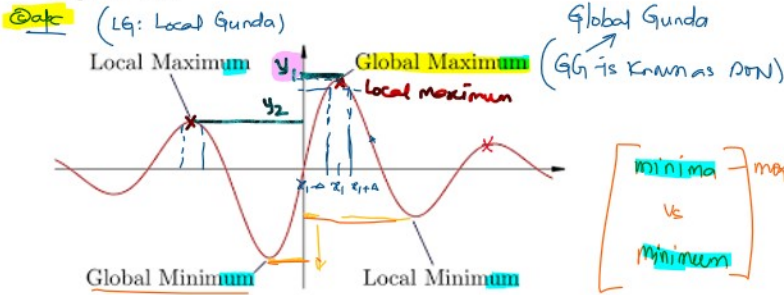


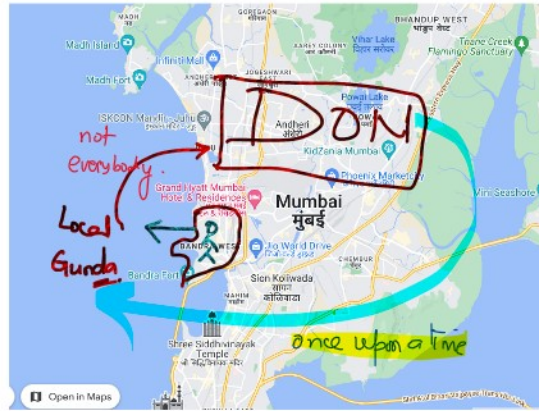
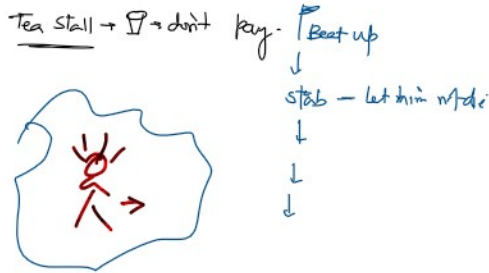
Types of GDA

03 August 2024 20:13



more than one point of minimum values
vs
Minimum

Story career roadmap for a Don



Stochastic Gradient Descent (SGD) → an optimization algorithm

refers to systems or processes that are random or probabilistic

$$LR \quad y = \beta_0 + \beta_1 x$$

NN: weights & biases.

In the context of machine learning/AI, stochastic describes algorithms or models that incorporate randomness in their operation

SGD: SGD is an optimization algorithm used to minimize the loss function. To do so, SGD computes the gradient for a single randomly chosen training data at each iteration



Leveraging SGD, gradient is computed for each and every sample in the dataset → hence it updates every sample in the dataset

Working of SGD:

#epochs = 100

$$100 \times 100 = 10,000$$

a) Initialization

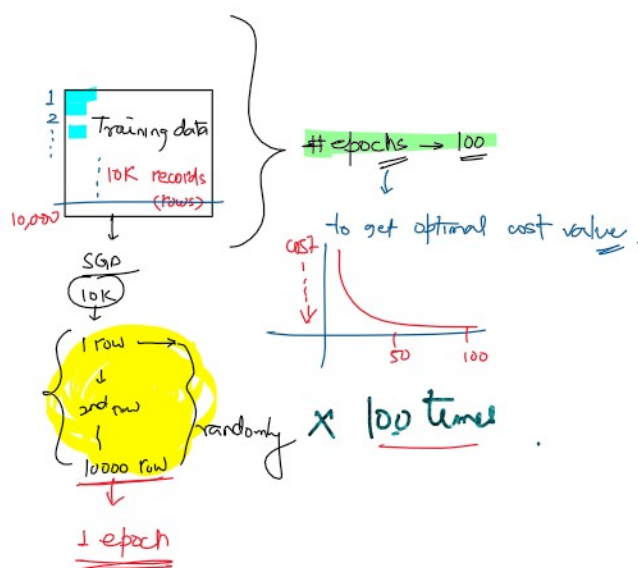
- Initialize the model parameters (weights and biases) with random values or zeroes.
- Set the learning rate (α or η)
 - ↓
 - set by modelers
 - $\alpha: (0.01 - 0.2)$ $\alpha = [0.01, 0.1, 0.15, 0.2]$
- Define the no. of epochs (iteration over the entire dataset)

In SGD — Training data : 10,000 records

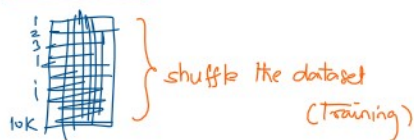
epochs # 100.

How many times \rightarrow model will update the weights & biases

Ans: For SGD: $100 \times 10,000 = 1 \times 10^6 \Rightarrow 1 \text{ million times}$.



b) Iteration



- shuffle the training dataset to ensure the random sampling
- For randomly selected single training data point (row):
 - * compute the gradient of the loss function
 - * Update the model parameters (weight & bias) using the computed gradient and learning rate

c) Update rule:

* For $j=0$

$$\theta_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} [J(\theta_0, \theta_1)]$$

Input \rightarrow 'by modeler'

\vdots $\times 100$ } time taking



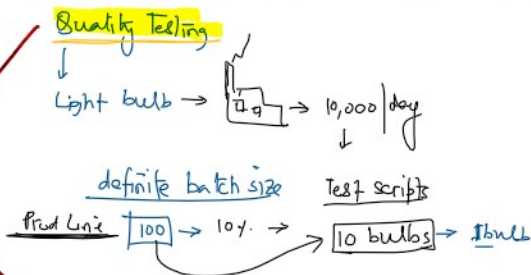
Vanilla \rightarrow No flavor
'Plain'

[datapoint / row / record / training example]

Batch Gradient Descent (BGD)

Batch Gradient Descent (BGD) computes the gradient of the loss function w.r.t. the parameters using the entire training dataset

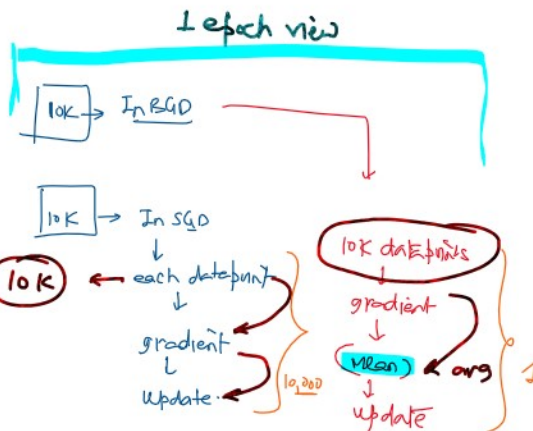
batch



In BGD, all the training data is taken into a consideration to take a single step

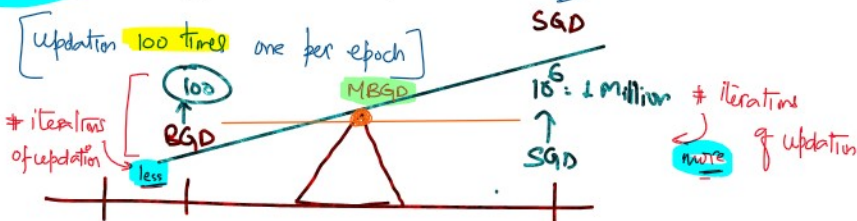
we would take the average of the gradients of all the training datapoints and use the mean gradient to update the model parameters (one number) / epoch

Just one step of gradient descent in one epoch.



epochs $\rightarrow 100$

to get optimal cost value.





more of update

Vanilla \rightarrow BGD : Entire dataset

learning rate \rightarrow Input for Update Formula

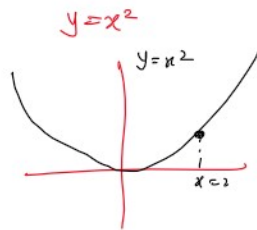
* For $j=0$

$$\theta_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} [J(\theta_0, \theta_1)]$$

$$2 - 0.01 \times 0.2$$

$$2 - 0.002$$

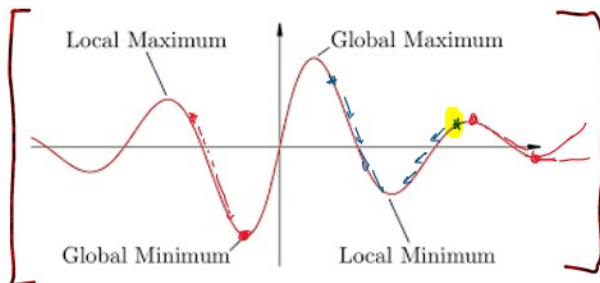
$$2 - 0.002 = 1.998$$



$$\frac{dy}{dx} = 2x$$

$$\left. \frac{dy}{dx} \right|_{x=2} = 2 \times 2 = 4$$

Newton - Raphson's Method



Task: Understand the importance of initializing GDA.

Mini-Batch Gradient Descent (MBGD)

MBGD is a trade-off between Batch Gradient and stochastic Gradient descents.

Training data 10,000 data records

batch size = 32 data records.

No. of batches : $10000/32 = 312.5 \approx 313$ ✓ 313 batches \rightarrow

Given # epochs = 100



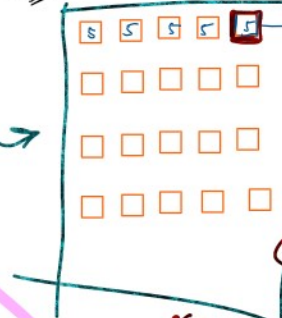


Q: For a dataset of 100 samples, if the batch size is 5, how many times updates will occur for 100 epochs.

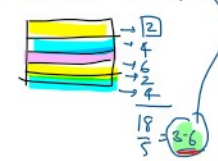
Ans - $\frac{\text{# batches} \cdot 100}{5} = 20 \text{ batches} \text{ (epoch)}$

$20 \times 100 = 2000 \text{ updates}$

20 batches MBGD.



Mean gradient

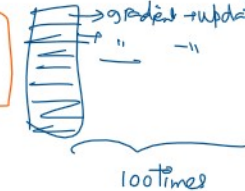
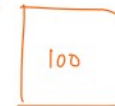


20 times

BGD

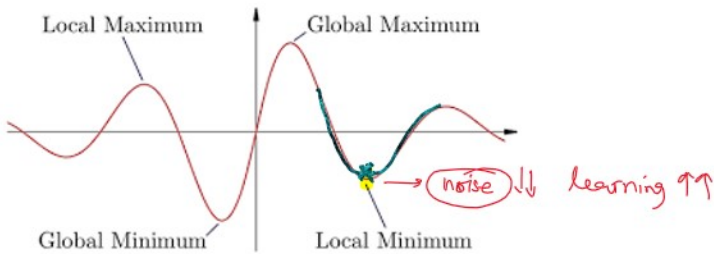


Mean Gradient \rightarrow \downarrow time



steps

1. Pick a mini-batch (batch-size = 32)
2. Calculate the mean gradient of the mini-batch.
3. Update the model parameters (weights & biases) using mean gradient
4. Repeat the steps 1-3 for each epoch.

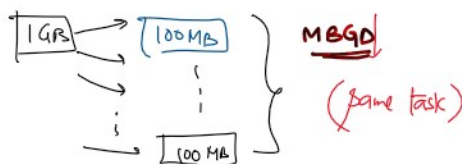


Comparison between SGD / BGD / MBGD

1) computationally expensive

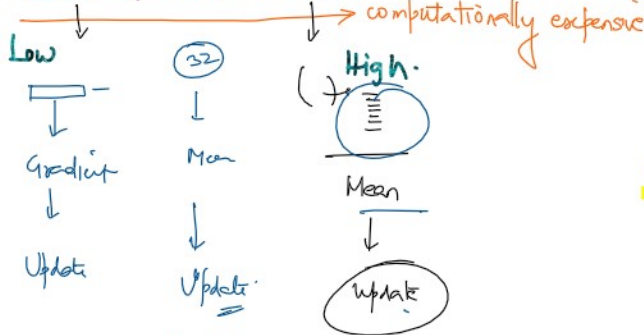
BGD \leftarrow 1GB : dataset \rightarrow 5 million rows

Group by \rightarrow Mean Gradient



$$\begin{bmatrix} 1 \text{ row} \\ 1 \\ \vdots \\ 51 \text{ rows} \end{bmatrix} \rightarrow \underline{\underline{\text{SGD}}}$$

SGD < MBGD < BGD



No. of ^{# for laps} iterations performed in implementing 3 GDs (refer the Google collab notebook)

$$m = 100$$

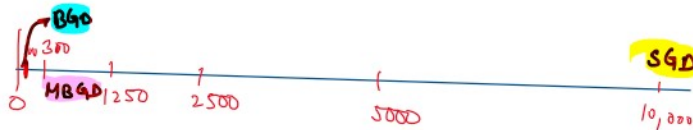
$$\text{epochs} = 100$$

$$n_{\text{batch}} = \left\lceil \frac{100}{32} \right\rceil = \left\lceil 3.125 \right\rceil = 3$$

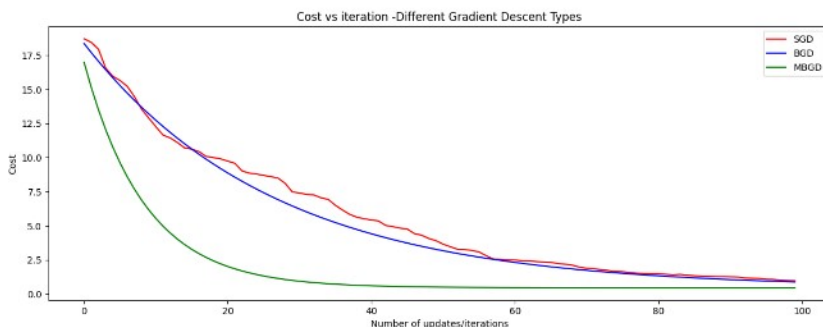
$$\text{SGD \#} (\text{training records} \times \text{epochs}) = m \times \text{epoch} = 100 \times 100 = 10,000$$

$$\text{BGD \#} = \text{epochs} = 100$$

$$\text{MBGD \#} = n_{\text{batch}} \times \text{epochs} = 3 \times 100 = 300$$



Aspect	SGD	BGD	MBGD
Update Frequency	Once per training example $m=100$	Once per epoch	Multiple times per epoch
Convergence Speed	Fast (every row)	Slow (entire dataset/epoch)	Fast (than SGD) \rightarrow batch-size (2)
Memory Usage	Low (single example)	High (entire dataset)	Moderate (mini-batch)
Update Stability	Less stable (one row)	Very stable	Moderately stable
Computational Cost	Low	High	Moderate
Scalability	Good for large datasets	Poor for large datasets	Good for large datasets
Noise in Updates	High (more noisy)	Low (less noisy)	Medium (less noisy than SGD)
Implementation Complexity	Simple	Simple	Slightly more complex



$$n_{\text{batch}} = 32$$

right batch size = ?

Hyper parameter tuning