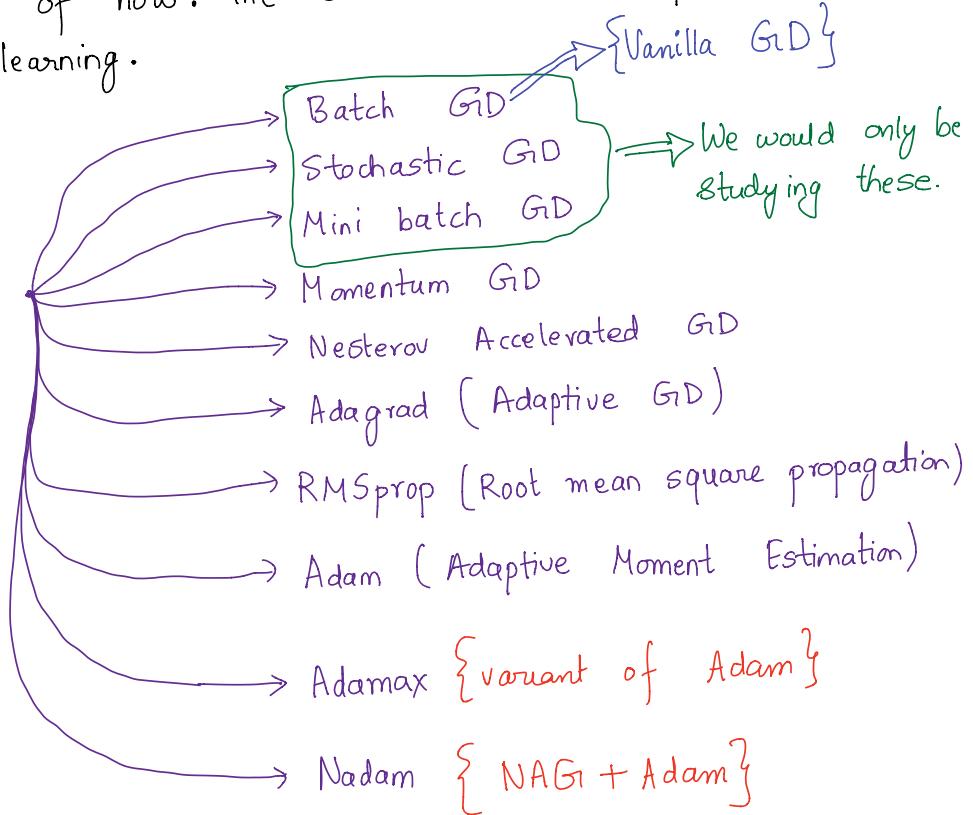


## ① {Types of gradient descent}

We basically have 10 types of GD. However, we would only be studying 3 as of now. The other variations of GD are mostly used in deep learning.



The only difference between these 3 types of GD is the # of data points required for single updatation of the coefficients.

The Batch Gradient descent is just the normal GD in which the coefficients are updated once after every epoch. So if we 1000 epochs, we would be updating the coefficients only 1000 times.

## ② {Understanding Batch GD with simple example}

Consider the dataset having 2 input features namely, CGPA and

IQ, and one output feature, namely package (lpa). For the sake of simplicity, we will only consider 2 data points i.e., only 2 rows.

$x_1$	$x_2$	$y$
cgpa	iq	lpa
8.5 ( $x_{11}$ )	93 ( $x_{12}$ )	3.7 $\Rightarrow y_1$
8.7 ( $x_{21}$ )	98 ( $x_{22}$ )	4.2 $\Rightarrow y_2$

We need to find 3 coefficients  $\beta_0, \beta_1$  and  $\beta_2$  in order to train our linear regression model.

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

Thus,

$$\hat{y}_1 = \beta_0 + \beta_1 x_{11} + \beta_2 x_{12}$$

$$\hat{y}_2 = \beta_0 + \beta_1 x_{21} + \beta_2 x_{22}$$

Suppose we have  $m$  columns (predictors), then the above equation could be written as :-

$$\begin{cases} \hat{y}_1 = \beta_0 + \beta_1 x_{11} + \beta_2 x_{12} + \dots + \beta_m x_{1m} \\ \hat{y}_2 = \beta_0 + \beta_1 x_{21} + \beta_2 x_{22} + \dots + \beta_m x_{2m} \end{cases}$$

Here, we have  $(m+1)$  coefficients to be found namely  $\langle \beta_0, \beta_1, \dots, \beta_m \rangle$  where  $m \triangleq \# \text{ of predictors}$

Now further suppose that we have  $n$  data points in our dataset. Then we would have  $n$   $\hat{y}'s$ .

$$\begin{aligned}\hat{y}_1 &= \beta_0 + \beta_1 x_{11} + \beta_2 x_{12} + \dots + \beta_m x_{1m} \\ &\vdots \\ \hat{y}_n &= \beta_0 + \beta_1 x_{n1} + \beta_2 x_{n2} + \dots + \beta_m x_{nm}\end{aligned}$$

This could also be written as -

$$\hat{y}_i = \beta_0 + \sum_{j=1}^m \beta_j x_{ij}$$

The error that our model makes on each input is

$$E = \hat{y}_i - y_i$$

So, the Loss function (cost function) can be written as

$$L = \frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{n}$$

This is a function of  $(m+1)$  variables. Thus  $L: \mathbb{R}^{m+1} \rightarrow \mathbb{R}$  can be viewed as a function of  $\beta_0, \beta_1, \dots, \beta_m$ .

$$L(\beta_0, \dots, \beta_m) = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

$$L(\beta_0, \dots, \beta_m) = \frac{1}{n} \sum_{i=1}^n \left( \beta_0 + \sum_{j=1}^m \beta_j x_{ij} - y_i \right)^2$$

In order to find the optimal value of  $\beta_0, \dots, \beta_m$  for which this function  $L$  takes minimum value, we need to use the method of GD which is a first order linear approximation technique. The new values of  $\beta_0, \dots, \beta_m$  are calculated by using the learning rate  $\eta$  and the value of gradient of  $L$  at the old values of  $\beta_0, \dots, \beta_m$ .

Thus,

$$(\beta_0, \dots, \beta_m)_{\text{new}} = (\beta_0, \dots, \beta_m)_{\text{old}} - \eta (\nabla L)_{\substack{\text{evaluated at} \\ (\beta_0, \dots, \beta_m)_{\text{old}}}}$$

For calculating  $\nabla L$  at  $(\beta_0, \dots, \beta_m)_{\text{old}}$ , we need to calculate

$$\frac{\partial L}{\partial \beta_0}, \dots, \frac{\partial L}{\partial \beta_m}$$

③ { Calculating  $\frac{\partial L}{\partial \beta_0}, \frac{\partial L}{\partial \beta_1}, \dots, \frac{\partial L}{\partial \beta_m}$  }

we have

$$L = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

Now  $\frac{\partial L}{\partial \beta_0} = \frac{1}{n} \sum_{i=1}^n 2(\hat{y}_i - y_i) \left( \frac{\partial \hat{y}_i}{\partial \beta_0} \right)$  { as  $\hat{y}_i$  is a function of  $\beta_0, \dots, \beta_m$  }

Similarly,  $\frac{\partial L}{\partial \beta_m} = \frac{1}{n} \sum_{i=1}^n 2(\hat{y}_i - y_i) \left( \frac{\partial \hat{y}_i}{\partial \beta_m} \right)$

These derivatives could be enclosed in a vector

$$\left( \frac{\partial L}{\partial \beta_0}, \dots, \frac{\partial L}{\partial \beta_m} \right) = \left( \frac{1}{n} \sum_{i=1}^n 2(\hat{y}_i - y_i) \frac{\partial \hat{y}_i}{\partial \beta_j} \right)_{j=0}^m$$

Thus,  $\frac{\partial L}{\partial \beta_j} = \frac{1}{n} \sum_{i=1}^n 2(\hat{y}_i - y_i) \frac{\partial \hat{y}_i}{\partial \beta_j}$ ,  $\forall j \in \{0, \dots, m\}$

Now since  $\hat{y}_i = \beta_0 + \sum_{j=1}^m \beta_j x_{ij}$

$$\Rightarrow \frac{\partial \hat{y}_i}{\partial \beta_j} = \begin{cases} 1 & \text{if } j=0 \\ x_{ij} & \text{if } j \neq 0 \end{cases}$$

Thus,  $\frac{\partial L}{\partial \beta_0} = \frac{1}{n} \sum_{i=1}^n 2(\hat{y}_i - y_i) \frac{\partial \hat{y}_i}{\partial \beta_0} \Rightarrow 1$

$$\boxed{\frac{\partial L}{\partial \beta_0} = \frac{2}{n} \sum_{i=1}^n (\hat{y}_i - y_i)} - ①$$

and  $\frac{\partial L}{\partial \beta_j} = \frac{1}{n} \sum_{i=1}^n 2(\hat{y}_i - y_i) \frac{\partial \hat{y}_i}{\partial \beta_j} \Rightarrow x_{ij} \text{ where } j \neq 0$

$$\Rightarrow \frac{\partial L}{\partial \beta_j} = \frac{2}{n} \sum_{i=1}^n (\hat{y}_i - y_i) x_{ij}, \boxed{j \neq 0}$$

Summarizing the above equations, we get

$$\frac{\partial L}{\partial \beta_j} = \begin{cases} \frac{2}{n} \sum_{i=1}^n (\hat{y}_i - y_i) , & j = 0 \\ \frac{2}{n} \sum_{i=1}^n (\hat{y}_i - y_i) x_{ij} , & j \neq 0 \end{cases}$$

and thus now the coefficients  $\beta_0, \dots, \beta_m$  can be updated as

$$(\beta_0, \dots, \beta_m)_{\text{new}} = (\beta_0, \dots, \beta_m)_{\text{old}} - \eta \left( \frac{\partial L}{\partial \beta_0}, \dots, \frac{\partial L}{\partial \beta_m} \right)_{\substack{\text{evaluated at} \\ (\beta_0, \dots, \beta_m)_{\text{old}}}}$$

It is very clear from the above equation, that for getting the new values of  $\beta$ 's we need to evaluate  $\frac{\partial L}{\partial \beta}$ 's at the old value of  $\beta$ 's. Each  $\frac{\partial L}{\partial \beta}$  has a summation which is running from 1 to n. Thus, while calculating each derivative in each epoch, we need to make use of all the data points of our dataset which makes the overall process of updation slower.

#### ④ { code for batch gradient descent }

we will make our own class of batch GD. In this we shall be making 2 methods. The `__init__()` constructor will take two parameters namely, `learning_rate` and `epochs`

```

def __init__(self, learning_rate, epochs):
    self.coef_ = None
    self.intercept_ = None
    self.lr = learning_rate
    self.epochs = epochs

```

The `predict()` method is fairly easy and takes in only 1 parameter namely  $X_{\text{test}}$ . Further, let us assume that we have  $n$  data points and  $m$  predictors. Thus,

$$X \in \mathbb{R}^{n \times m} \quad X_{\text{test}} \in \mathbb{R}^{n_1 \times m}$$

$y \in \mathbb{R}^{1 \times n}$

$\downarrow$   
 # of rows  
 in  $X_{\text{test}}$

$$\text{self.coef}_- \in \mathbb{R}^{1 \times m} \triangleq (\beta_1, \dots, \beta_m)$$

```

def predict(self, X_test):
    return X_test @ self.coef_.T + self.intercept_

```

$$\left\{
 \begin{array}{l}
 X_{\text{test}} \in \mathbb{R}^{n_1 \times m} \text{ and } \text{self.coef}_- \in \mathbb{R}^{1 \times m}. \text{ So, } \text{self.coef}_-.T \\
 \in \mathbb{R}^{m \times 1}. \text{ Thus, } X_{\text{test}} @ \text{self.coef}_-.T \in \mathbb{R}^{n_1 \times 1}.
 \end{array}
 \right\}$$

The final and the most important method is `fit()`. This is the most important method and takes 2 parameters,  $X_{\text{train}}$

and  $y_{\text{train}}$ :

$$X_{\text{train}} \in \mathbb{R}^{n_0 \times m}$$

$$y_{\text{train}} \in \mathbb{R}^{1 \times n}$$

```
def fit(self, X_train, y_train):
```

$$(\beta_1, \dots, \beta_m) \{ \text{self.coef\_} = \text{np.ones}(X_{\text{train}}.shape[1]) \# \begin{array}{l} \text{initialize} \\ \beta_1, \dots, \beta_m \\ = 1 \end{array}$$
$$\underbrace{\text{self.intercept\_}}_{\beta_0} = 0 \# \text{initialize } \beta_0 = 0$$

```
for _ in range(epochs):
```

$$y_{\text{pred}} = \text{self.intercept\_} + X_{\text{train}} @ \text{self.coef\_}$$

$$\text{del_beta_not} = [2 * \underbrace{\text{np.sum}(y_{\text{pred}} - y)}_{\cdot \text{shape}[0]}] / X_{\text{train}}$$
$$\# \frac{\partial L}{\partial \beta_0} = \frac{2}{n_0} \sum_{i=1}^{n_0} (\hat{y}_i - y_i)$$

$$\text{others} = 2 * \left( \underbrace{(y_{\text{pred}} - y)}_{n_0 \times 1} @ X_{\text{train}} \right) / \underbrace{X_{\text{train}}.shape[0]}_{n_0}$$

$$\left( \frac{\partial L}{\partial \beta_1}, \dots, \frac{\partial L}{\partial \beta_m} \right) = \frac{2}{n_0} (\hat{y} - y) X_{\text{train}}$$

$$\text{self.intercept\_} = \text{self.intercept\_} - \text{self.lr} * \text{del_beta\_not}$$

$$\text{self.coef\_} = \text{self.coef\_} - \text{self.lr} * \text{others}$$

## ⑤ {Problems with Batch gradient descent}

### a) {Computational Expanse}

For large datasets, calculating the gradient of the cost function over the entire dataset can be computationally expansive and slow.

### b) {Memory Usage}

It requires loading the entire dataset into memory, which can be impractical for very large datasets.

### c) {Convergence Issues}

It may converge slowly if the cost function has many local minima or saddle points

### d) {Stuck in local minima}

In case of convex functions, the batch GD works fine. However, when we study deep learning, the loss function is often NOT a simple convex function. It consists of multiple local minima's and global minima. The batch GD algorithm is very likely to be trapped inside a local minimum compared to other algorithms like Stochastic GD or mini batch gradient descent.

