

Week 3 Optimizing In Neural Networks

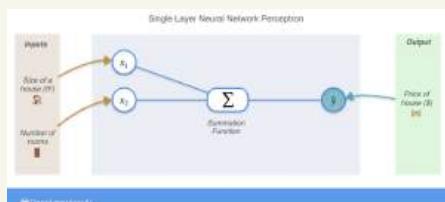
Regression with a perceptron



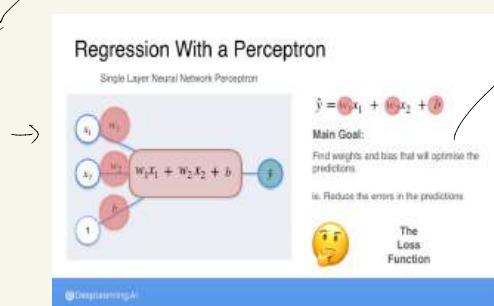
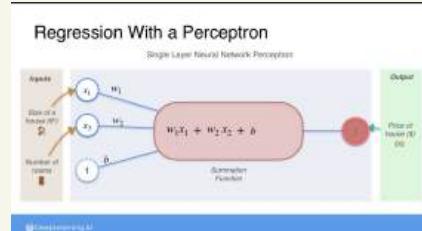
The goal is to find the best possible line
we have 1 feature here
what if we add another feature

Size of a house (m²)	Number of rooms	Price of house (\$)
1000	2	50,000
2000	4	100,000
3000	6	150,000

Now this becomes more complicated. What if we had 10 features or 100.
This is where the perceptron comes in.



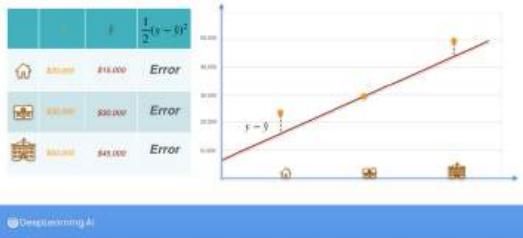
Now what is happening in the summation function
each input is multiplied by a weight to determine
how important it is to the output.
let's say we have w_1, w_2 & a bias term that gets
added to this.



Therefore
 \hat{y} is represented
as a function of
 w_1, w_2, b

The loss function tells us
now far we are from a
prediction.

Mean Squared Error



In this $(y - \hat{y})$ is the distance from accurate prediction or the error.

Now why do we take $(y - \hat{y})^2$

- 1) to quantify big outliers
- 2) the error could be +ve or -ve these numbers when summed can be zero or even small w.r.t. which becomes harder to compute.
- 3) $\frac{1}{2}$ is simply to remove the 2 after differentiating

$$\hat{y} = w_1x_1 + w_2x_2 + b, \text{ Loss function: } L(y, \hat{y}) = \frac{1}{2}(y - \hat{y})^2$$

Main goal \rightarrow Find w_1, w_2, b that gives \hat{y} with the least error

To optimize w_1, w_2, b you need gradient descent

$$w_1 \rightarrow w_1 - \alpha \frac{\partial L}{\partial w_1}, \quad w_2 \rightarrow w_2 - \alpha \frac{\partial L}{\partial w_2}, \quad b = b - \alpha \frac{\partial L}{\partial b}$$

Some initial starting values

$$\text{Prediction function: } \hat{y} = w_1x_1 + w_2x_2 + b \quad \text{Loss function: } L(y, \hat{y}) = \frac{1}{2}(y - \hat{y})^2$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial b} > \text{chain rule} \quad \frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_1} > \text{rule} \quad \frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_2}$$

$$\begin{aligned} \frac{\partial L}{\partial \hat{y}} &= \frac{1}{2} \times 2x(-1) \cdot (y - \hat{y}) & \frac{\partial \hat{y}}{\partial b} &= 1 & \frac{\partial \hat{y}}{\partial w_1} &= x_1 & \frac{\partial \hat{y}}{\partial w_2} &= x_2 \\ &= -(y - \hat{y}) \end{aligned}$$

$$\begin{aligned} \frac{\partial L}{\partial b} &= -(y - \hat{y}) \cdot 1 & \frac{\partial L}{\partial w_1} &= \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_1} = -(y - \hat{y})x_1 & \frac{\partial L}{\partial w_2} &= -(y - \hat{y})x_2 \\ &\rightarrow ① & &\rightarrow ② & &\rightarrow ③ \end{aligned}$$

We use gradient descent

$$\begin{aligned} w_1 &= w_1 - \alpha \frac{\partial L}{\partial w_1} = w_1 - \alpha (-x_1(y - \hat{y})) & \rightarrow & \text{If we do this many times} \\ w_2 &\rightarrow w_2 - \alpha (-x_2(y - \hat{y})) & \rightarrow & \text{we get good values to} \\ b &\rightarrow b - \alpha (-(y - \hat{y})) & \rightarrow & \text{minimize loss} \end{aligned}$$

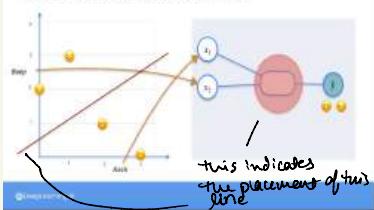
Classification Problem Motivation

Classification Problem Motivation

Review	Ack	Deep	Not
Awkward smiley	2	0	Happy 😊
Boring	0	2	Sad 😢
Awk sleep sleep smiley	1	3	Sad 😢
Awk level smiley	2	1	Happy 😊

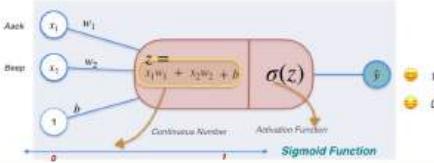
Dataset for the sentiment analysis
 This can be plotted & known as a perceptron

Classification Problem Motivation



Classification With a Perceptron

Single Layer Neural Network Perception



z is a continuous no. that can be anywhere in the number line. We want that in linear regression because we want any no. to be the value but in this we have happy (sad so we want the number to be 1 or 0) So we use a sigmoid function. to turn the entire no. line b/w 1 or 0.

@DeepLearningAI

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

If you look at the graph you will notice everything b/w 0 & 1 for y or $0/1$ whereas the whole number line for input never really touches 1 or 0.

Derivative

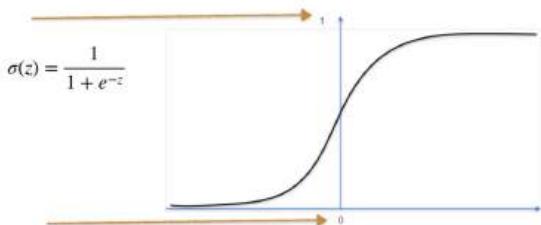
$$\sigma(z) = \frac{1}{1+e^{-z}} \rightarrow \sigma(z) \Rightarrow (1+e^{-z})^{-1} \rightarrow \frac{\partial}{\partial z} \sigma(z) = \frac{\partial}{\partial z} (1+e^{-z})^{-1}$$

$$\rightarrow \frac{\partial}{\partial z} \sigma(z) = -1(1+e^{-z})^{-2} \left(\frac{\partial}{\partial z} (1) + \frac{\partial}{\partial z} (e^{-z}) \right)$$

$$\Rightarrow -1(1+e^{-z})^{-2} (0 + e^{-z} \left(\frac{\partial}{\partial z} (-z) \right))$$

$$\Rightarrow -1(1+e^{-z})^{-2} (e^{-z})(-1)$$

Sigmoid Function



@DeepLearningAI

$$\begin{aligned}\frac{\partial}{\partial z} \sigma(z) &= \cancel{(1+e^{-z})^{-2}} \cdot (e^{-z}) \cancel{(1+e^{-z})^2} \\ &= (1+e^{-z})^{-2} (e^{-z}) = 1/(1+e^{-z})^2 \cdot (e^{-z}) \\ &= \frac{e^{-z}}{(1+e^{-z})^2}\end{aligned}$$

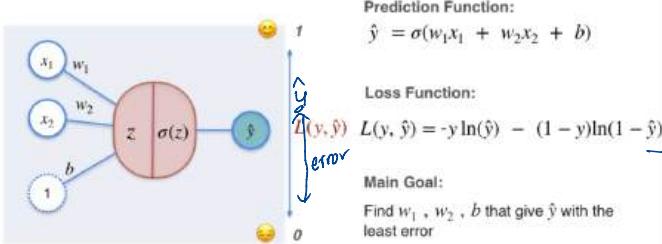
Add & subtract 1

$$\rightarrow \frac{1}{(1+e^{-z})} - \frac{1}{(1+e^{-z})^2} \Rightarrow \frac{1}{1+e^{-z}} \left(1 - \frac{1}{(1+e^{-z})} \right)$$

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

Since it's easy to calculate the differentiation during it is imp for ml.

Classification With a Perceptron



For classification what works best is log loss.

$L(y, \hat{y})$, we use this to update our variables to bring loss down.

This is the loss function.

$L(y, \hat{y})$ is large if y & \hat{y} are far.

$L(y, \hat{y})$ small if y & \hat{y} are close.

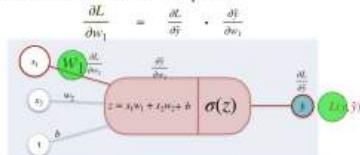
So now for optimal values of w_1, w_2, b we need gradient descent

$$w_1 \rightarrow w_1 - \alpha \frac{\partial L}{\partial w_1} \quad w_2 \rightarrow w_2 - \alpha \frac{\partial L}{\partial w_2} \quad b \rightarrow b - \alpha \frac{\partial L}{\partial b}$$

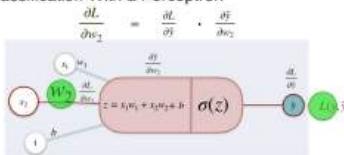
So we need to find all the derivatives & perform gradient descent.

We need to

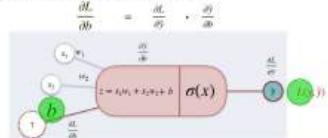
Classification With a Perceptron



Classification With a Perceptron



Classification With a Perceptron



Classification with a Perceptron –

Key Concepts

Network Flow

- Inputs: x_1, x_2
- Linear combin $z = x_1 w_1 + b$
- Activation: $y = \sigma(z) = \frac{1}{1 + e^{-z}}$
- Output: $\hat{y} \in (0, 1)$

Loss Function

- Used for binarification: $L(y, \hat{y}) = -[y \ln(\hat{y}) + (1-y) \ln(1-\hat{y})]$
- Compares predicted probability \hat{y} to actual label y

Why Log Loss Comes After Sigmoid

- Sigmoid gives the probability
- Log loss evaluates how close \hat{y} is to the true label
- The model's goal: minimize this loss

Backpropagation (Gradient Flow)

- Chain rule: $\frac{\partial L}{\partial b} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial b}$
- This helps compute how much to adjust w bias to reduce loss

Classification With a Perceptron

$$\begin{aligned}\frac{\partial L}{\partial b} &= \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial b} \\ \frac{\partial L}{\partial w_1} &= \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_1} \\ \frac{\partial L}{\partial w_2} &= \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_2}\end{aligned}$$

$\hat{y} = \sigma(w_1 x_1 + w_2 x_2 + b)$
 $L(y, \hat{y}) = -y \ln(\hat{y}) - (1-y) \ln(1-\hat{y})$

Classification With a Perceptron

$$\frac{\partial L}{\partial \hat{y}} = \frac{-(y - \hat{y})}{\hat{y}(1-\hat{y})}$$

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

$$\frac{\partial L}{\partial \hat{y}} = \frac{-y}{\hat{y}} + \frac{1-y}{1-\hat{y}}$$

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

$$= \frac{-(y - \hat{y})}{\hat{y}(1-\hat{y})}$$

$$\frac{\partial \hat{y}}{\partial w_1} = \hat{y} (1 - \hat{y}) x_1$$

$$\frac{\partial \hat{y}}{\partial w_2} = \hat{y} (1 - \hat{y}) x_2$$

$$\frac{\partial \hat{y}}{\partial b} = \hat{y} (1 - \hat{y})$$

due to derivative of sigmoid $\rightarrow \sigma'(1-\sigma)$

that's where the $\hat{y} (1 - \hat{y})$ come in

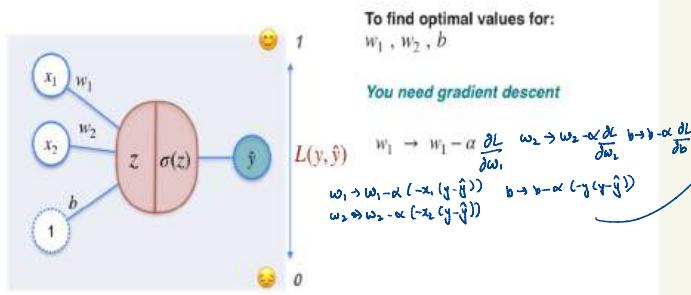
Classification With a Perceptron

$$\begin{aligned}\frac{\partial L}{\partial b} &= \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial b} = \frac{-(y - \hat{y})}{\hat{y}(1-\hat{y})} \\ \frac{\partial L}{\partial w_1} &= \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_1} = \frac{-(y - \hat{y})}{\hat{y}(1-\hat{y})} \\ \frac{\partial L}{\partial w_2} &= \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_2} = \frac{-(y - \hat{y})}{\hat{y}(1-\hat{y})}\end{aligned}$$

$$\frac{\partial L}{\partial b} = -(y - \hat{y}) \frac{\partial L}{\partial \hat{y}} = -(y - \hat{y})x_0 \frac{\partial L}{\partial z} = -(y - \hat{y})x_0$$

so for a good prediction $(y - \hat{y})$ minimal.

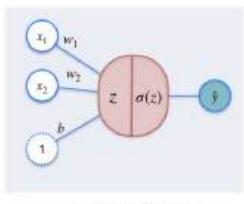
Classification With a Perceptron



Now we saw a perceptron, a neural network is a network of perceptrons in many layers.

Classification Problem Motivation

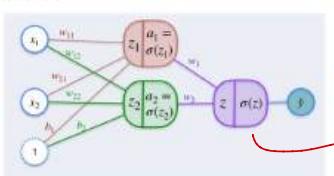
Barkmen	Aack	Beep	Mood
Aack aack aack!	3	0	Happy 😊
Beep beep!	0	2	Sad 😢
Aack beep beep beep!	1	3	Sad 😢
Aack beep aack!	2	1	Happy 😊



DeepLearning.AI

2.2.1 Neural Network

- one input layer
- one hidden layer
- one output layer



- This builds a linear boundary
 However it sometimes is more complicated than just 1 line.
 Then we use multiple perceptrons.

$$\begin{aligned}z_1 &= w_{11}x_1 + w_{12}x_2 + b_1 \\ z_2 &= w_{21}x_1 + w_{22}x_2 + b_2\end{aligned}$$

$$a_1 \sigma(z_1) \quad a_2 \sigma(z_2)$$

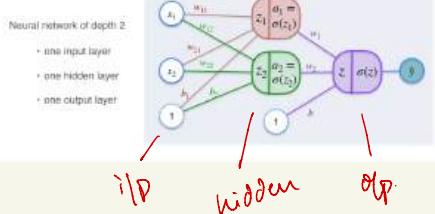
$$\text{then } z = a_1w_1 + a_2w_2$$

$$\sigma(z) \Rightarrow \sigma(z)$$

DeepLearning.AI

We add a bias to output layer

2,2,1 Neural Network



2,2,1 Neural Network

$$a_1 = \sigma(z_1)$$

$$z_1 = x_1 w_{11} + x_2 w_{21} + b_1$$

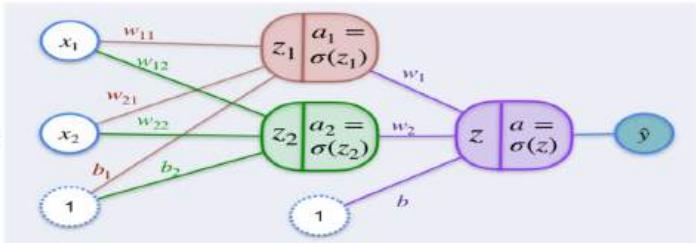
$$a_2 = \sigma(z_2)$$

$$z_2 = x_1 w_{12} + x_2 w_{22} + b_2$$

$$\hat{y} = \sigma(z)$$

$$z = a_1 w_1 + a_2 w_2 + b$$

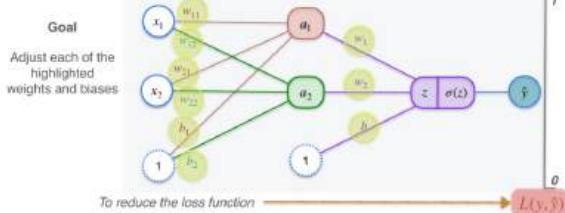
$$L(y, \hat{y}) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$



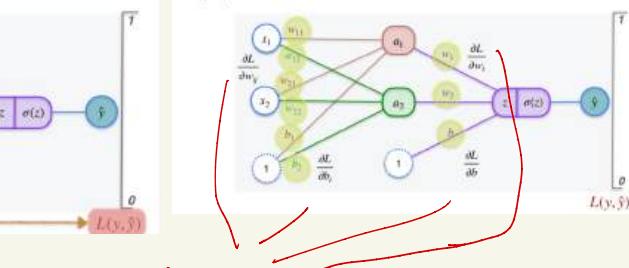
↳ This is the order in which a NN works.

Let's train them now:

2,2,1 Neural Network

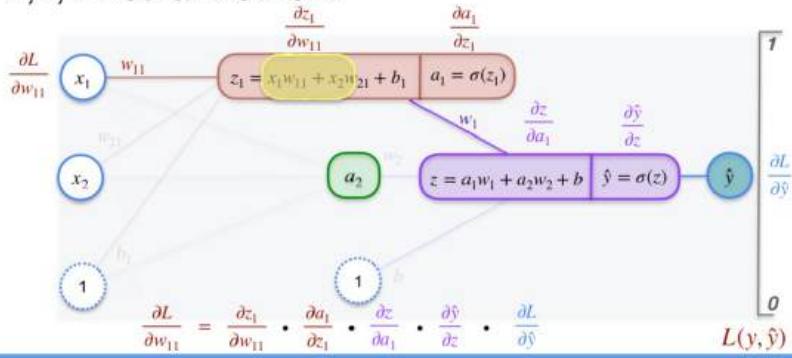


2,2,1 Neural Network



These partial derivatives tell us exactly which direction to move each w & b to reduce the loss function.

2.2.1 Neural Network



DeepLearning.AI

$\frac{\partial L}{\partial \hat{y}}$ depends on $\frac{\partial \hat{y}}{\partial z}$ which depends on $\frac{\partial z}{\partial a_1}$ on $\frac{\partial a_1}{\partial z_1}$ on $\frac{\partial z_1}{\partial w_{11}}$.

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial z_1}{\partial w_{11}} \cdot \frac{\partial a_1}{\partial z_1} \cdot \frac{\partial z}{\partial a_1} \cdot \frac{\partial L}{\partial \hat{y}}$$

$$\Rightarrow x_1 \cdot a_1(1-a_1) \cdot w_1 \cdot \hat{y}(1-\hat{y}) \cdot \frac{-(y-\hat{y})}{\hat{y}(1-\hat{y})}$$

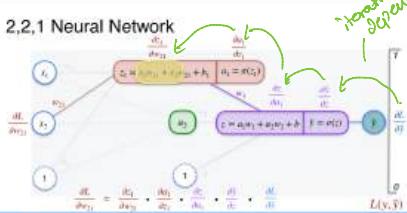
$$\Rightarrow -x_1 w_1 a_1 (1-a_1)(y-\hat{y})$$

$$\begin{cases} L(y, \hat{y}) = -y \log(\hat{y}) - (1-y) \log(1-\hat{y}) \\ \hat{y} = \sigma(z) \\ z = a_1 w_1 + a_2 w_2 + b \\ a_1 = \sigma(z_1) \\ z_1 = x_1 w_{11} + x_2 w_{21} + b_1 \end{cases}$$

Now perform gradient descent.

$$w_{11} \rightarrow w_{11} - \alpha \cdot x_1 w_1 a_1 (1-a_1)(y-\hat{y})$$

for w_{12}

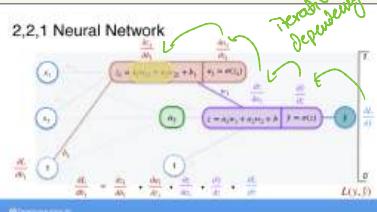


2.2.1 Neural Network

$$\begin{aligned} \frac{\partial L}{\partial w_{21}} &= \frac{\partial z_2}{\partial w_{21}} \cdot \frac{\partial a_2}{\partial z_1} \cdot \frac{\partial z_1}{\partial a_1} \cdot \frac{\partial L}{\partial \hat{y}} \\ \hat{y} &= \sigma(z) \\ z &= a_1 w_1 + a_2 w_2 + b \\ a_1 &= \sigma(z_1) \\ z_1 &= x_1 w_{11} + x_2 w_{21} + b_1 \end{aligned}$$

Perform gradient descent with
 $w_{21} \rightarrow w_{21} - \alpha \cdot x_2 w_2 a_1 (1-a_1)(y - \hat{y})$ to find optimal value of w_{21} that gives the least error

For b_1

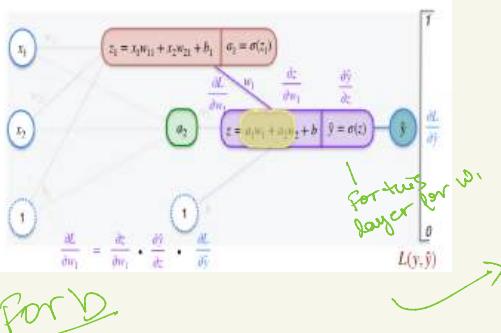


2.2.1 Neural Network

$$\begin{aligned} \frac{\partial L}{\partial b_1} &= \frac{\partial z_1}{\partial b_1} \cdot \frac{\partial a_1}{\partial z_1} \cdot \frac{\partial z}{\partial a_1} \cdot \frac{\partial L}{\partial \hat{y}} \\ \hat{y} &= \sigma(z) \\ z &= a_1 w_1 + a_2 w_2 + b \\ a_1 &= \sigma(z_1) \\ z_1 &= x_1 w_{11} + x_2 w_{21} + b_1 \end{aligned}$$

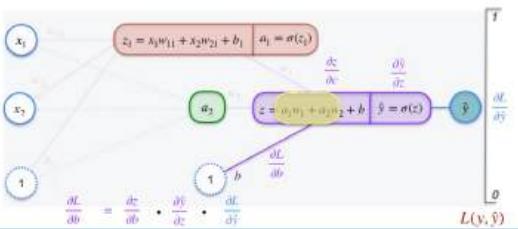
Perform gradient descent with
 $b_1 \rightarrow b_1 - \alpha \cdot a_1 (1-a_1)(y - \hat{y})$ to find optimal value of b_1 that gives the least error

2.2.1 Neural Network



For D

2.2.1 Neural Network



2.2.1 Neural Network

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

$$\frac{\partial L}{\partial w_1} = \frac{\partial z}{\partial w_1} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial L}{\partial \hat{y}}$$

$$\hat{y} = \sigma(z)$$

$$z = a_1 w_1 + a_2 w_2 + b$$

$$\frac{\partial L}{\partial w_1} = a_1 \cdot \hat{\sigma}(1 - \hat{y}) \cdot \frac{(y - \hat{y})}{\hat{\sigma}(1 - \hat{y})} = -a_1(y - \hat{y})$$

Perform gradient descent with

$$w_1 \rightarrow w_1 - \alpha_1(y - \hat{y})$$

to find optimal value of w_1 that gives the least error

2.2.1 Neural Network

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

$$\frac{\partial L}{\partial b} = \frac{\partial z}{\partial b} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial L}{\partial \hat{y}}$$

$$\hat{y} = \sigma(z)$$

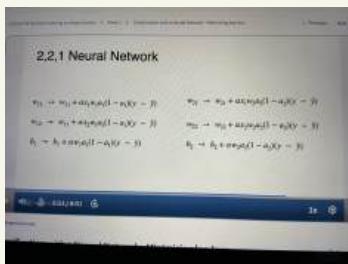
$$z = a_1 w_1 + a_2 w_2 + b$$

$$\frac{\partial L}{\partial b} = 1 \cdot \hat{\sigma}(1 - \hat{y}) \cdot \frac{(y - \hat{y})}{\hat{\sigma}(1 - \hat{y})} = -(y - \hat{y})$$

Perform gradient descent with

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

to find optimal value of b that gives the least error



These are the formulas

2.2.1 Neural Network

$$w_1 \rightarrow w_1 + \alpha a_1(y - \hat{y}))$$

$$w_2 \rightarrow w_2 + \alpha a_2(y - \hat{y}))$$

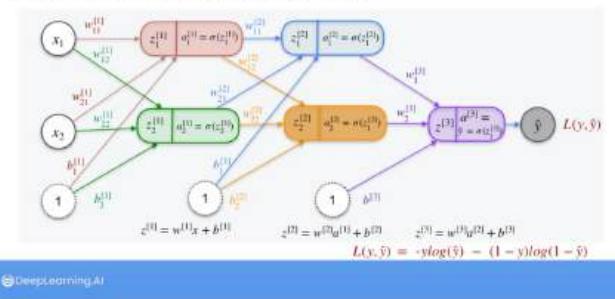
$$b \rightarrow b + \alpha(y - \hat{y}))$$

Front-tuis we are able to derive

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

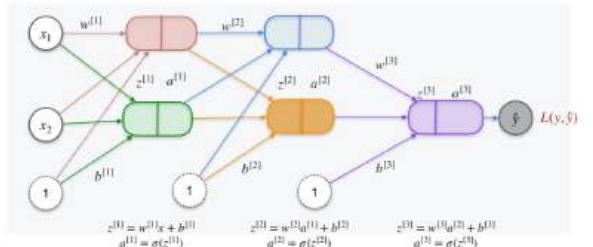
And we try to minimize it

Back Propagation Introduction



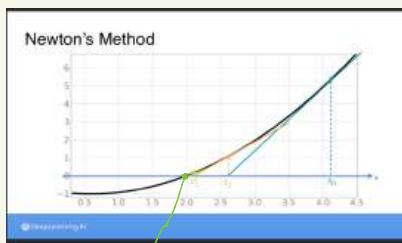
DeepLearning.AI

Back Propagation Introduction



Newton's Method

In principle newton's method is used to find zeroes of a function.

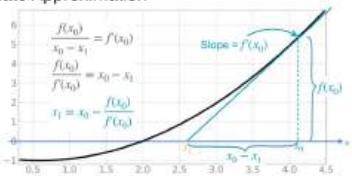


goal is to
find where
function = 0
 $f(x) = 0$

we take x_0 randomly
take the derivative & draw the tangent
see that the place it intersects
the x-axis at x_1 , which is closer to
the zero of function.
we do something at x_1 , & the
slope intersects now at x_2 which
is very close.

Found 0 at x_2

Update Approximation

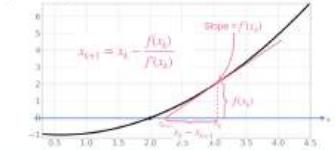


here we can see that

$$\text{Slope} = f'(x_0) = \frac{\text{rise}}{\text{run}} = \frac{f(x_0)}{x_0 - x_1}$$

$$\frac{f(x_0)}{f'(x_0)} = x_0 - x_1, \quad x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}$$

Update Approximation



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

Newton's Method for Optimization

Newton's method

Goal: find a zero of $f(x)$



1) Start with some x_0

2) Update:

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

3) Repeat 2) until you find the root.

DeepLearning.AI

NM for Optimization

Goal: minimize $g(x) \rightarrow$ find zeros of $g'(x)$

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

1) Start with some x_0

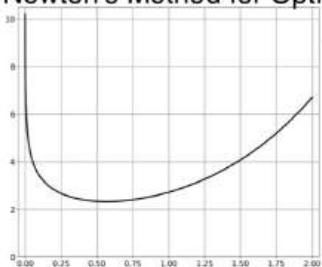
2) Update:

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

3) Repeat 2) until you find the root.

Newton's Method for Optimisation

Newton's Method for Optimization



$$g(x) = e^x - \log(x)$$

$$g'(x) = e^x - \frac{1}{x} \rightarrow f(x) \rightarrow g'(x)$$

$$(g'(x))' = e^x + \frac{1}{x^2} \rightarrow f'(x) \rightarrow (g'(x))'$$

Now to find minimum using NM we need to use it to find the zeroes of the derivative.

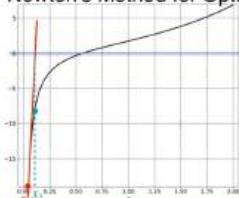
This is the graph of the derivative. we need to find it's zero.

Newton's Method for Optimization

$$g(x) = e^x - \log(x) \quad g'(x) = e^x - \frac{1}{x}$$

Minimum: $x^* = 0.567$

Newton's Method for Optimization



$$\begin{aligned} g(x) &= e^x - \log(x) & f(x) \\ g'(x) &= e^x - \frac{1}{x} \\ \text{Minimum: } x^* &= 0.5671 \\ (g'(x))' &= e^x + \frac{1}{x^2} \\ x_0 &= 0.05 & g'(x_0) \\ x_1 &= x_0 - \frac{g'(x_0)}{(g'(x_0))'} \\ &= 0.05 - \frac{\left(e^{0.05} - \frac{1}{0.05}\right)}{\left(e^{0.05} + \frac{1}{0.05}\right)} = 0.097 \end{aligned}$$

we iterate again with
again we get

x_2

Newton's Method for Optimization

$$g(x) = e^x - \log(x) \quad g'(x) = e^x - \frac{1}{x}$$

Minimum: $x^* = 0.5671$

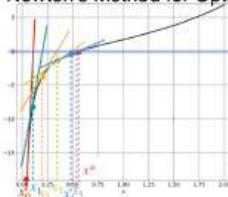
$$(g'(x))' = e^x + \frac{1}{x^2}$$

$x_1 = 0.097$

$$x_2 = x_1 - \frac{g'(x_1)}{(g'(x_1))'}$$

$$= 0.097 - \frac{\left(e^{0.097} - \frac{1}{0.097}\right)}{\left(e^{0.097} + \frac{1}{0.097}\right)} = 0.183$$

Newton's Method for Optimization



$$\begin{aligned} g(x) &= e^x - \log(x) & g'(x) = e^x - \frac{1}{x} \\ \text{Minimum: } x^* &= 0.567 \\ (g'(x))' &= e^x + \frac{1}{x^2} \\ x_0 &= 0.558 \\ x^* &= x_0 - \frac{g'(x_0)}{(g'(x_0))'} \\ &= 0.558 - \frac{\left(e^{0.558} - \frac{1}{0.558}\right)}{\left(e^{0.558} + \frac{1}{0.558}\right)} = 0.567 \end{aligned}$$

on the 5th iteration we reached the minimum value which is $\rightarrow 0.567$.

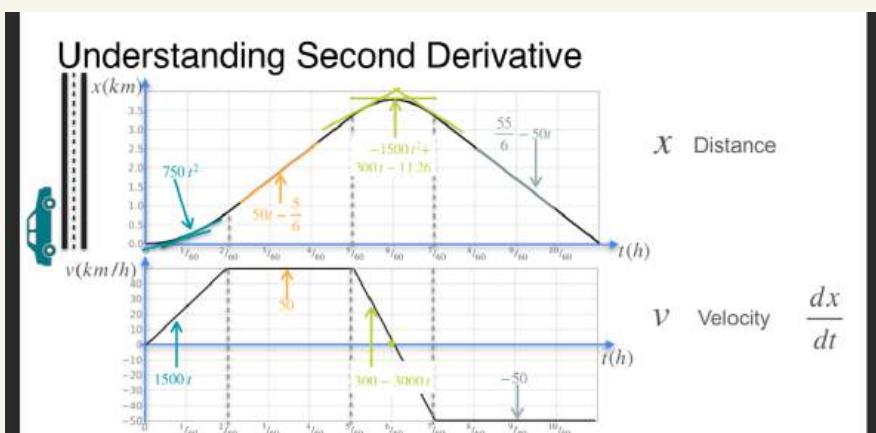
$$\text{Second Derivative} \rightarrow x_{k+1} = x_k - \frac{g'(x_k)}{(g'(x_k))'} \rightarrow \text{second derivative}$$

Lagrange notation: $\frac{d^2 f(x)}{dx^2} = \frac{d}{dx} \left(\frac{df(x)}{dx} \right)$ Lagrange notation: $f''(x)$

Understanding Second derivative

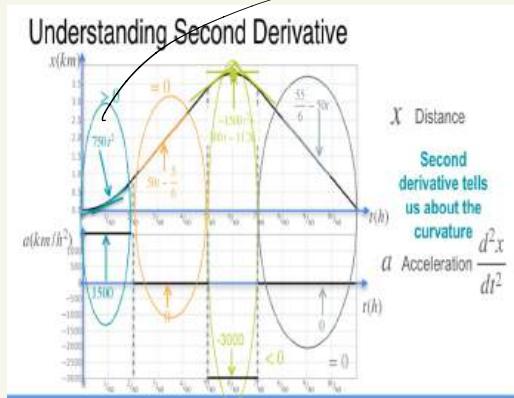
1st derivative \rightarrow velocity for x distance $\rightarrow \frac{dx}{dt}$ \rightarrow rate of change of distance w.r.t time

2nd derivative \rightarrow acceleration for time $t \rightarrow \frac{dv}{dt} = \frac{d^2 x}{dt^2}$ \rightarrow acceleration tells gives increasing velocity
 -ve \rightarrow decreasing velocity.
 $0 \rightarrow$ constant



Comparison
W/ distance
velocity graph

constant increase
in vel \Rightarrow acceleration > 0
at constant velocity.



again from 7-10 \rightarrow The distance decreases constantly
 $\Rightarrow 0$ acceleration

\rightarrow Notice in the first 2 mins when the distance increases at an increasing rate the acceleration

2-5 constant distance covered
 \Rightarrow acceleration $= 0$

5-7 true speed or distance slows down & decreases.
 \Rightarrow acceleration < 0 .

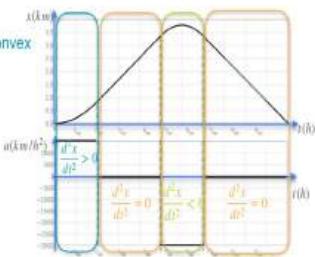
The second derivative gives a measure of the amount the curve deviates from being a straight line.

Curvature

$$\frac{d^2x}{dt^2} > 0 \quad \text{Concave up or convex}$$

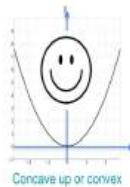
$$\frac{d^2x}{dt^2} < 0 \quad \text{Concave down}$$

$$\frac{d^2x}{dt^2} = 0 \quad \text{Need more information}$$

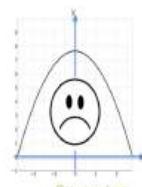


DeepLearningAI

Curvature



Concave up or convex
 $f''(0) > 0$



Concave down
 $f''(0) < 0$

DeepLearningAI

U → smile you are up

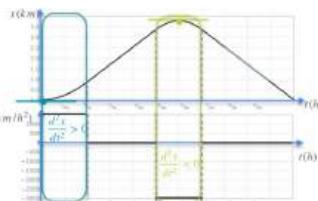
∩ → frown concave down

Second Derivative and Optimization

$$\frac{d^2x}{dt^2} > 0 \quad (\text{Local}) \text{ Minimum}$$

$$\frac{d^2x}{dt^2} < 0 \quad (\text{Local}) \text{ maximum}$$

$$\frac{d^2x}{dt^2} = 0 \quad \text{Inconclusive}$$



DeepLearningAI

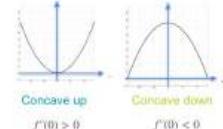
Curvature

First derivative



Increasing
 $f'(t) > 0$

Second derivative



Concave up
 $f''(t) > 0$

increasing function
decreasing function

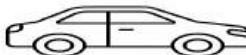
Second Derivative in Optimization - Summary for ML

- First Derivative $f'(x)$**
 - Tells the slope (increasing or decreasing)
 - Used to find stationary points (where slope = 0)
- Second Derivative $f''(x)$**
 - Helps classify those stationary points:

Condition	Interpretation
$f'(x) = 0, f''(x) > 0$	Local Minimum 😊 (concave up)
$f'(x) = 0, f''(x) < 0$	Local Maximum 😢 (concave down)
$f''(x) = 0$	Inconclusive – check more

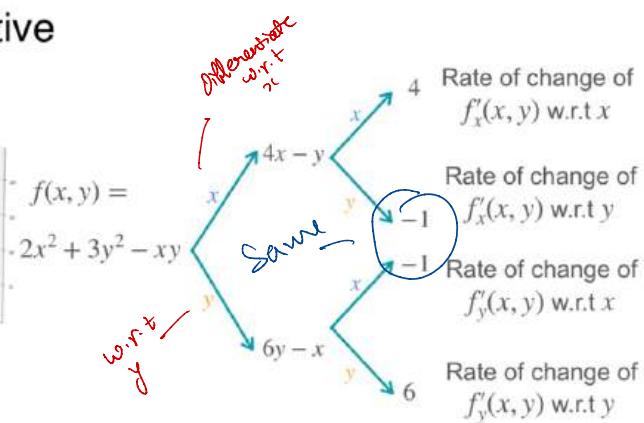
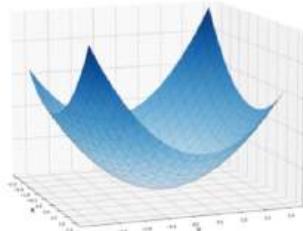
Analogy (Car Example)

- Position $f(x)$
- Velocity $f'(x)$
- Acceleration $f''(x)$
- Positive acceleration = speeding up (curve bends up)
- Negative acceleration = slowing down (curve bends down)



Second derivative in multiple variables

Second Derivative



DeepLearning.AI

What Do These Mean?

Rate of change of $f_x(x, y)$ w.r.t x

Rate of change of $f_x(x, y)$ w.r.t y

Rate of change of $f_y(x, y)$ w.r.t y

Rate of change of $f_y(x, y)$ w.r.t x

Change in the change in the function w.r.t tiny changes in x and y

Same idea as with one variable!

1. Change in the slope along one coordinate axis w.r.t tiny changes along an orthogonal coordinate axis

the first two w.r.t x
w.r.t y

DeepLearning.AI

What Do These Mean?

Rate of change of $f_x(x, y)$ w.r.t x

Rate of change of $f_x(x, y)$ w.r.t y

Rate of change of $f_y(x, y)$ w.r.t y

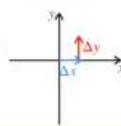
Rate of change of $f_y(x, y)$ w.r.t x

Change in the change in the function w.r.t tiny changes in x and y

Same idea as with one variable!

1. Change in the slope along one coordinate axis w.r.t tiny changes along an orthogonal coordinate axis

2. They are the same!
(in most cases)



Notation

Lafbin's notation

Rate of change of $f'_x(x, y)$ w.r.t x

Rate of change of $f'_x(x, y)$ w.r.t y

Rate of change of $f'_y(x, y)$ w.r.t y

Rate of change of $f'_y(x, y)$ w.r.t x

Lafgrave's notation

$f_{xx}(x, y)$

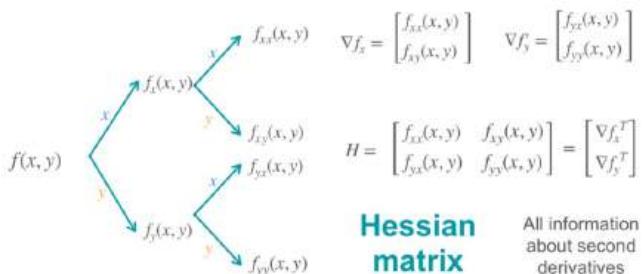
$f_{yy}(x, y)$

$f_{xy}(x, y)$

$f_{yx}(x, y)$

DeepLearning.AI

Hessian Matrix - General Case

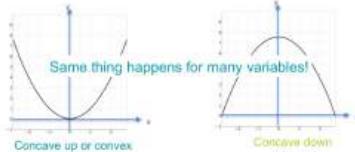


DeepLearning.AI

Question $\rightarrow f(x, y) = x^2 + y^2$ Hessian matrix=?

$$\begin{aligned} \frac{\partial f}{\partial x} &= 2x & \frac{\partial f}{\partial y} &= 2y \\ \frac{\partial^2 f}{\partial x^2} &= 2 & \frac{\partial^2 f}{\partial x \partial y} &= 0 & \frac{\partial^2 f}{\partial y^2} &= 2 & \frac{\partial^2 f}{\partial y \partial x} &= 0 \\ & \left[\begin{array}{cc} 2 & 0 \\ 0 & 2 \end{array} \right] \end{aligned}$$

Remember...



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

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

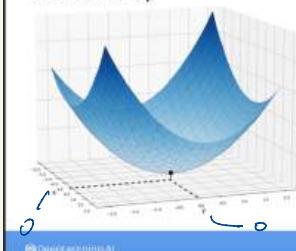
$$\det(H(0,0) - \lambda I) = \det \begin{bmatrix} 4-\lambda & -1 \\ -1 & 6-\lambda \end{bmatrix}$$

$$\rightarrow (4-\lambda)(6-\lambda) - (-1)(-1)$$

$$= \lambda^2 + 10\lambda + 23 \quad \lambda_1 = 6.41 \quad \lambda_2 = 3.59$$

$(0,0)$ is min as
these are the functions with
concave up.

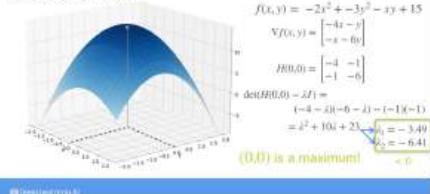
Concave Up



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

For this now do we check if the function is positive or negative we have a matrix. For that we use eigen values

Concave Down



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

$$\nabla f(x, y) = \begin{bmatrix} -4x-y \\ -x-6y \end{bmatrix}$$

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

$$\det(H(0,0) - \lambda I) =$$

$$(4-\lambda)(6-\lambda) - (-1)(-1)$$

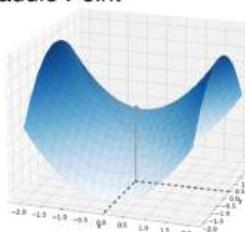
$$= \lambda^2 + 10\lambda + 23$$

$$\lambda_1 = -3.49$$

$$\lambda_2 = -6.41$$

$(0,0)$ is a maximum!

Saddle Point



$$f(x, y) = 2x^2 - 2y^2$$

$$\nabla f(x, y) = \begin{bmatrix} 4x \\ -4y \end{bmatrix}$$

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

$$\det(H(0,0) - \lambda I) =$$

$$(4-\lambda)(-4-\lambda) - 0$$

$$(0,0) \text{ is saddle point}$$

$$\lambda_1 = -4 \quad < 0$$

$$\lambda_2 = 4 \quad > 0$$

DeepLearning.AI

Summary

	1 variable $f(x)$	2 variables $f(x, y)$	More variables $f(x_1, x_2, \dots, x_n)$
(Local) minima	Happy face $f''(x) > 0$	Upper paraboloid $\lambda_1 > 0 \& \lambda_2 > 0$	All $\lambda_i > 0$
(Local) maxima	Sad face $f''(x) < 0$	Down paraboloid $\lambda_1 < 0 \& \lambda_2 < 0$	All $\lambda_i < 0$
Need more information	$f''(x) = 0$	Saddle point $\lambda_1 > 0 \& \lambda_2 < 0$ $\lambda_1 < 0 \& \lambda_2 > 0$ Or some $\lambda_i = 0$	Some $\lambda_i > 0$ and some $\lambda_j < 0$ OR At least one $\lambda_i = 0$

Eigenvectors and Eigenvalues

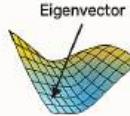
The eigenvectors of a matrix indicate directions in which a transformation acts by stretching or compressing, while the eigenvalues determine the scaling factor.

Eigenvector



Curves up in
all directions

Local minimum



Up in one, down
in another

Saddle point



Curves down in
all directions

Local maximum

DeepLearning.AI

Newton's Method

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

2 variables $\begin{bmatrix} x_{k+1} \\ y_{k+1} \end{bmatrix} = \begin{bmatrix} x_k \\ y_k \end{bmatrix} - \frac{H^{-1}(x_k, y_k)}{2 \times 2} \cdot \frac{\nabla F(x_k, y_k)}{2 \times 1}$

$$f(x, y) = x^4 + 0.8y^4 + 4x^2 + 2y^2 - xy - 0.2x^2y.$$

$$\frac{\partial F}{\partial x} = \begin{cases} 4x^3 + 8x - y - 0.4xy \\ 3.2y^3 + 4y - x - 0.2x^2 \end{cases}$$

$$\frac{\partial F}{\partial y} = \begin{cases} 0 \\ 12x^2 + 8 - 0.4xy \end{cases}$$

$$\frac{\partial^2 F}{\partial x^2} = -1 - 0.4x$$

$$\frac{\partial^2 F}{\partial y^2} = 9.6y^2 + 4$$

gradient ($\nabla F(x, y)$)

Now let's start at some point $\begin{bmatrix} x_0 \\ y_0 \end{bmatrix} = \begin{bmatrix} 4 \\ 4 \end{bmatrix}$

$$\nabla F(4, 4) = \begin{bmatrix} 277.6 \\ 213.6 \end{bmatrix} H(4, 4) = \begin{bmatrix} 198.4 & -2.6 \\ -2.6 & 157.6 \end{bmatrix}$$

$$\begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = \begin{bmatrix} 4 \\ 4 \end{bmatrix} - \begin{bmatrix} 198.4 & -2.6 \\ -2.6 & 157.6 \end{bmatrix}^{-1} \begin{bmatrix} 277.6 \\ 213.6 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} 2.58 \\ 2.62 \end{bmatrix} \rightarrow \text{closer to } 0, 0.$$

Now 2nd iteration $\begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = \begin{bmatrix} 2.58 \\ 2.62 \end{bmatrix}$

$$\nabla F(2.58, 2.62) = \begin{bmatrix} 84.25 \\ 63.4 \end{bmatrix} H(2.58, 2.62) = \begin{bmatrix} 86.83 & -2.032 \\ -2.032 & 69.39 \end{bmatrix}, \begin{bmatrix} x_2 \\ y_2 \end{bmatrix} = \begin{bmatrix} 2.58 \\ 2.62 \end{bmatrix} - \begin{bmatrix} 86.83 & -2.032 \\ -2.032 & 69.39 \end{bmatrix}^{-1} \begin{bmatrix} 84.25 \\ 63.4 \end{bmatrix}$$

$$= \begin{bmatrix} 1.59 \\ 1.67 \end{bmatrix} \rightarrow \text{repeat until you reach zero.}$$