

Gradient Descent

Prerequisite - OLS, differentiation

OLS - Gives value of m & c using formula. No need for iteration!

$$c = \bar{y} - m\bar{x}$$

$$m = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

where \bar{y} = mean of y ; \bar{x} = mean of x .

derivation.

We want $\frac{dJ}{dc} = 0$ & $\frac{dJ}{dm} = 0$. where $J = \frac{1}{2n} \sum_{i=1}^n (y_i - mx_i - c)^2$

So, $\frac{dJ}{dc} \Rightarrow \frac{1}{n} \sum_{i=1}^n (y_i - mx_i - c) = 0$

$= \bar{y} - m\bar{x} - \frac{1}{n} \sum_{i=1}^n c = 0 \Rightarrow c = \bar{y} - m\bar{x}$

Putting this in $\frac{dJ}{dm} = 0$.

Disadv. of OLS:-
It is computationally expensive
with high no. of features. So,
grad. descent is introduced.

as big matrix multiplication \rightarrow more would have to be done (see OLS for multiple reg)

$\frac{d}{dm} \left(\frac{1}{2n} \sum_{i=1}^n (y_i - mx_i - \bar{y} + m\bar{x})^2 \right) = 0$

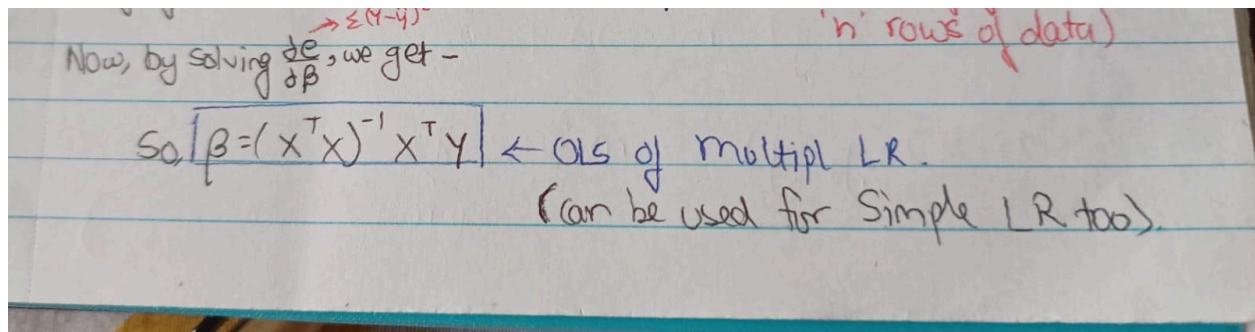
$\Rightarrow \sum_{i=1}^n 2 (y_i - mx_i - \bar{y} + m\bar{x}) (-x_i + \bar{x}) = 0$

$\Rightarrow \sum_{i=1}^n - (x_i - \bar{x}) (y_i - \bar{y} - m(x_i - \bar{x})) = 0$

$\Rightarrow \sum_{i=1}^n [(x_i - \bar{x}) (y_i - \bar{y}) - m(x_i - \bar{x})^2] = 0$

$\Rightarrow \sum_{i=1}^n (x_i - \bar{x}) (y_i - \bar{y}) = m \sum_{i=1}^n (x_i - \bar{x})^2$

$\Rightarrow m = \frac{\sum_{i=1}^n (x_i - \bar{x}) (y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$



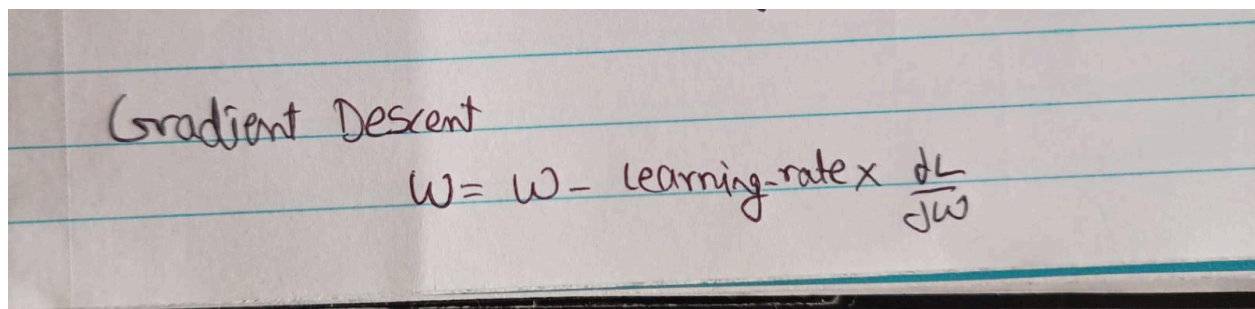
Problem with OLS (Ordinary least squares) - OLS is computationally expensive when high number of features due to more matrix inverse calculation and matrix multiplication involved.

So, Gradient Descent introduced.

Formula -

(how it works ? +ve an -ve slope)

Learning rate - why ? what ?

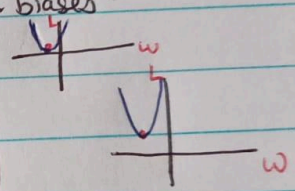


So, gradient descent

For, mse & logloss L is a function of weights & biases

mse: $L_y = (3 - (3w + 4))^2$

logloss: $L_y = -\left[\ln\left(\frac{1}{1+e^{-(4w+9)}}\right) + \ln\left(1 - \frac{1}{1+e^{-(4w+9)}}\right)\right]$



It form a parabola & we want to find such value of w where L is minimum (red dot on graphs).

Cost function (MSE)

Solving this

$$\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_j} = \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Differentiation -

when $j=0$,

$$\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0} = \frac{\partial}{\partial \theta_0} \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2$$

$$\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0} = \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})$$

when $j=1$,

$$\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1} = \frac{\partial}{\partial \theta_1} \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2$$

$$\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1} = \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)}) x^{(i)}$$

So, Repeat until convergence

look out for when cost becomes more or less stable/constant (or) do early stopping

$$\theta_0 = \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$$

$$\theta_1 = \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)}$$

When to stop?

- When the weights becomes more or less constant (doesn't change with more iterations)
- When total number of iterations decided is reached

Visualizations - [link](#) @58:00 - gradient descent

[Link](#) @1:37:00 - visualization (look at the code then watch this)

[Link](#) @1:44:00 - effect of learning rate

Types of gradient descents -

Case of Local minima - IN case of non convex functions, Batch gradient descent gets stuck in local minima. So stochastic and mini batch gradient descent used. These are used when the data is 'big' (as it converges faster with less no. of epochs)

Types of Grad. descent.

Grad. descent \rightarrow way to minimize loss function.

① Gradient descent (Batch) \rightarrow Entire data for updation ^{of wt} in one epoch.
1 update in 1 epoch.

② Stochastic GD - wt. update for each data row in each epoch.
 $\xrightarrow{\text{row / no. of eg.}} n' \text{ update in 1 epoch.}$

③ miniBatch GD - Make small batches, ~~2th~~ say 4 batches, so for each epoch, 4 updation would take place.

Batch.

SGD.

① Time - Faster

① Time - Slower.

② For equal epoch, this will converge slower.

② For equal epoch, this will converge to proper solution faster - more updates.

③ $\text{loss} \downarrow$ smooth loss decrease

③ Bounce behavior $\text{loss} \downarrow$

④ Can move in local minime.

④ Helps in moving out of local minime.

⑤ Exact sol.

⑤ Exact solution not found, due to bounce behaviour.