

WEEK 2 [Week 1 we learnt about functions with 1 variable]

Now 2 variables or more.

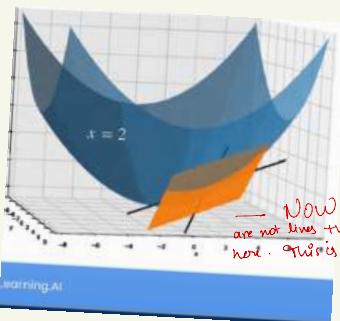
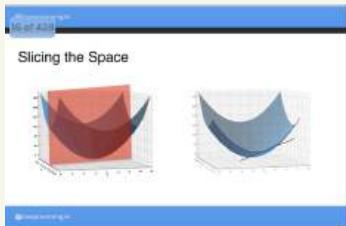
Optimizing tangent planes on 2-D is hard similarly optimizing two or more functions can get very complicated even for a computer

To speed this up

1) Gradient descent.

Functions of two Variables

$$F(x,y) = x^2 + y^2 \quad \text{2 inputs 1 output.} \rightarrow \text{So 3-D space}$$

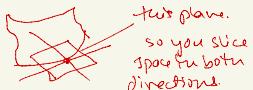


So how do we get a tangent plane. We cut the space into planes then calculate tangent on the planes

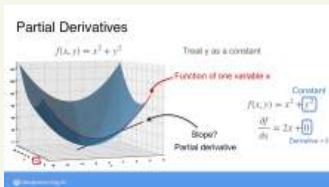
Say we cut by $y=4$
 $F(x,4) = x^2 + 16$
 $\frac{\partial F(x,y)}{\partial x} = 2x$

Fix $x=2$
 $F(2,y) = 2^2 + y^2$
 $\frac{\partial}{\partial y} (F(2,y)) = 2y$

so the lines $2x, 2y$ form the tangent plane.



Partial Derivatives



The slope of the parabola we get when fixing y as a constant.

$$\begin{aligned} & f(x,y) = x^2 + y^2 \quad \text{Function of one variable } x \\ & \text{treat } y \text{ as a constant.} \quad \text{Constant } y \\ & f(x,y) = x^2 + \boxed{y^2} \quad \frac{\partial f}{\partial x} = 2x \quad (\text{Derivative } + 0) \\ & \frac{\partial f}{\partial x} = 2x \quad \text{treat } y \text{ as a constant} \quad \text{treat } y \text{ as a constant} \\ & \frac{\partial f}{\partial y} = 0 \quad 2y \end{aligned}$$

$$\begin{aligned} & f(x,y) \\ & f_x = \frac{\partial f}{\partial x} \quad f_y = \frac{\partial f}{\partial y} \\ & \text{partial derivative of } f \text{ with respect to } x \quad \text{partial derivative of } f \text{ with respect to } y. \end{aligned}$$

$$F(x,y) = x^2 + y^2 \rightarrow \frac{\partial F}{\partial x} = 2x \quad (\text{in relation to } x \text{ } y \text{ is constant})$$

$$\frac{\partial F}{\partial y} = 2y \quad (\text{if } \quad \text{if } y \text{ is a constant})$$

$$\frac{\partial F(x)}{\partial x} = 3x^2 y^3 \quad \text{so if we take the relation to } x \text{ we treat } y \text{ as a constant.}$$

$$\frac{\partial F(x)}{\partial x} = ? \quad 3 \cdot y^3 \cdot \frac{\partial (x^2)}{\partial x} \rightarrow 3y^3 \cdot 2x = \underline{6xy^3}$$

$$\text{Now } \frac{\partial(3x^2y^3)}{\partial y}, x \text{ is now a constant so } \frac{\partial(3y^3)}{\partial y} \cdot 3x^2 = 3y^2 \cdot 3x^2 = 9x^2y^2$$

Gradient $\rightarrow F(x^2+y^2) \rightarrow$ You can slice in two ways each gives you a partial derivative
The gradient is simply the vector containing these partial derivatives.

$$\text{so for } F(x,y) \rightarrow \text{gradient} = \begin{bmatrix} 2x \\ 2y \end{bmatrix} \text{ or } \nabla F = \begin{bmatrix} \frac{\partial F}{\partial x} \\ \frac{\partial F}{\partial y} \end{bmatrix}$$

This is because $F(x,y)$ has 2 variables
if a function has 17 variables we will have 17 partial derivatives.

So these gradients $\begin{bmatrix} 2x \\ 2y \end{bmatrix}$ is a pretty good description of the tangent plane, because it gives the slopes of the two lines that form the tangent plane.

$$\text{Q1 } F(x,y) = x^2 + y^2 \quad \text{calculate the gradient } F, \nabla F \text{ at } (2,3)$$

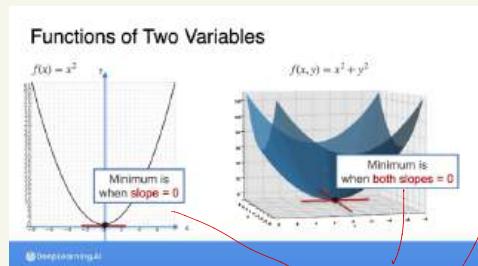
gradient = $\begin{bmatrix} \frac{\partial F}{\partial x} \\ \frac{\partial F}{\partial y} \end{bmatrix}$

$$\frac{\partial F}{\partial x} = 2x \quad \frac{\partial F}{\partial y} = 2y$$

$$\Rightarrow \begin{bmatrix} 2x \\ 2y \end{bmatrix} \text{ at } (2,3) \Rightarrow \begin{bmatrix} 4 \\ 6 \end{bmatrix}$$

The gradient is useful to minimize / maximize a function of two or more variables in the same way a derivative for 1 variable.

Functions of two variables



imp

on the left $\frac{\partial f}{\partial x} = 0$

$$2x = 0 \Rightarrow x = 0$$

on the right

$$\frac{\partial f}{\partial x} = 0 \quad \& \quad \frac{\partial f}{\partial y} = 0$$

$$2x = 0 \quad 2y = 0 \\ \text{same as } (x,y) = (0,0)$$

So basically to find the min or max you have to set the partial derivatives as zero & solve that linear equation.

Motivation for optimization in two variables

say you're in a sauna & you can move in any direction. You want to find the best place. You want to take steps to find the coldest place in the room.



HOW? → i) You look for the place where if you move in any direction from there you feel an increase in temp. Another way to look at it is if you were to take a tangent plane to a function, this tangent plane should be parallel to the floor.

Parallel to the floor.

The two partial derivatives should be zero.

Exercise: $T = F(x,y) = 85 - \frac{1}{4}x^2(x-6)y^2(y-6)$
Proactive source: given this function find $\frac{\partial T}{\partial x}$ & $\frac{\partial T}{\partial y}$

$$P(x,y) = 85 - \frac{1}{90}x^3y^3 + \frac{1}{15}x^3y^2 + \frac{1}{10}xy^3 - \frac{2}{5}x^2y^2$$

$$\rightarrow \frac{\partial f}{\partial x} \Rightarrow -\frac{1}{30}x^2y^3 + \frac{1}{5}x^2y^2 + \frac{2}{11}xy^3 - \frac{4}{5}xy^2$$

$$\Rightarrow \frac{\partial F}{\partial y} = -\frac{1}{30}x^3y^2 + \frac{2}{15}x^3y + \frac{1}{5}x^2y^2 - \frac{4}{5}x^2y$$

Then we get u to

$$\frac{\partial F}{\partial x} = 0 \quad \text{and} \quad \frac{\partial F}{\partial y} \neq 0$$

Then we get a lot of values
then we go through all pairs

$$\frac{\partial F}{\partial x} = \frac{-xy^2}{30} (x(y-6) - 4(y-6))$$

Now
 $\frac{\partial F}{\partial x} =$

$$\frac{\partial F}{\partial x} = \frac{-xy^2(x-4)(y-6)}{30} \quad (1)$$

$$\frac{\partial F_2}{\partial y} = -\frac{1}{30}x^2y(x-y-6y+24) \quad \text{Ansatz}$$

Now if
 $\frac{\partial f}{\partial x} = 0 \quad \text{and} \quad \frac{\partial f}{\partial y} = 0$
 so 3 products > 0
 that means.

$$\textcircled{1} - \frac{-xy^2}{30} (x-4)(y-6) \quad \textcircled{2} - \frac{-2xy}{30} (y-4)(x-6)$$

$$\textcircled{3} 0 = xy^2 \geq 0, x-4=0, y-6=0 \quad \textcircled{4} 0 = 2xy = 0, y-4=0, x-6=0$$

$$x=0, y=0, x=4, y=6 \quad \textcircled{5} x=0, y=0, y=4, x=6$$

From $\textcircled{3} \times \textcircled{5}$ we get the following candidates:

$$\begin{array}{ll} \textcircled{6} x=0 & \textcircled{7} x=4 \\ y=0 & y=4 \\ \textcircled{8} x=0, y=0 & \textcircled{9} x=0, y=4 \\ \textcircled{10} x=4, y=0 & \textcircled{11} x=4, y=4 \\ \textcircled{12} x=4, y=6 & \textcircled{13} x=6, y=4 \end{array}$$

red out of $\textcircled{12}$

As you can see $\textcircled{12}$ is 0 to 5 in either direction

from the candidates you see
 $x=4, y=4$ is minimal.



Let's use the power line model again (Linear regression)

But 2-D.

The goal here is to find the optimal place for a fiber connection that goes to a straight line in such a way that you reduce the total cost of connecting to the 3 power lines. This can be done by connecting wires to the power lines but one wire

must be parallel to the y-axis.
 Now if we focus on the blue power line $(1, y)$ if we make parallel to the y-axis & connect we connect at point $(1, mx+b)$
 cost of connection \rightarrow square of length of wires
 $y = mx+b$ is the fiber line we need to find m & b so that distance is minimized

goal: Minimize sum of squares cost.
 blue $\rightarrow (1, 2)$ green $\rightarrow (3, 3)$ orange $\rightarrow (2, 5)$

So you can see the distance

$$\text{For the blue line: } (m+b-2)^2$$

$$\text{for the green: } (3m+b-3)^2$$

$$\text{for the orange: } (2m+b-5)^2$$

why squared \rightarrow function is x^2

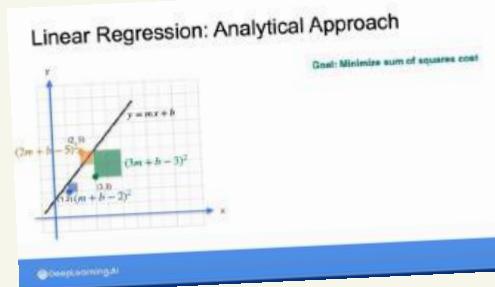
$$\text{total cost} \rightarrow (m+b-2)^2 + (3m+b-3)^2 + (2m+b-5)^2$$

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

$$\frac{\partial E}{\partial m} = 0, \quad \frac{\partial E}{\partial b} = 0 \quad \frac{\partial E}{\partial m} = 28m + 0 + 12b + 42 \quad \rightarrow 28m + 12b + 42 = 0$$

$$\frac{\partial E}{\partial b} = 6b + 12m - 20 = 0 \quad 6b + 6 - 20 = 0 \quad 6b - 14 = 0 \quad b = \frac{7}{3}$$

$$\begin{aligned} 28m + 12b - 42 &= 0 \\ -24m - 12b + 40 &= 0 \\ 4m - 2 &= 0 \\ m &= \frac{1}{2} \end{aligned}$$



$$m=1/2 \quad b=7/3 \quad E(m=1/2, b=7/3) = 4.167$$

To solve the linear equations can we do something? Yes Gradient descent

Ques. $\nabla F(x, y) = x^2y + 3x^2 \quad \frac{\partial F}{\partial x} = 2xy + 6x \quad \text{gradient } F(x, y) = ?$

gradient = $\begin{bmatrix} \frac{\partial F}{\partial x} \\ \frac{\partial F}{\partial y} \end{bmatrix} \rightarrow \begin{bmatrix} y^2 + 2 \\ 2xy + 3 \end{bmatrix}$

3) $F(x, y) = x^2 + 2y^2 + 8y \quad \text{min of } F \text{ is:}$

$\frac{\partial F}{\partial x} = 2x = 0 \quad \frac{\partial F}{\partial y} = 4y + 8$

$(x=0) \quad y = -2 \quad (y = -2)$

so if we put $F(x=0, y=-2) = 4 - 16 = -12$

Q) gradient of $x^2 + 2xy + 2z^2$

$\frac{\partial}{\partial x} = 2x + 2yz \quad \frac{\partial}{\partial y} = 2xz \quad \frac{\partial}{\partial z} = 2xy + 2z$

$2x + 2yz = 0$
 $2xz = 0$
 $2xy + 2z = 0$

Gradient Descent

When we try to solve the gradient or derivatives of functions especially ones with many variables it is very difficult

So we use Gradient Descent, - iterative way to find min/max

e.g. $f(x) = e^x - \log(x)$ minimum?

$$f'(x) = e^x - \frac{1}{x} = 0$$

$$e^x = \frac{1}{x}$$



We move in both directions, thus which is smaller we move there

- This is the learning rate.

Now is there any information, yes. The slope of any current points



If slope -ve move to right
 $x = x + \text{step}$

If new point is sl.

old point.

$$x_1 = x_0 - F'(x_0)$$

↳ slope of x_0

Big steps can be chaotic we need small steps.

$$x_1 = x_0 - 0.01 F'(x_0)$$

learning rate

You can choose any learning rate you want.

$$x_1 = x_0 - \alpha F'(x_0)$$

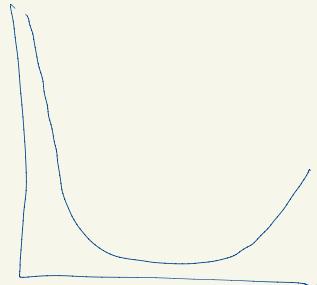
↳ learning rate

$$F(x) = e^x - \log(x) \quad F'(x) = e^x - 1/x$$

$$\text{Start } x = 0.05 \quad \text{Rate: } \alpha = 0.05$$

$$F'(0.05) = -18.9$$

$$\text{move by } -0.005 F'(0.05)$$



Now let's try gradient descent

1) start $x = 0.05, \alpha = 0.005$

$$F'(0.05) = -18.9$$

$$\text{move by } -0.005 \cdot F'(0.05) \quad x_1 = 0.1447$$

Find

$$F'(0.1447) = -5.7552$$

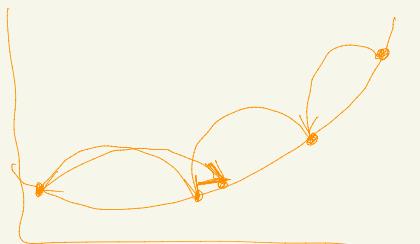
$$\text{move by } -0.005 F'(0.1447)$$

$$x \rightarrow 0.1735$$

Notice we never
needed to solve
derivative = 0
we only need to know
the derivative & apply
it to the upcoming step

What is a Good Learning Rate.

Too large



You may never reach the minimum because the steps are too big.

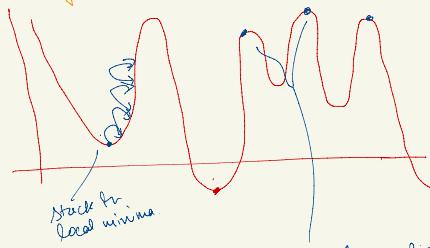
Too small



takes too long, or you might never reach it.

Very hard to get the best learning rate

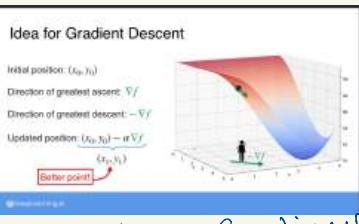
Drawbacks



Stuck for local minima.

To fix this you can take multiple initial starting points for gradient descent

Gradient descent with Heat Example. [More than 1 variable]



So in this the gradient ∇f

If you move in ∇f \rightarrow greatest ascent

so you move with ∇f direction.

In one variable we used derivative of two

we use gradient.

lets see how it works $\rightarrow T = f(x, y) = 85 - \frac{1}{90} x^2 (x-6) y^2 (y-6)$

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix} \quad \nabla f = \begin{bmatrix} -\frac{1}{90} x^2 (3x-12) y^2 (y-6) \\ -\frac{1}{90} x^2 (x-6) y (3y-12) \end{bmatrix}$$

Start $x=0.5$, $y=0.6$

$$\nabla f(0.5, 0.6) = \begin{bmatrix} -0.1134 \\ -0.0935 \end{bmatrix} \quad \text{move by } -0.05 \nabla f(0.5, 0.6)$$

Now $x = 0.5057$
 $y = 0.6047$

Now run with new x, y .

Find $\nabla f(0.5057, 0.6047) = \begin{bmatrix} -0.1162 \\ -0.0961 \end{bmatrix}$ move by $-0.05 \nabla f(0.5057, 0.6047)$

If you repeat this many times you get to the minimum $\begin{bmatrix} x \rightarrow 0.5115 \\ y \rightarrow 0.6095 \end{bmatrix}$

Functionality $\rightarrow P(x, y)$ Goal: Finding minimum of $P(x, y)$

Step 1: Define a learning rate α , choose a starting point (x_0, y_0)

Step 2: Update: $\begin{bmatrix} x_k \\ y_k \end{bmatrix} = \begin{bmatrix} x_{k-1} \\ y_{k-1} \end{bmatrix} - \alpha \nabla f(x_{k-1}, y_{k-1})$

Step 3 Repeat step 2 until you are close enough to true minimum.

You know when the true min is happening, because you move really really slowly or not at all.

Note: This still has drawbacks of local/global minimum
you can avoid this by using different places.

We solved a Power lines problem

Let's solve using gradient descent

We had the following:

(2,5), (1,2), (3,3)

The problem became

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

We solved $E(m, b)$ by finding $\frac{\partial E}{\partial m} = 0$ & $\frac{\partial E}{\partial b} = 0$

Now let's try to find gradient descent.

The gradient was

$$\begin{bmatrix} 28m + 12b - 42 \\ 6b + 12m - 20 \end{bmatrix}$$

Linear Regression: Gradient Descent

Goal: Minimize sum of squares cost

$$\nabla E = [28m + 12b - 42, 6b + 12m - 20]$$

$$m = ?$$

The points m, b such that the cost is minimum

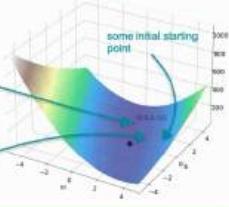
Steps:

Start with (m_0, b_0)

descend until you find the minimum

iterate

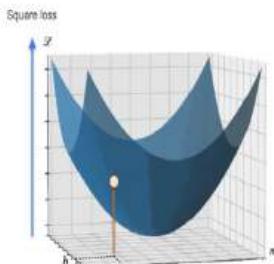
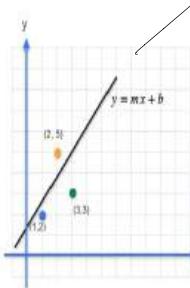
$$(m_{k+1}, b_{k+1}) = (m_k, b_k) - \alpha \nabla E(m_k, b_k)$$



@DeepLearningAI

for a good line

Gradient Descent

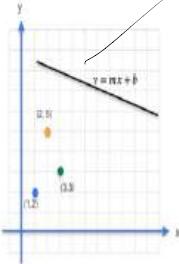


@DeepLearningAI

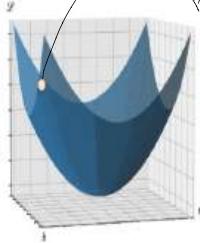
We map m, b in the m, b plane

We are stacking up the cost function on y -axis or vertical axis.

Gradient Descent

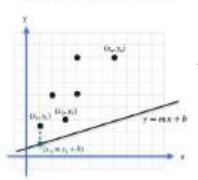


for a terrible observation
gradient descent plots at wrong point



@DeepLearningAI

Gradient Descent



Now we add all losses of all points & divide by n to get average

$$L(m, b) = \frac{1}{n} \sum_{i=1}^n (mx_i + b - y_i)^2$$

This extra 2 is that when we take derivative exponent multiplies by 2
It doesn't really matter but still

Now we create a descent.

We start at a random $\begin{bmatrix} m_0 \\ b_0 \end{bmatrix} \rightarrow \begin{bmatrix} m_1 \\ b_1 \end{bmatrix}$

$$\begin{bmatrix} m_1 \\ b_1 \end{bmatrix} = \begin{bmatrix} m_0 \\ b_0 \end{bmatrix} - \alpha \nabla L_1(m_0, b_0)$$

$$\begin{bmatrix} m_n \\ b_n \end{bmatrix} = \begin{bmatrix} m_{n-1} \\ b_{n-1} \end{bmatrix} - \alpha \nabla L_2(m_{n-1}, b_{n-1})$$

formal version of gradient descent