

GRADIENT DESCENT RULE FOR MEAN SQUARED ERROR

Let's take a model as hypothesis initially

$$\bar{y} = mx + c$$

where $m \rightarrow$ weight \leftarrow $c \rightarrow$ bias \leftarrow $\bar{y} \rightarrow$ output

Let y be the actual value

for a training model (x^0, y^0) calculate \bar{y}^0

$$\bar{y}^0 = mx^0 + c \quad \text{if } m = (w_0, w_1) \text{ then } \sum_{i=1}^m \frac{w_i}{m} = w_1$$

$y^0 - \bar{y}^0$ gives us loss

Squaring it to amplify the loss & make it positive

$$(y^0 - \bar{y}^0)^2 = L^0 \quad \text{if } L^0 = (w_0, w_1) \geq 0$$

Doing it for m training models and adding them up

$$(y^0 - \bar{y}^0)^2 + (y^1 - \bar{y}^1)^2 + \dots + (y^m - \bar{y}^m)^2 = L^0 + L^1 + \dots + L^m$$

$$\sum_{i=0}^m (y^{(i)} - \bar{y}^{(i)})^2 \rightarrow \text{find minimum of } L \text{ for } i \text{ from } 0 \text{ to } m$$

Dividing both sides with m to find mean of loss L

$$\frac{1}{m} \sum_{i=0}^m (y^{(i)} - \bar{y}^{(i)})^2 = L(x) \text{ (mean of loss)}$$

$L(x)$ is to find with regards to weight and bias

This gives us the loss as a function of x

where lesser the loss better the fit is

Let's do the initial case of finding minimum

the less the loss the better the fit is during fitting

We then like to minimize this loss function to minimize loss after

by tweaking the weight and bias

We are finding the minimum of a bowl shaped function where w_0 and w_1 are w and b respectively

Partial differentiation w.r.t. ω in $J(\omega)$ about ω^*

$m \rightarrow \omega$, $c \rightarrow b$

$$0 + x_m = \bar{b}$$

$$1 \sum_{i=0}^m (y^{(i)} - \omega_i x^{(i)} - b)^2 = \text{Lagrangian cost function}$$

subject to certain constraint

Differentiating w.r.t. ω we have gradient of J

$$-2 \sum_{i=0}^m x^{(i)} (y^{(i)} - \bar{y}^{(i)}) = \frac{\partial J}{\partial \omega} = \frac{\partial \mathcal{L}}{\partial \omega} + \frac{\partial \mathcal{L}}{\partial \omega} = \frac{\partial \mathcal{L}}{\partial \omega} - \frac{\partial \mathcal{L}}{\partial \omega}$$

Differentiating w.r.t. b we have gradient of J

$$-2 \sum_{i=0}^m (y^{(i)} - \bar{y}^{(i)}) = \frac{\partial J}{\partial b} = \frac{\partial \mathcal{L}}{\partial b} - \frac{\partial \mathcal{L}}{\partial b}$$

at m th iteration we have gradient in each of the points.

The gradient descent update rule

The goal is to minimize cost function

In gradient descent,

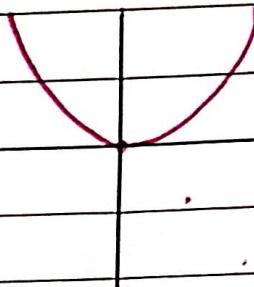
1. We first start with some random values of ω and b [Typically: $0] = \mathbb{E}(y^{(0)}, b^{(0)})$]

2. We keep on changing the values of ω and b to reduce $J(\omega)$ and $J(b)$ until we get minimum

This is basically taking small steps to reach minimum possible elevation. The one point to note with gradient descent is depending on which point you are starting from you can end up at different local minima. Since we are dealing with linear regression and the function is in $J(\omega, b)$ which represents a bowl like graph which only has one minimum which is

In this case, the function is in $J(\omega, b)$ which represents a bowl like graph which only has one minimum which is

global minima.



We subtract the the gradient descent value because if we test it and check if we add we end up moving away from global minima.

We also need to choose ' α ' a learning rate to determine the size of the step taken towards global minimum. We need to make sure ' α ' is not too large to make sure it doesn't overshoot past the global minima & Not too small to make sure it doesn't take too many iterations. In practice we take $\alpha = 0.01$ But in Theory, we test α with multiple value & see which brings down the fastest. That's where we get the formula

$$\omega = \omega - \alpha \nabla J(\omega) \quad b = b - \alpha \nabla J(b)$$

$\nabla J(\omega) \rightarrow$ represents gradient descent

$-\nabla J(\omega) \rightarrow$ represents gradient descent

$$\omega = \omega - \alpha \nabla J(\omega) \quad b = b - \alpha \nabla J(b)$$

are the formulae used in my Python program