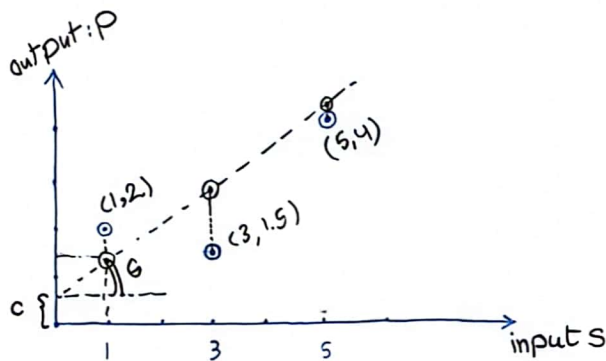


Linear Regression

using Gradient descent

starting with a simple problem.



Q: what is the best fitting line?

let's assume slope is already known and we want to get the intercept.

$$\text{so } y = w_0 + x \cdot w_1, \quad w_1 = 1$$

so to get the best value for the intercept we will construct a loss function (least square) and then use Gradient-descent to optimize it

$$\begin{aligned} \text{Loss Function} &= \sum (P - \underbrace{P(s)}_{\text{predicted}})^2 \\ \text{sol. f} &= \sum_{i=1}^n \left[P - \left(\underbrace{w_0}_{\text{intercept}} + \underbrace{w_1}_{\text{slope}=1} \cdot \underbrace{s_i}_{\text{real}} \right) \right]^2 \\ &= [2 - (w_0 + 1)]^2 + [3 - (w_0 + 1.5)]^2 + [4 - (w_0 + 5)]^2 \end{aligned}$$

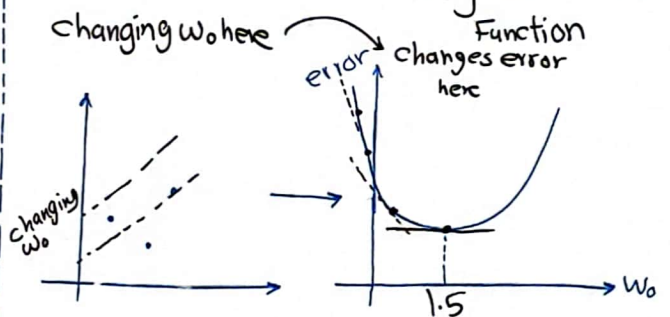
so now we take the $\frac{dL.F}{dw_0}$ to get the ^{gradient} slope of the loss function with respect to w_0

$$\begin{aligned} \text{so } \frac{dL.F}{dw_0} &= -2 [2 - (w_0 + 1)] \\ &\quad -2 [3 - (w_0 + 1.5)] \\ &\quad -2 [4 - (w_0 + 5)] \end{aligned}$$

now using this equation

$$w_{0_New} = w_{0_old} - \text{slope} * (\text{learning rate})$$

using code or calculator taking look at l.s



so some iteration the answer the answer for $w_0 = 1.5$

this was for a single parameter w_0

let's take a look at w_0 & w_1

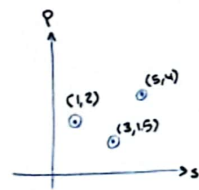
Linear Regression

using Gradient descent

so now working with Two Parameters

$$y = w_0 + w_1 x$$

Notice that
Number of parameters
equals $N.O \text{ Features} + 1$



using the same input points

$n = N.O \text{ examples}$
"points"

$$\begin{aligned} \text{Loss Function} &= \sum_1^n (P - P(s_i))^2 \\ &= (2 - (w_0 + w_1 \cdot 1))^2 + (1.5 - (w_0 + w_1 \cdot 3))^2 \\ &\quad + (4 - (w_0 + w_1 \cdot 5))^2 \end{aligned}$$

Now we have a Function of 2 Variables
so we should take the derivative with respect
to each one of them

$$\frac{\partial L.F}{\partial w_0} = -2[2 - (w_0 + w_1)] - 2[1.5 - (w_0 + 3w_1)] - 2[4 - (w_0 + 5w_1)]$$

$$\frac{\partial L.F}{\partial w_1} = -2[2 - (w_0 + w_1)] - 2 \times 3[1.5 - (w_0 + 3w_1)] - 2 \times 5[4 - (w_0 + 5w_1)]$$

so now we need to Calculate
 w_{0_new} & w_{1_new}

$$w_{0_new} = w_{0_old} - \frac{\partial L.F}{\partial w_0} \times \text{learningRate}$$

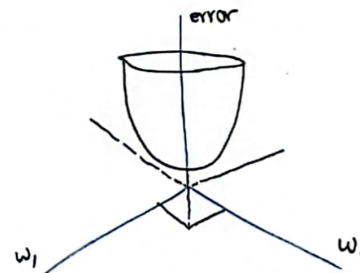
$$w_{1_new} = w_{1_old} - \frac{\partial L.F}{\partial w_1} \times \text{learningRate}$$

we can Turn it into Vector
operations

$$\begin{bmatrix} w_0 \\ w_1 \end{bmatrix} = \begin{bmatrix} w_{0_old} \\ w_{1_old} \end{bmatrix} - \begin{bmatrix} \frac{\partial L.F}{\partial w_0} \\ \frac{\partial L.F}{\partial w_1} \end{bmatrix} \times \text{Learning Rate}$$

using code for iterating & updating
we will reach the minimum error

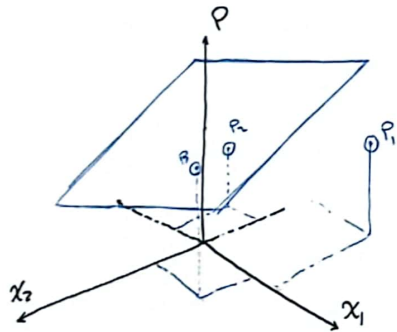
Looking at the loss Function in 3d
and BTW this is the limit of Visualizing



Linear Regression

using Gradient descent

let's take a look if we increased the Number of the inputs to 2:
we will have something like this



now the problem will be Fitting this plane to these points

$$y = w_0 x_0^n + w_1 x_1^n + w_2 x_2^n$$

the equation of the plane

we will do the same as previous

but now we will try to Vectorize these operations

where $x_0^n = 1$
 n is number of point where each point t target (x_1, x_2, P)

working with multi variate problems:

$$y = \sum_{i=0}^D (w_i x_i^{(n)})$$

D is the Number of Features
so $D+1$ param
 $0 \rightarrow D$ 0 indexed

and as known

$$\begin{bmatrix} w_0 & w_1 & w_2 & w_3 \end{bmatrix}_{1 \times D+1} \cdot \begin{bmatrix} x_0^n & x_1^n & x_2^n \\ x_3^n & x_4^n & x_5^n \end{bmatrix}_{D+1, n} = y(w_i, x_i^n)$$

$w^T \cdot x^n$

$$y = w^T \cdot \underbrace{x^n}_{1 \times n}$$

So this is the hyper-plane Function after vectorizing it

so the new loss Function

$$\text{cost}(w_i) = \frac{1}{2N} \sum_{n=1}^N (y(x_i^n, w_i) - t^n)^2$$

$$\frac{\partial \text{cost}}{\partial w_i} = \frac{1}{N} \sum_{n=1}^N (y(x_i^n, w_i) - t^n) \cdot x_i^n$$

N = number of points

Vectorizing the $\frac{\partial \text{cost}}{\partial w_i}$

$$\frac{\partial \text{cost}}{\partial w_i} = \sum_{n=1}^N (y(X_i^n, w_i) - t^n) \cdot X_i^n$$

$$y(X_i^n, w_i) - t^n$$

$$\begin{bmatrix} \text{sum of errors}_1 & \text{sum of errors}_2 & \dots & N \end{bmatrix} \begin{bmatrix} 1 & 1 \\ x_1^0 & x_1^1 \\ x_2^0 & x_2^1 & \dots & N \\ x_3^0 & x_3^1 \end{bmatrix}$$

$$1 \times N$$

$$D+1 \times N$$

$$X \odot [y(X_i^n, w_i^n) - t^n]^T$$

$$D+1 \times N \quad N \times 1$$

let's call it

$$E$$

so now

$$\frac{\partial \text{cost}}{\partial w_i} = \begin{bmatrix} \frac{\partial \text{cost}}{\partial w_0} \\ \frac{\partial \text{cost}}{\partial w_1} \\ \vdots \\ D+1 \end{bmatrix} = X \odot E^T$$

Transpose

تم