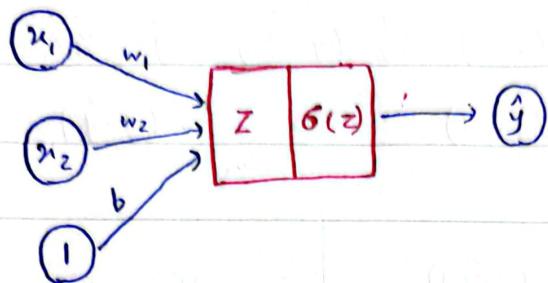


* Sigmoid func takes the entire number line and crunches it into interval zero, one

$$\sigma(z) = \frac{1}{1 + e^{-z}} = (1 + e^{-z})^{-1}$$



$$\frac{d}{dz} \sigma(z) = \sigma(z)(1 - \sigma(z))$$



Prediction func:

$$\hat{y} = \sigma(w_1 n_1 + w_2 n_2 + b)$$

Loss func:

$$L(y, \hat{y}) = -y \ln(\hat{y}) - (1-y) \ln(1-\hat{y})$$

★ we use Log-loss:

Goals:

- the math works out really nicely

Find w_1, w_2, b that give \hat{y} with the

- classification has a probabilistic nature

least error

Using gradient descent

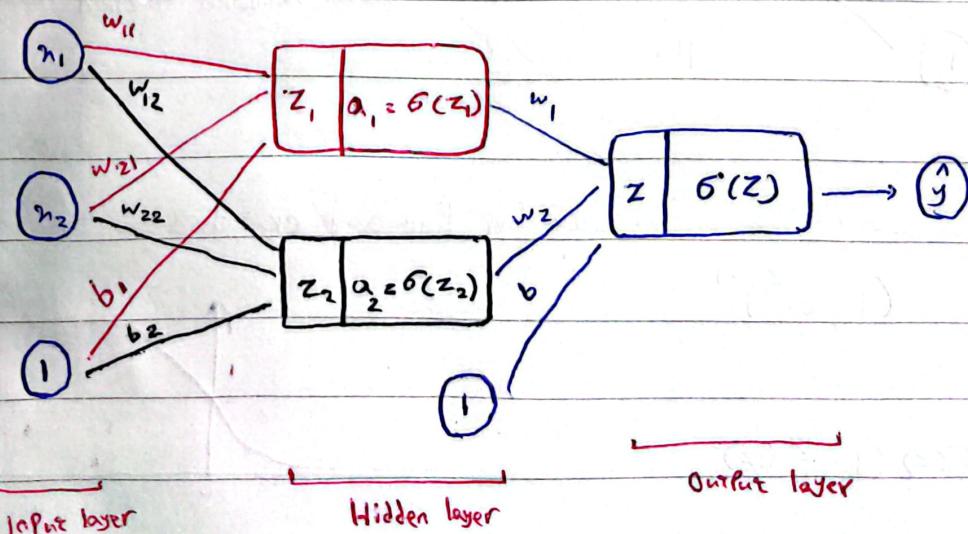
$$w_1 \rightarrow w_1 - \alpha \frac{\partial L}{\partial w_1} = w_1 - \alpha (-n_1(y - \hat{y}))$$

$$w_2 \rightarrow w_2 - \alpha \frac{\partial L}{\partial w_2} = w_2 - \alpha (-n_2(y - \hat{y}))$$

$$b \rightarrow b - \alpha \frac{\partial L}{\partial b} = b - \alpha (-(y - \hat{y}))$$

★ Perceptron ≈ Perception + Neuron: A specific type of artificial neuron designed to perform binary classification

★ Neural Network: A bunch of perceptrons organized in layers



L1

$$w_{ij} \rightarrow w_{ij} + \alpha_{ij} w_i a_i (1 - a_i) (y - \hat{y}) \quad b_i \rightarrow b_i + \alpha w_i a_i (1 - a_i) (y - \hat{y})$$

L2

$$w_i \rightarrow w_i + \alpha a_i (y - \hat{y})$$

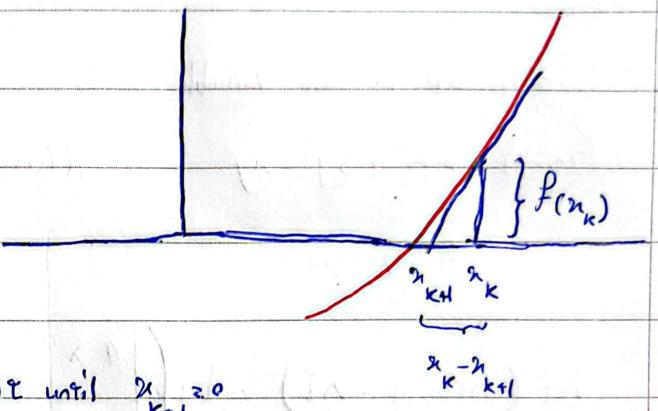
$$b \rightarrow b + \alpha (y - \hat{y})$$

Newton's Method: is used to find the zeros of a function

$$\Rightarrow \frac{f(x_k)}{x_k - x_{k+1}} = f'(x_k)$$

$$\rightarrow \frac{f(x_k)}{f'(x_k)} = x_k - x_{k+1}$$

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)} \rightarrow \text{update it until } x_{k+1} \approx 0$$



NM for optimization: minimize $g(x)$ = find zeros of $g'(x)$

$$f(x) = g(x) \quad f'(x) = g'(x)$$

1) Start with some x_0 .

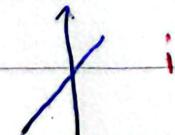
2) Update

$$x_{k+1} = x_k - \frac{g(x)}{g'(x)}$$

3) Repeat until you find the candidate for minimum

Curvature

First derivative

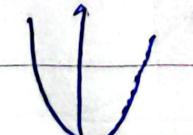


increasing

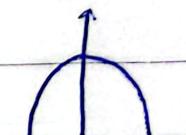


Decreasing

Second derivative



concave up



concave down

* $f'(x_0) = 0 \text{ & } f''(x_0) > 0 \rightarrow x_0 \text{ is a local min}$

* $f'(x_0) = 0 \text{ & } f''(x_0) < 0 \rightarrow x_0 \text{ is a local max}$

Hessian Matrix : $H(x,y) =$

$$\begin{bmatrix} f_{xx}(x,y) & f_{xy}(x,y) \\ f_{yx}(x,y) & f_{yy}(x,y) \end{bmatrix}$$

* in most cases:

$$f_{xy}(x,y) = f_{yx}(x,y)$$

Concave up in two variables

$$f(x,y) = 2x^2 + 3y^2 - xy$$

$$H(0,0) = \begin{bmatrix} 4 & -1 \\ -1 & 6 \end{bmatrix}$$

Finding eigen values

$$\det(H(0,0) - \lambda I) = \det \left(\begin{bmatrix} 4-\lambda & -1 \\ -1 & 6-\lambda \end{bmatrix} \right) = \lambda^2 - 10\lambda + 23 \quad \begin{array}{l} \lambda_1 = 6.41 \\ \lambda_2 = 3.59 \end{array}$$

→ both are $> 0 \rightarrow (0,0)$ is a min

Concave down in two variables

$$f(x,y) = -2x^2 - 3y^2 - xy + 15$$

$$H(0,0) = \begin{bmatrix} -4 & -1 \\ -1 & -6 \end{bmatrix}$$

$$\det(H(0,0) - \lambda I) = \lambda^2 + 10\lambda + 23 \quad \begin{array}{l} \lambda_1 = -3.59 \\ \lambda_2 = -6.41 \end{array}$$

→ both are $< 0 \rightarrow (0,0)$ is a max

* if $\lambda_1 > 0, \lambda_2 < 0$ or $\lambda_1 < 0, \lambda_2 > 0 \rightarrow$ we can't conclude anything!

or $\lambda_1 = 0$

Newton's method for 2 vars

$$1 \text{ var} \quad u_{k+1} = u_k - \frac{f'(u_k)}{f''(u_k)} \quad \text{or} \quad u_k = f(u_k) - \frac{f'(u_k)}{f''(u_k)}$$

$$2 \text{ vars} \quad \begin{bmatrix} u_{k+1} \\ y_{k+1} \end{bmatrix} = \begin{bmatrix} u_k \\ y_k \end{bmatrix} - H^{-1}(u_k, y_k) \cdot \nabla f(u_k, y_k)$$

2×2 2×1