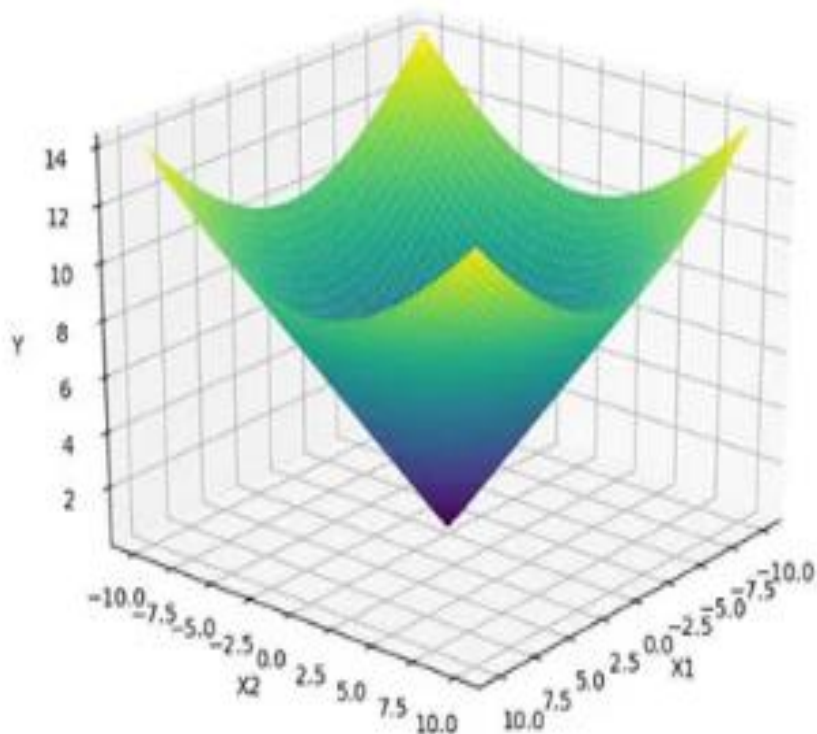


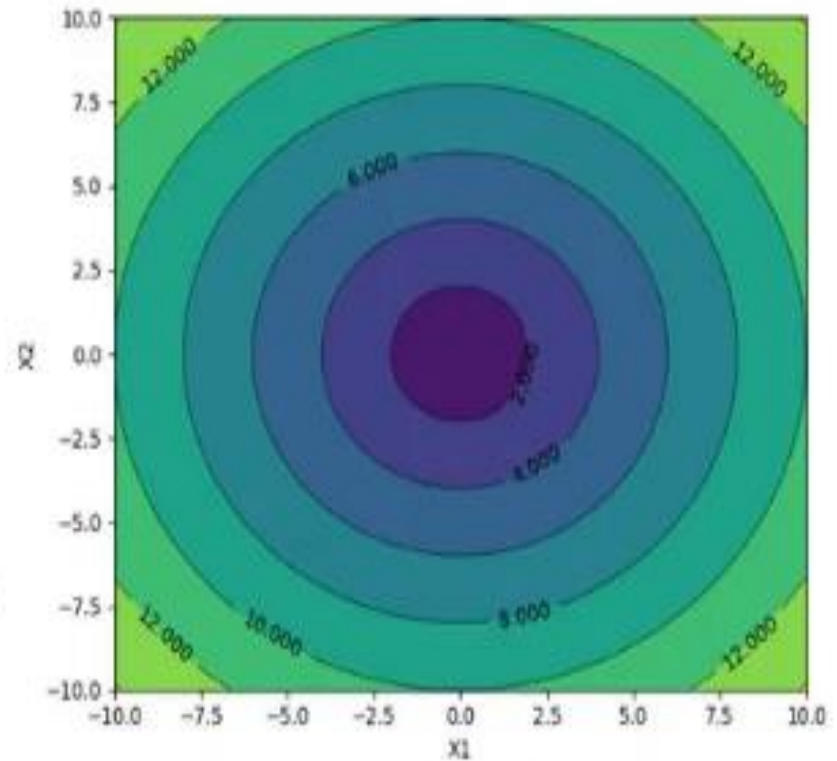
ANN AND DEEP LEARNING (CS636)

BY: Nidhi S. Periwal,
Teaching Assistant,
COED, SVNIT, Surat

Representation of Gradient in form of 2D Contour



3D Plot

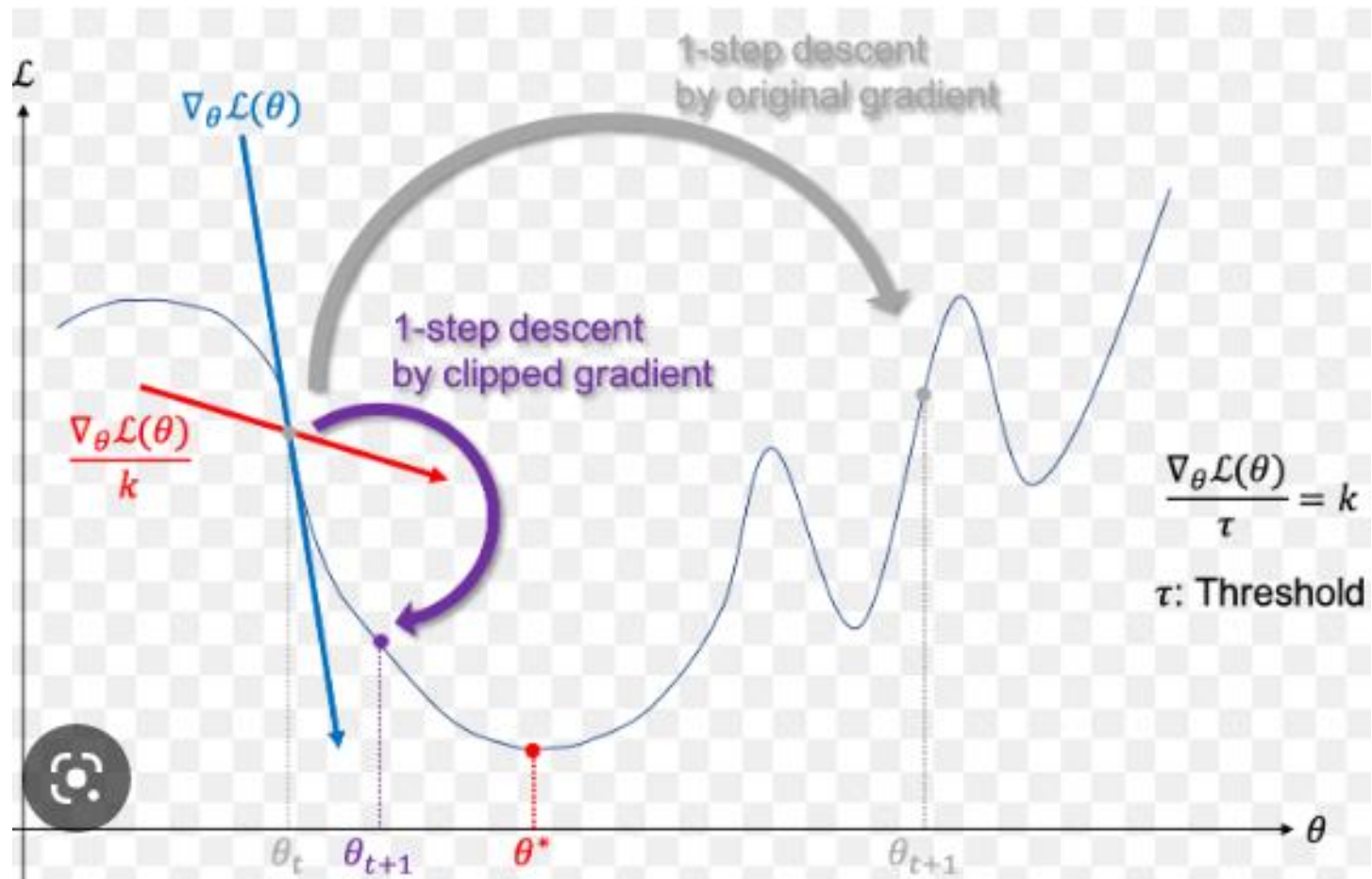


Contour Plot

Gradient Clipping

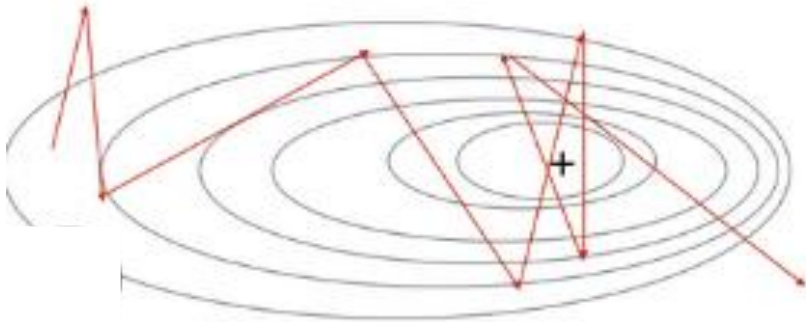
- Gradient Clipping is a method where the error derivative is changed or clipped to a threshold during backward propagation through the network, and using the clipped gradients to update the weights.
- **Gradient clipping** is a technique to prevent exploding gradients in very deep networks, usually in recurrent neural networks.
- *When the traditional gradient descent algorithm proposes to make a **very large step**, the **gradient clipping heuristic** intervenes to reduce the step size to be small enough that it is less likely to go outside the region where the gradient indicates the direction of approximately steepest descent.*
- It is a method that only addresses the **numerical stability of training deep neural network** models and does **not** offer any general improvement in **performance**.

Gradient Clipping

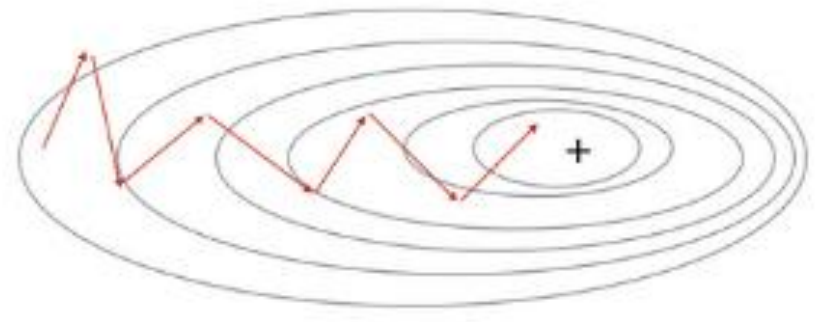


Gradient Clipping

Without gradient clipping



With gradient clipping

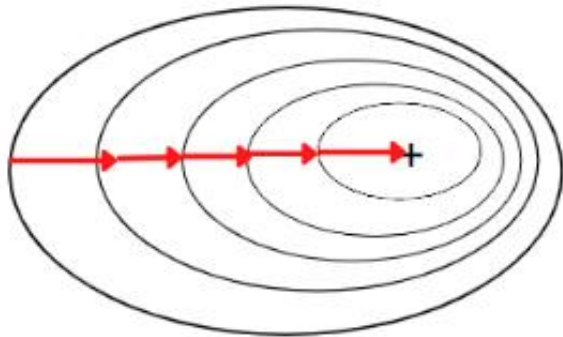


Challenging Optimization

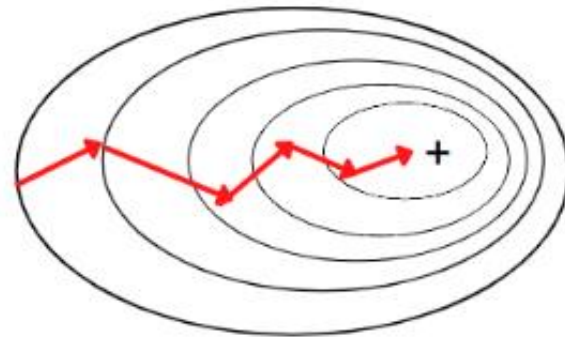
- Training deep learning neural networks is very challenging.
- The best general algorithm known for solving this problem is stochastic gradient descent, where model weights are updated each iteration using the backpropagation of error algorithm.

Types of Gradient Descent

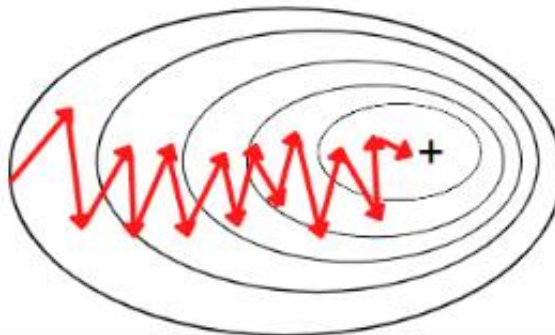
Batch Gradient Descent

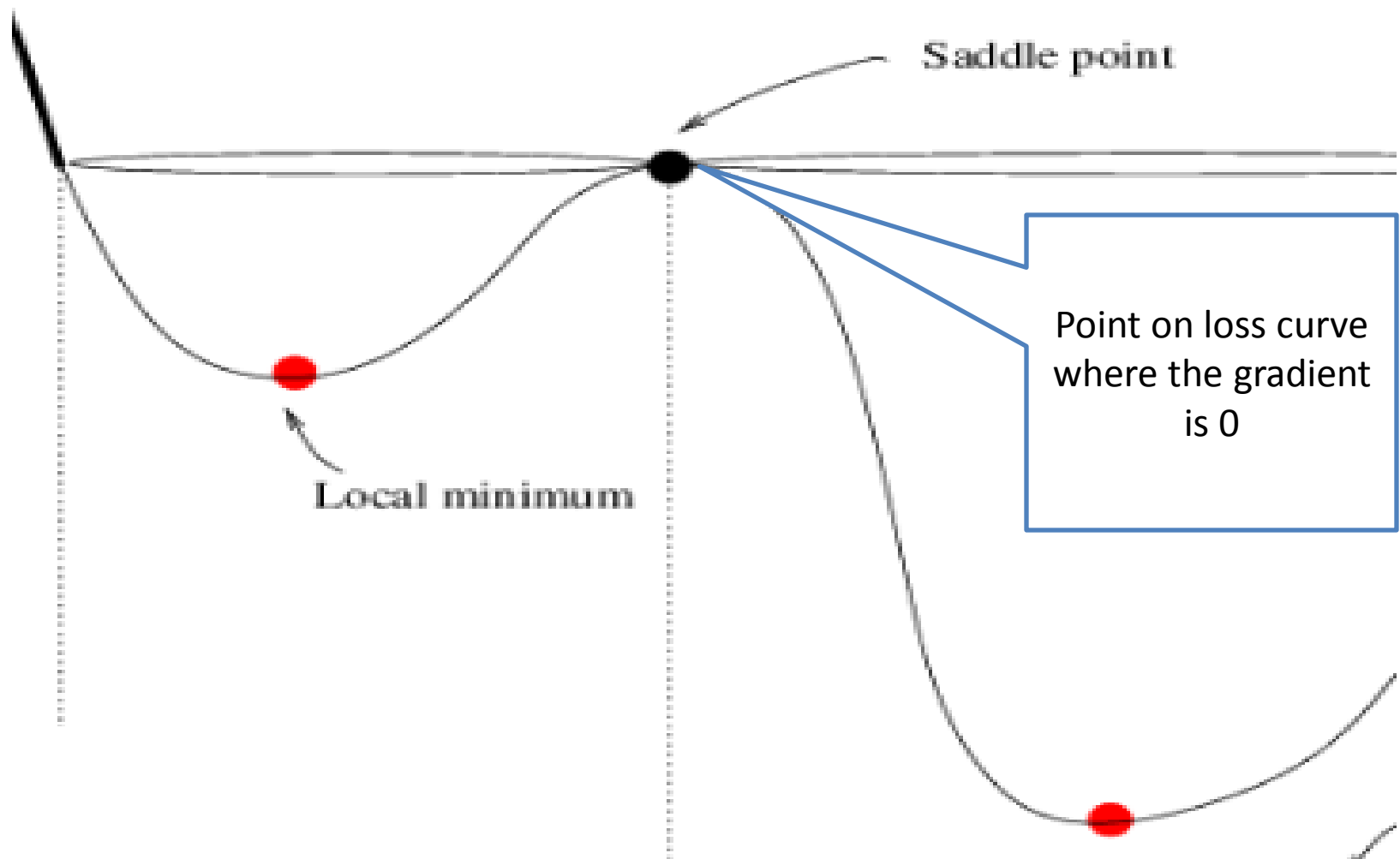


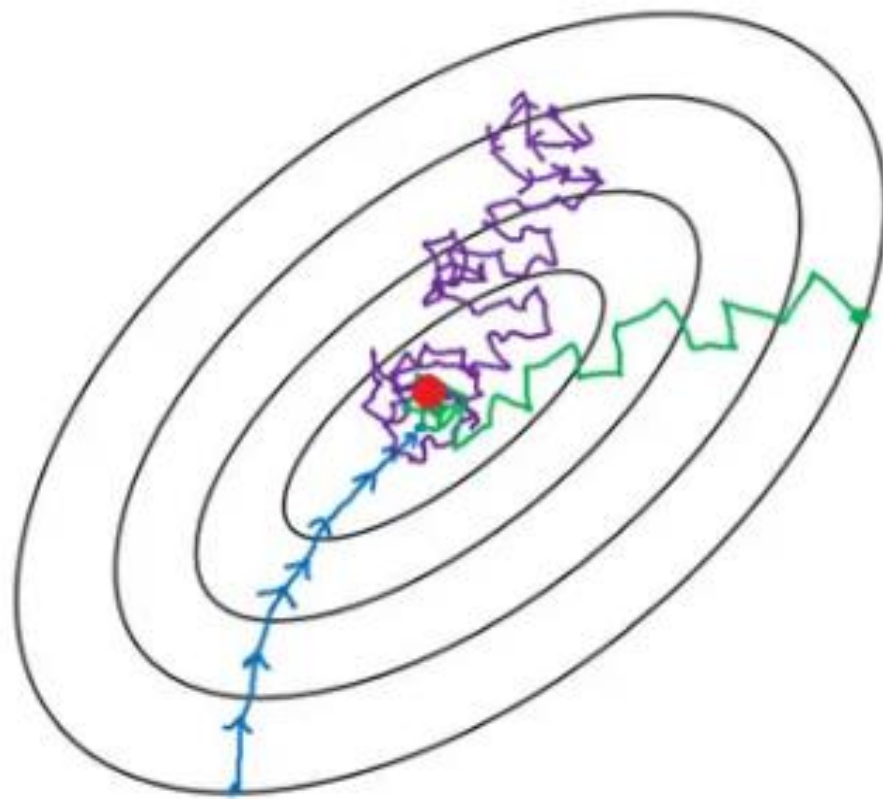
Mini-Batch Gradient Descent



Stochastic Gradient Descent







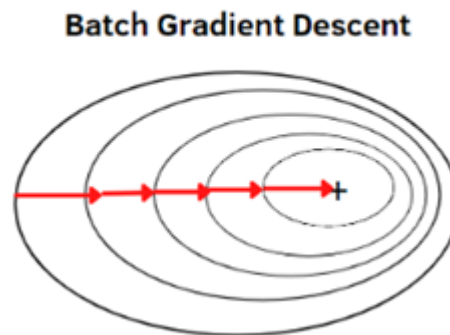
- Batch gradient descent
- Mini-batch gradient Descent
- Stochastic gradient descent

GRADIENT DESCENT

- Gradient descent is an **iterative algorithm whose purpose is to make changes to a set of parameters in hopes of reaching an optimal set of parameters that leads to the lowest loss function value possible.**
- The gradient descent algorithm is an **optimization algorithm mostly used in machine learning and deep learning.**
- Gradient descent adjusts parameters to minimize particular functions to local minima.
- The algorithm objective is to identify model parameters like weight and bias that reduce model error on training data.
- **A gradient measures how much the output of a function changes if you change the inputs a little bit.**
- In machine learning, a gradient is a derivative of a function that has more than one input variable. Known as the slope of a function in mathematical terms, the **gradient simply measures the change in all weights about the change in error.**

BATCH GRADIENT DESCENT

- Batch gradient descent, also known as vanilla gradient descent, calculates the error for each example within the training dataset.
- Batch gradient descent calculates error. Still, the model is not changed until every training sample has been assessed. The entire procedure is referred to as a cycle and a training epoch.
- **Batch gradient descent calculates the error based on each training record but updates the parameters after evaluating all training records.**
- Benefits of batch are its **computational efficiency, which produces a stable error gradient and a stable convergence.**



BATCH GRADIENT DESCENT

Advantages

- Fewer model updates mean that this variant of the steepest descent method is more **computationally efficient** than the stochastic gradient descent method.
- Reducing the update frequency provides a **more stable error gradient and a more stable convergence for some problems.**
- Separating forecast error calculations and model updates provides a **parallel processing-based algorithm implementation.**

Disadvantages

- A more stable error gradient can cause the model to **prematurely converge to a suboptimal set of parameters.**
- End-of-training epoch updates require the additional complexity of accumulating prediction errors across all training examples- **Local Minima**
- **Saddle Point Problem.**
- The batch gradient descent method typically requires the entire training dataset in **memory** and is implemented for use in the algorithm.
- **Large datasets** can result in very slow model updates or training speeds.
- Slow and require more computational power.

STOCHASTIC GRADIENT DESCENT (SGD)

- Stochastic gradient descent (SGD) **changes the parameters for each training sample one at a time for each training example in the dataset.** Depending on the issue, this can make SGD faster than batch gradient descent.
- One benefit is that the regular updates give us a **fairly accurate idea of the rate of improvement.**
- However, the batch approach is less **computationally expensive** than the frequent updates. The frequency of such updates can also produce **noisy gradients**, which could cause the error rate to fluctuate rather than gradually go down.

STOCHASTIC GRADIENT DESCENT (SGD)

Advantages

- **Overcome Saddle Point issue to some extent as compared to Gradient Descent.**
- Model's performance and improvement rates with frequent updates.
- This variant of the steepest descent method is probably the easiest to understand and implement.
- Increasing the frequency of model updates will allow you to learn more about some issues faster.
- **Faster and require less computational power.**
- **Suitable for the larger dataset.**

STOCHASTIC GRADIENT DESCENT (SGD)

Disadvantages

- Frequent model updates are more **computationally intensive than other steepest descent configurations**, and it takes **considerable time to train the model with large datasets**.
- Frequent updates can result in noisy gradient signals. This can result in model parameters and cause errors to fly around (more variance across the training epoch).
- A noisy learning process along the error gradient can also make it difficult for the algorithm to commit to the model's minimum error.
- Suffers from problem of **oscillating trajectory of descent**.

MINI-BATCH GRADIENT DESCENT

- Mini-batch gradient descent combines the ideas of batch gradient descent with SGD, it is the preferred technique.
- It divides the training dataset into manageable groups and updates each separately. **This strikes a balance between batch gradient descent's effectiveness and stochastic gradient descent's durability.**
- Mini-batch sizes typically range from 50 to 256, although, like with other machine learning techniques, there is no set standard because **it depends on the application.**
- The most popular kind in deep learning, this method is used when training a neural network.

MINI-BATCH GRADIENT DESCENT

Advantages

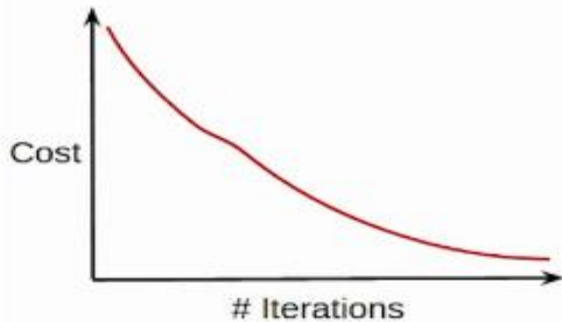
- The model is updated more frequently than the gradient descent method, allowing for **more robust convergence and avoiding local minima.**
- Batch updates provide a **more computationally efficient process than stochastic gradient descent.**
- Batch processing allows for both the efficiency of not having all the training data in memory and implementing the algorithm.

Disadvantages

- **Mini-batch requires additional hyperparameters “mini-batch size” to be set for the learning algorithm.**
- Error information should be accumulated over a mini-batch of training samples, such as batch gradient descent.
- **It will generate complex functions.**

Batch Gradient Descent

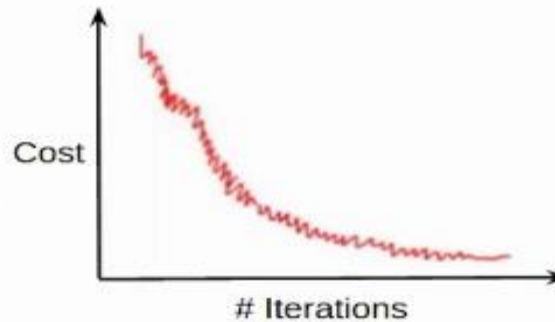
- Entire dataset for updation
- Cost function reduces smoothly



- Computation cost is very high

Stochastic Gradient Descent (SGD)

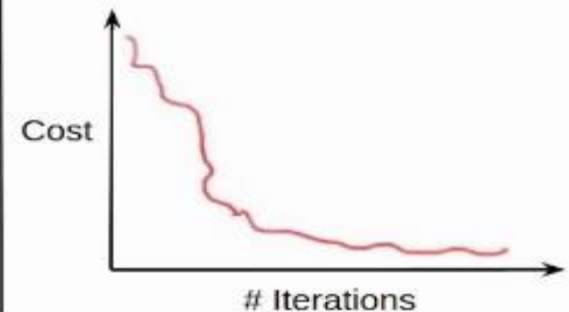
- Single observation for updation
- Lot of variations in cost function



- Computation time is more

Mini-Batch Gradient Descent

- Subset of data for updation
- Smoother cost function as compared to SGD



- Computation time is lesser than SGD
- Computation cost is lesser than Batch Gradient Descent

Sum it up...

- **The mini-batch steepest descent method is the recommended method because it combines the concept of batch steepest descent with SGD.** Simply divide your training dataset into manageable groups and update each individually. This balances the effectiveness of batch gradient descent with the durability of stochastic gradient descent.
- When using batch gradient descent, adjustments are made after calculating the error for a certain batch. **One advantage of the batch gradient descent method is its computational efficiency, which produces a stable error gradient and a stable convergence.**
- **Stochastic Gradient Descent (SGD) sequentially modifies the parameters of each training sample in each training sample of the dataset. This allows SGD to be faster than batch gradient descent. One benefit is that the regular updates give us a fairly accurate idea of the rate of improvement.**

Real World...

- Gradient Descent may never reach the exact minima. It reaches to points close to minima, hence loss value does not change much.
- Hence, networks are generally designed to run for a specific number of iterations (epochs) but stopped **when the loss values stop to improve.**

Other Optimization Algorithms

- Gradient Descent with Momentum
- Adagrad
- Adadelata
- RMSprop
- Adam

Gradient Descent with Momentum

- Momentum is an extension to the gradient descent optimization algorithm, often referred to as **gradient descent with momentum**.
- **It is designed to accelerate the optimization process, e.g. decrease the number of function evaluations required to reach the optima, or to improve the capability of the optimization algorithm, e.g. result in a better final result.**
- A problem with the gradient descent algorithm is that the progression of the search can **bounce around the search space based on the gradient**.
- For example, the search may progress downhill towards the minima, but during this progression, it may move in another direction, even uphill, depending on the gradient of specific points (sets of parameters) encountered during the search.
- This can **slow down** the progress of the search, especially for those optimization problems where the broader trend or shape of the search space is more useful than specific gradients along the way.

Gradient Descent with Momentum

- **Agenda:** To add history to the parameter update equation based on the gradient encountered in the previous updates.
- This change is based on the metaphor of momentum from physics where **acceleration in a direction can be accumulated from past updates**.
- Momentum involves adding an additional hyperparameter that controls the amount of history (momentum) to include in the update equation, i.e. the step to a new point in the search space.
- The value for the hyperparameter is defined in the range 0.0 to 1.0 and often has a value close to 1.0, such as 0.8, 0.9, or 0.99. A momentum of 0.0 is the same as gradient descent without momentum.

Gradient Descent with Momentum

- The target function $f()$ returns a score for a given set of inputs, and the derivative function $f'()$ gives the derivative of the target function for a given set of inputs. A starting point (x)
- We move against the gradient
$$x = x - \text{step_size} * f'(x)$$
- Gradient descent update equation down into two parts: **the calculation of the change to the position and the update of the old position to the new position.**
- The change in the parameters is calculated as the gradient for the point scaled by the step size.
$$\text{change_x} = \text{step_size} * f'(x)$$
- The new position is calculated by simply subtracting the change from the current point
$$x = x - \text{change_x}$$

Gradient Descent with Momentum

- The update at the current iteration (t) will add the change used at the previous iteration ($t-1$) weighted by the momentum hyperparameter, as follows:

$$\text{change_x}(t) = \text{step_size} * f'(x(t-1)) + \text{momentum} * \text{change_x}(t-1)$$

- The update to the position is then performed as before.
$$x(t) = x(t-1) - \text{change_x}(t)$$
- The change in the position accumulates magnitude and direction of changes over the iterations of the search, proportional to the size of the **momentum hyperparameter**.
- Eg: A large momentum (e.g. 0.9) will mean that the update is strongly influenced by the previous update, whereas a modest momentum (0.2) will mean very little influence.

Gradient Descent with Momentum

- Momentum has the effect of dampening down the change in the gradient and, in turn, the step size with each new point in the search space.
- Momentum is most useful in optimization problems where the objective function has a **large amount of curvature (e.g. changes a lot)**, meaning that the gradient may change a lot over relatively small regions of the search space.
- Momentum is helpful when the search space is **flat or nearly flat, e.g. zero gradient**.

Gradient Descent with Momentum

For gradient descent with momentum, the weight update in the i^{th} iteration is done as follows:

$$m_{wi} = \beta_1 m_{wi-1} + (1 - \beta_1) \partial W$$

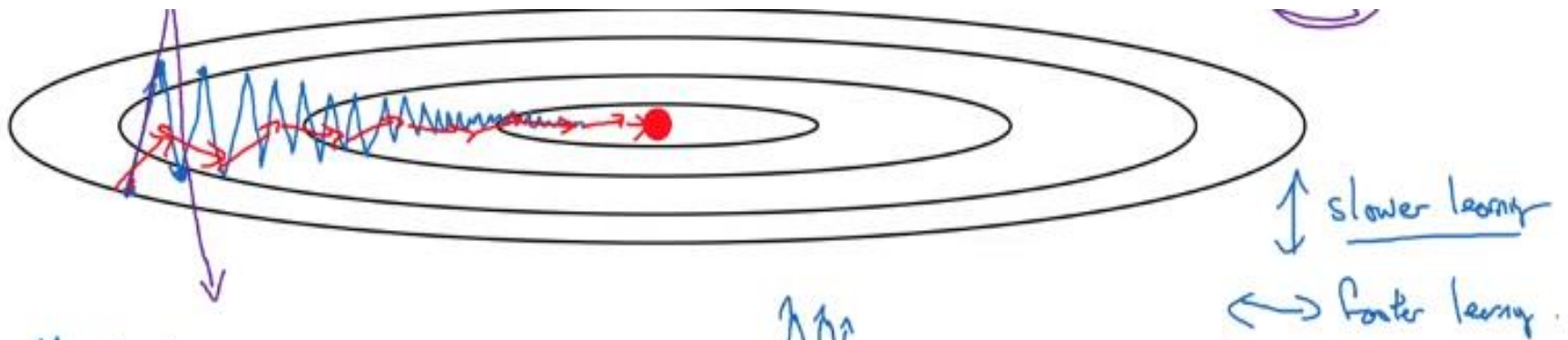
$$W_i = W_{i-1} - \alpha m_{wi}$$

Similarly, the bias update in the i^{th} iteration is done using the equations below:

$$m_{bi} = \beta_1 m_{bi-1} + (1 - \beta_1) \partial b$$

$$b_i = b_{i-1} - \alpha m_{bi}$$

- β_1 – Momentum hyperparameter
- **Generally- 0.9**



Momentum:

On iteration t :

Compute $\Delta W, \Delta b$ on current mini-batch.

$$V_{\Delta W} = \beta V_{\Delta W} + (1-\beta) \Delta W$$

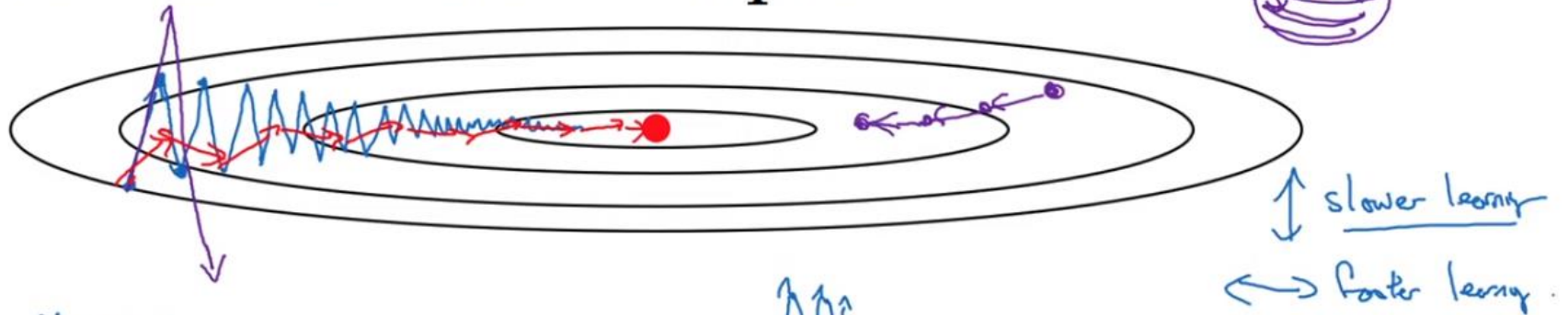
$$V_{\Delta b} = \beta V_{\Delta b} + (1-\beta) \Delta b$$

$$W := W - \alpha V_{\Delta W}, \quad b := b - \alpha V_{\Delta b}$$

$$V_{\theta} = \beta V_{\theta} + (1-\beta) \theta_t$$



Gradient descent example



Momentum:

On iteration t :

Compute $\Delta W, \Delta b$ on current mini-batch.

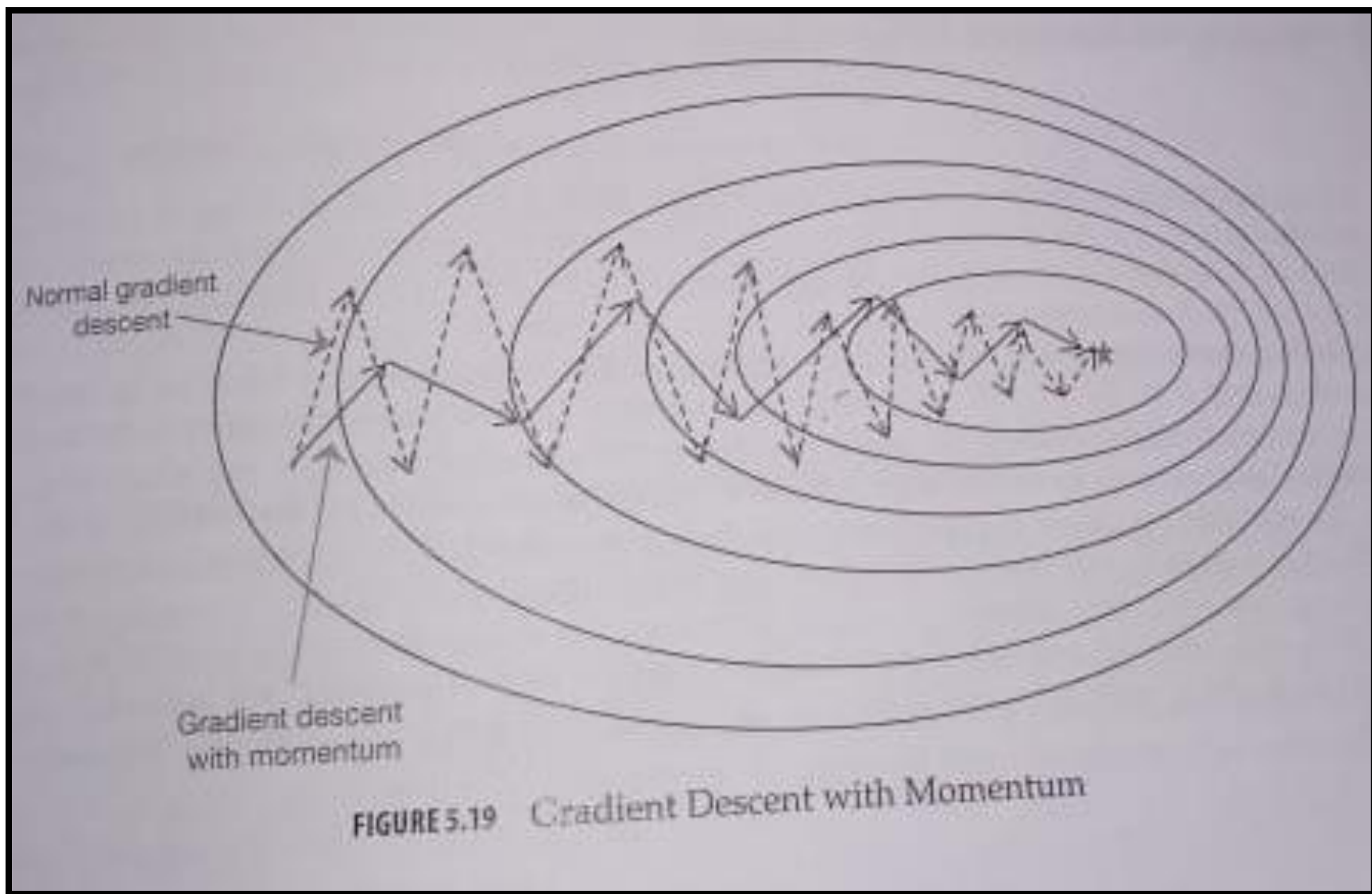
$$V_{\Delta W} = \beta V_{\Delta W} + (1-\beta) \frac{\Delta W}{\alpha}$$

$$V_{\Delta b} = \beta V_{\Delta b} + (1-\beta) \Delta b$$

Friction \rightarrow β velocity \leftarrow acceleration

$$W := W - \alpha V_{\Delta W}, \quad b := b - \alpha V_{\Delta b}$$

$$V_{\theta} = \beta V_{\theta} + (1-\beta) \theta_t$$



- GD with Momentum

$$\theta_{t+1} = \theta_t - \boxed{\eta} \cdot g_t$$

Learning rate always Constant

- Adagrad

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\alpha_t + \epsilon}} \cdot g_t$$

Any Loss Function $f(w, b)$



Partial Differentiation

w.r.t weights $(\frac{\partial f}{\partial w})$



$$w_{\text{new}} = w_{\text{old}} - lr * \frac{\partial f}{\partial w}$$

Partial Differentiation

w.r.t biases $(\frac{\partial f}{\partial b})$



$$b_{\text{new}} = b_{\text{old}} - lr * \frac{\partial f}{\partial b}$$

Adagrad Optimizer

- **Adagrad** stands for **Adaptive Gradient Optimizer**.
- **Issue in GD with Momentum:** Learning Rate is same. For higher Learning rates, there is possibility of overshooting minima. If LR is lower, then takes long time to train.
- Hence, during initial iterations of GD, LR needs to be higher. As GD iterations progress and come closer to minima, the LR needs to be reduced.
- **Adagrad helps in automating the LR tuning process** by using an adaptive LR tuning process using an adaptive LR method based on an accumulative value of **historical** gradients.

Adagrad Optimizer

$$\alpha_i = \frac{\alpha}{\sqrt{\delta + \sum_{k=1}^{i-1} \partial W_k^2}}$$

$$W_i = W_{i-1} - \alpha_i \partial W$$

$$\alpha_i = \frac{\alpha}{\sqrt{\delta + \sum_{k=1}^{i-1} \partial b_k^2}}$$

$$b_i = b_{i-1} - \alpha_i \partial b$$

- Value of delta is very small 10^{-5} to 10^{-7} to avoid divide by 0
- It is suitable for **sparse data** .

$$-W_t = W_{t-1} - \boxed{\eta_t} \times \frac{\partial L}{\partial W_{t-1}}$$

$$\eta'_t = \frac{\eta}{\sqrt{d_t + \epsilon}}$$

$d_t = \sum_{i=1}^t \left(\frac{\partial L}{\partial W_i} \right)^2$

ϵ → small +ve number

g_t^2

In Every Iteration we add square of gradients ($\frac{\partial f^2}{\partial \theta}$) in alpha so it will have history of all past gradients

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\alpha_t + \epsilon}} \cdot g_t$$

Then we can divide learning rate with this term. Now, learning rate will vary with every iteration

$\epsilon = 10^{-8}$ is a term added to alpha to avoid zero division error

Ada grad

$$w_{i,t} = w_{i,t-1} - \eta'_t g_t$$

$$\eta'_t = \frac{\eta}{\sqrt{K + \epsilon}} \quad \eta = 0.1$$

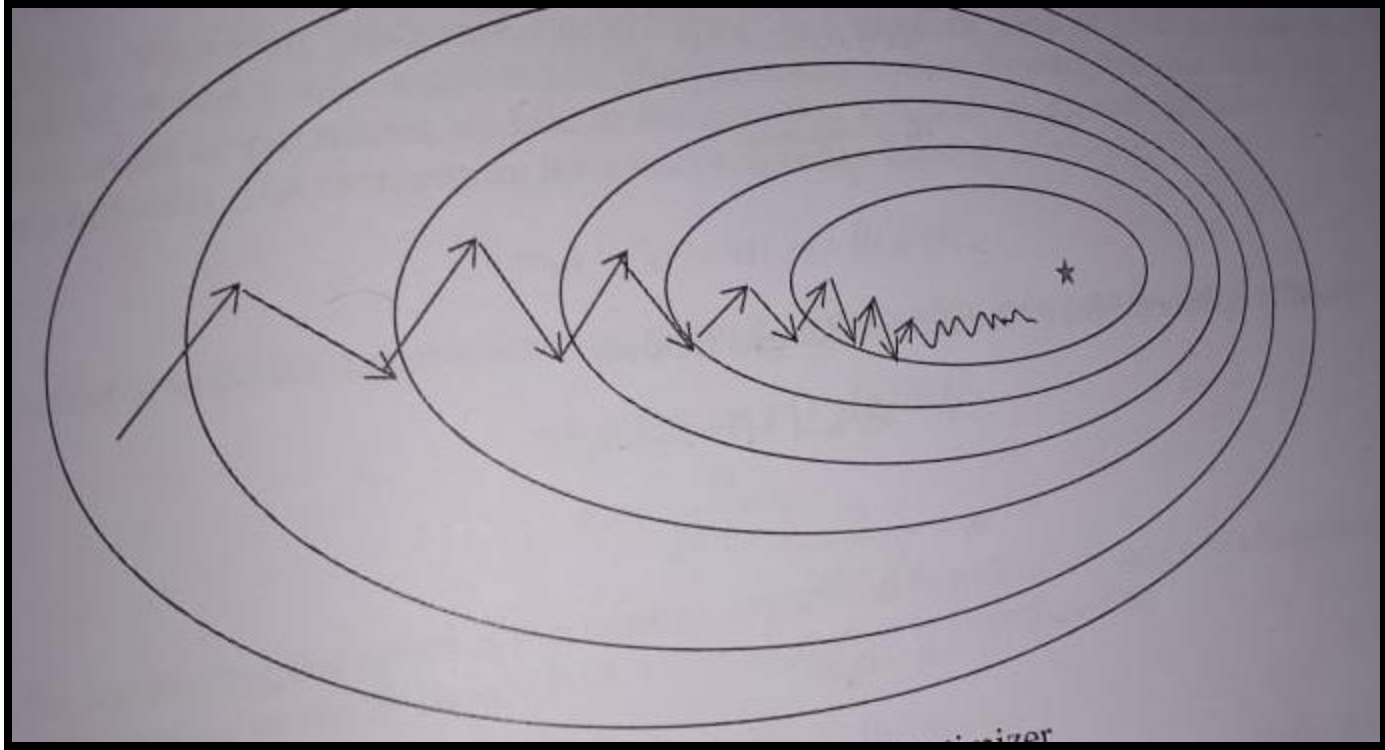
$\epsilon = \text{Small Value}$

$$K = \sum_{i=1}^{t-1} g_i^2$$

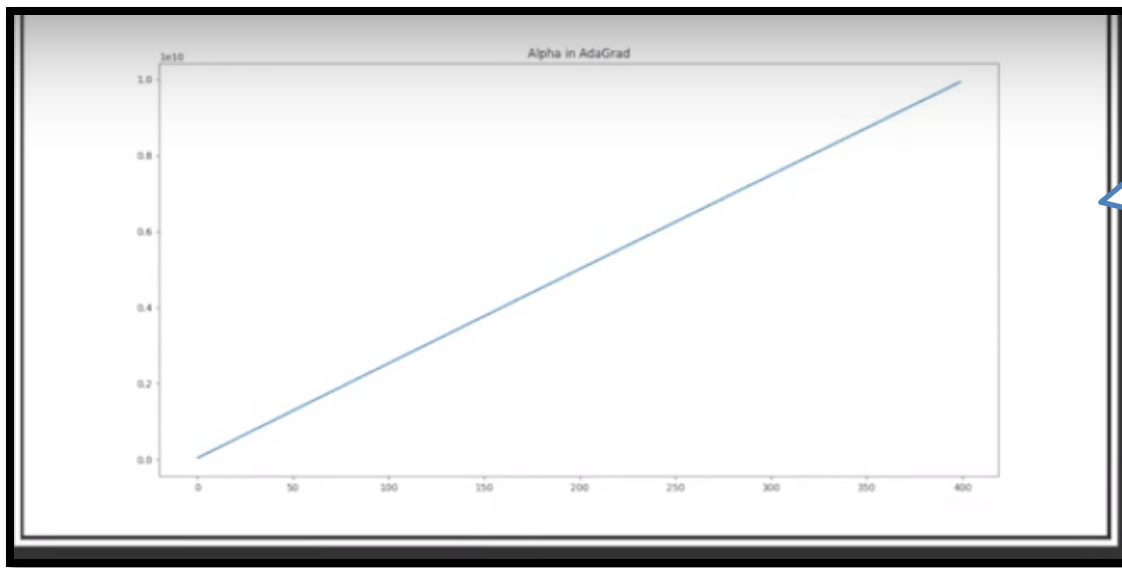
$$\theta_{t+1} = \theta_t - \boxed{\frac{\eta}{\sqrt{\alpha_t + \epsilon}}} \cdot g_t$$

This term gets smaller with every iteration as α is increasing. Thus, this continual decay of learning rate is the weakness of AdaGrad

Adagrad Optimizer

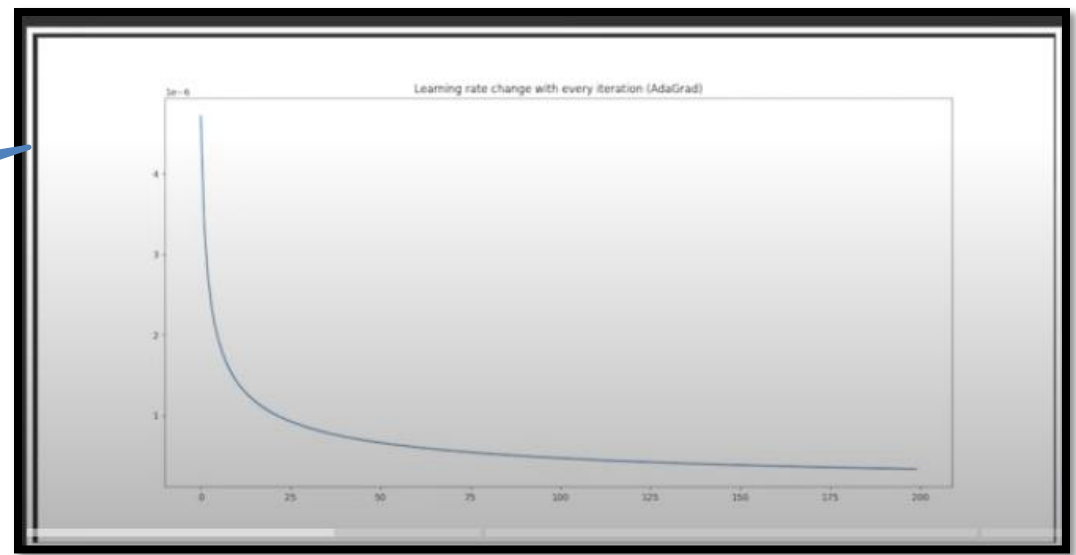


- Limitation: **Learning rate decay over iterations due to sum of squared gradients in denominator.** Hence, after certain iterations algorithm remains stuck forever.



Alpha in Adagrad

**Learning Rate decay over iterations
in Adagrad**



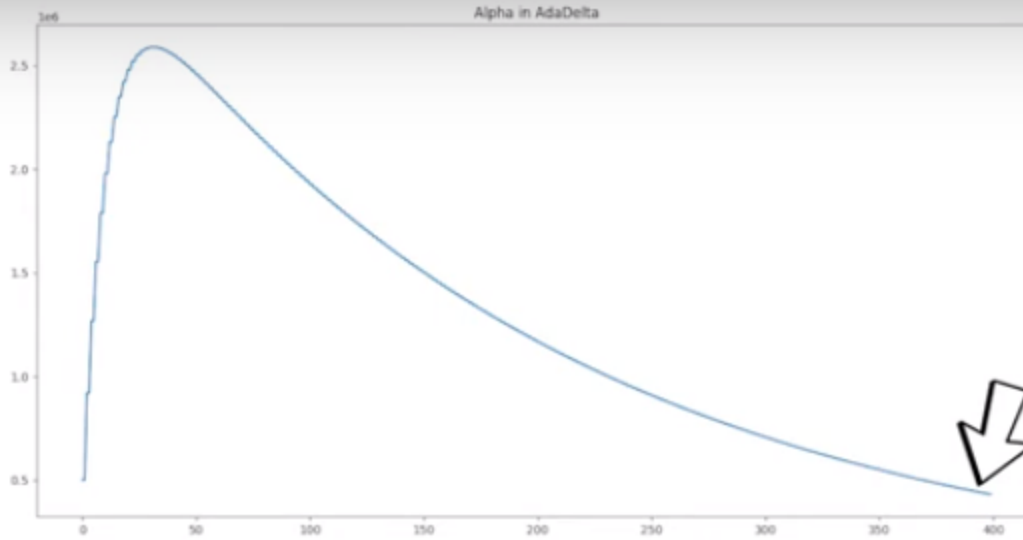
AdaDelta Algo

- Adadelta is similar to Adagrad Algorithm.
- It makes a small modification in calculation of weights across different iterations.
- **Instead of squared gradients of all past iterations, it takes the gradients for a specific windows or number of iterations.**
- Eg: If window size used is 3, it will take squared gradients only for 3 preceding iterations.
- **This will help in addressing the issue faced by Adagrad, i.e drastic decay of learning rate over the iterations.**

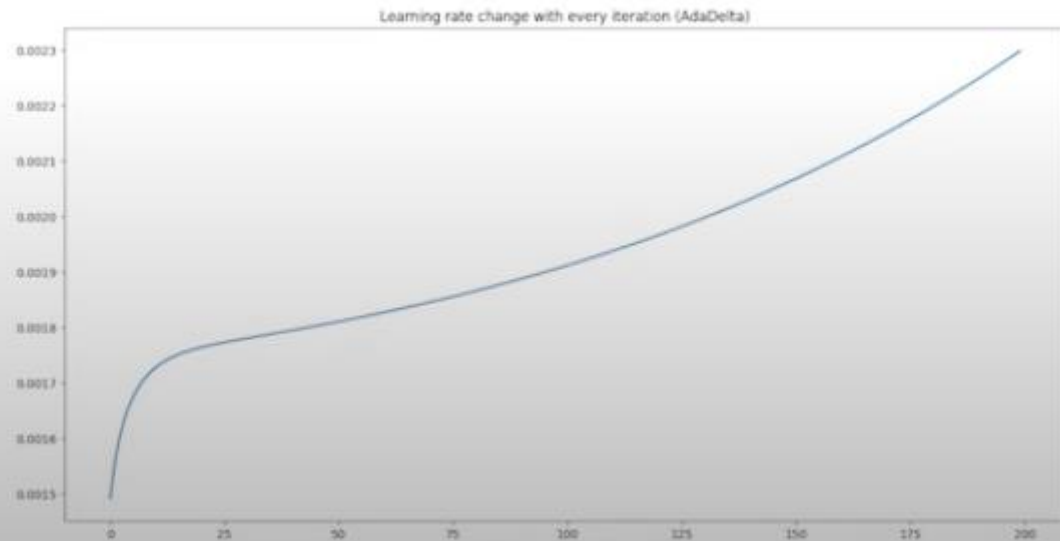
AdaDelta Algo

- Adadelta optimization is a stochastic gradient descent method that is based on adaptive learning rate per dimension to address following drawback:
 - **The continual decay of learning rates throughout training.**
- Adadelta is a more **robust** extension of Adagrad that adapts learning rates based on a moving window of gradient updates, instead of accumulating all past gradients.
- **This way, Adadelta continues learning even when many updates have been done.**

Alpha in Adadelta



**Learning Rate in
Adadelta**



RMSProp Algorithm

- The **RMSProp** algorithm modifies AdaGrad to perform better in the **non-convex** setting by changing the gradient accumulation into an exponentially weighted moving average.
- **Adagrad Disadvantages**
 - Designed to converge rapidly when applied to a convex function.
 - When applied to a non-convex function to train a neural network, the learning trajectory may pass through many different structures and eventually arrive at a region that is a **locally** convex bowl.
 - AdaGrad **shrinks the learning rate according to the entire history of the squared gradient** and may have made the learning rate too small before arriving at such a convex structure.

RMSProp Algorithm

- This algorithm **doesnot take a flat cumulative value of past gradients.**
- It uses a decay factor β_2 which helps to give more weightage to recent gradients and less weights to older gradients. β_2 value is often 0.999.

RMSProp Algorithm

For RMSProp algorithm, the weight update for the i^{th} iteration is done as

$$s_{wi} = \beta_2 s_{wi-1} + (1 - \beta_2) \partial W_i^2$$

$$\alpha_i = \frac{\alpha}{\sqrt{\delta + s_{wi}}}$$

$$W_i = W_i - \alpha_i \partial W$$

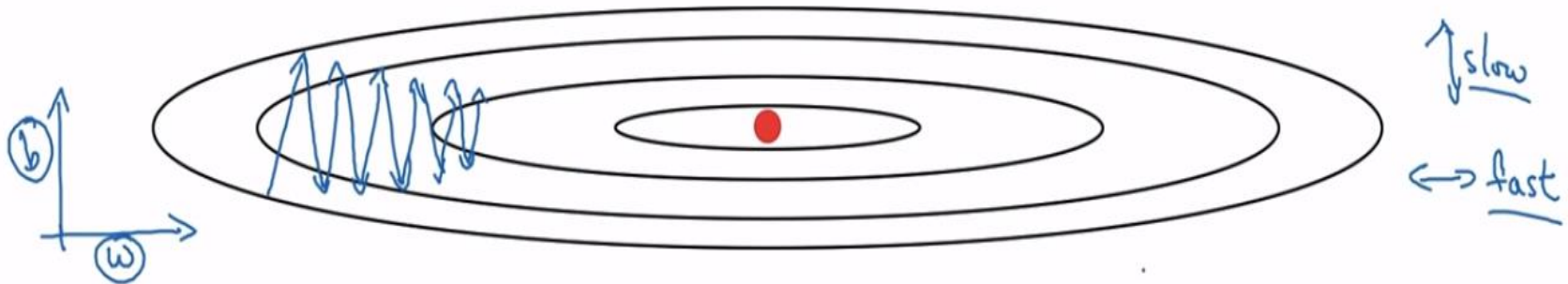
The bias update for the i^{th} iteration is done as follows:

$$s_{bi} = \beta_2 s_{bi-1} + (1 - \beta_2) \partial b_i^2$$

$$\alpha_i = \frac{\alpha}{\sqrt{\delta + s_{bi}}}$$

$$b_i = b_i - \alpha_i \partial b$$

RMSprop



On iteration t :

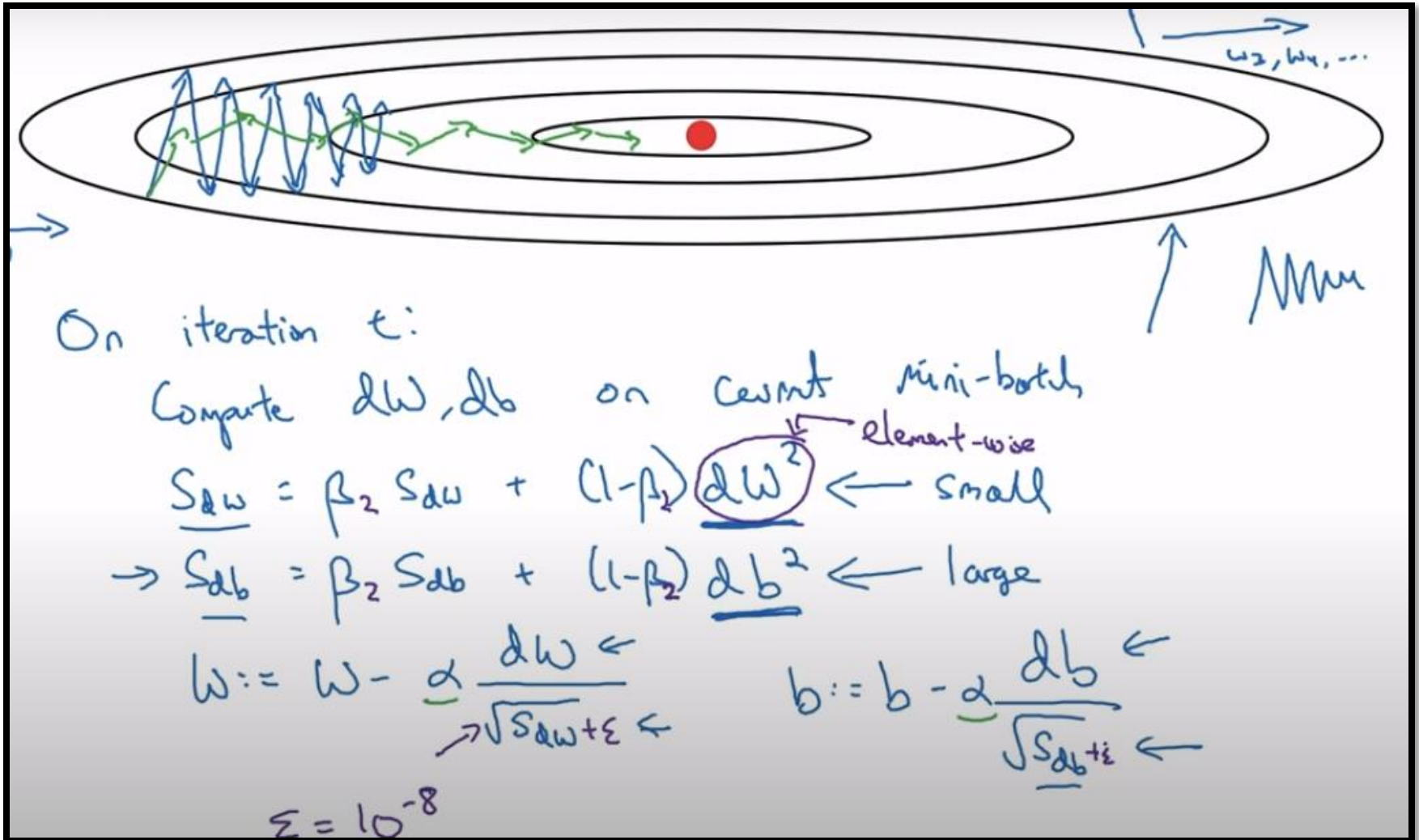
Compute dw, db on current mini-batch

$$S_{dw} = \beta S_{dw} + (1-\beta) \underbrace{dw^2}_{\text{element-wise}} \leftarrow \text{small}$$

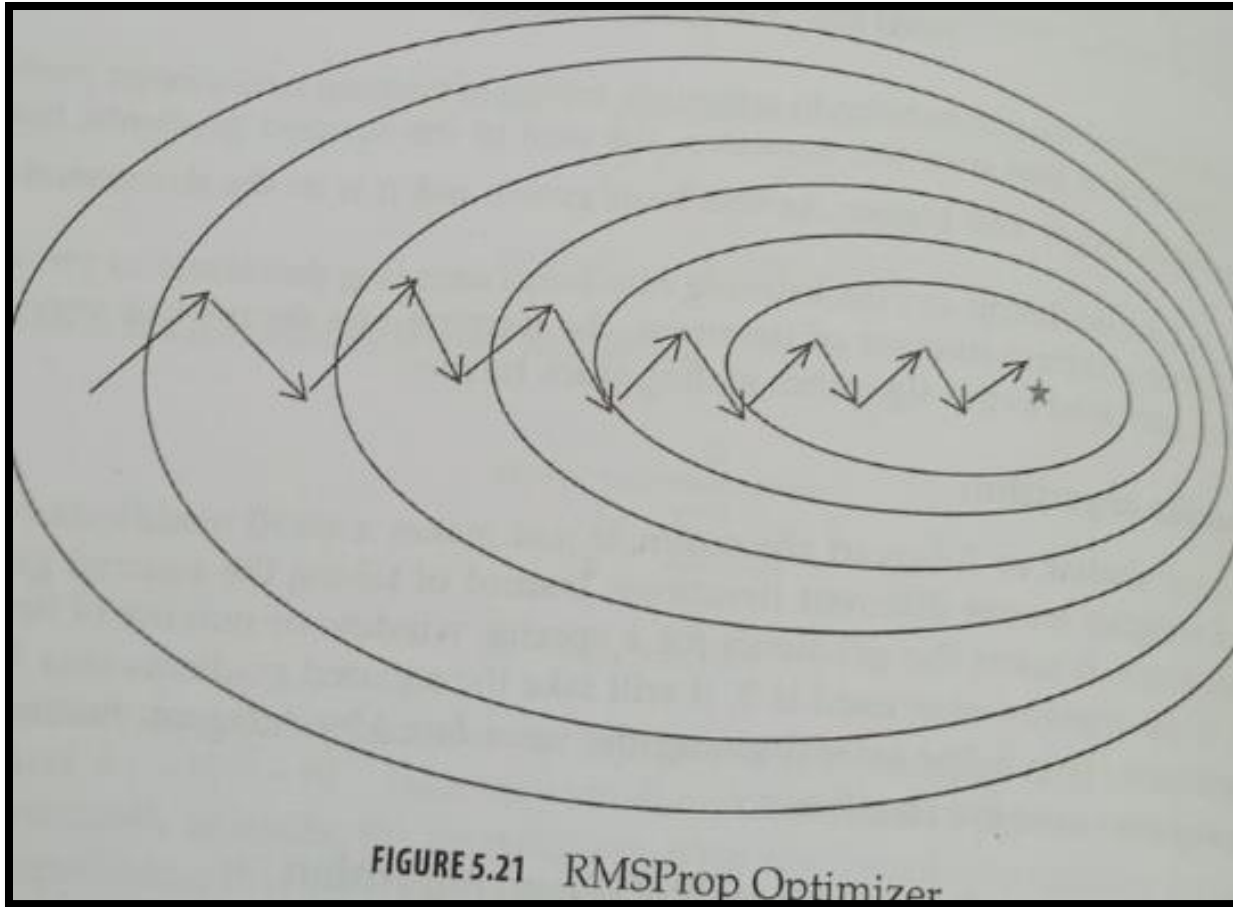
$$S_{db} = \beta S_{db} + (1-\beta) db^2 \leftarrow \text{large}$$

$$w := w - \alpha \frac{dw}{\sqrt{S_{dw}}} \leftarrow \quad b := b - \alpha \frac{db}{\sqrt{S_{db}}} \leftarrow$$

RMSProp Algorithm



RMSProp Algorithm



RMSProp Algorithm

- RMSProp uses an exponentially decaying average to discard history from the extreme past so that it can converge rapidly after finding a convex bowl, as if it were an instance of the AdaGrad algorithm initialized within that bowl.
- **RMSProp has been shown to be an effective and practical optimization algorithm for deep neural networks. It is currently one of the go-to optimization methods being employed routinely by deep learning practitioners.**

Adam

- Adaptive Moment Estimation is an algorithm for optimization technique for gradient descent.
- **The method is really efficient when working with large problem involving a lot of data or parameters.**
- **It requires less memory and is efficient.**
- Intuitively, it is a combination of the 'gradient descent with momentum' algorithm and the 'RMSP' algorithm.
- It also exponentially smooths the first-order gradient in order to incorporate momentum into the update.
- It also directly addresses the bias inherent in exponential smoothing when the running estimate of a smoothed value is unrealistically initialized to 0.

Adam optimization algorithm

$$V_{dw} = 0, S_{dw} = 0, V_{db} = 0, S_{db} = 0$$

On iteration t :

Compute dw, db using current mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) dw, V_{db} = \beta_1 V_{db} + (1 - \beta_1) db \leftarrow \text{"momentum"} \beta_1$$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) dw^2, S_{db} = \beta_2 S_{db} + (1 - \beta_2) db^2 \leftarrow \text{"RMSprop"} \beta_2$$

$$V_{dw}^{\text{corrected}} = V_{dw} / (1 - \beta_1^t), V_{db}^{\text{corrected}} = V_{db} / (1 - \beta_1^t)$$

$$S_{dw}^{\text{corrected}} = S_{dw} / (1 - \beta_2^t), S_{db}^{\text{corrected}} = S_{db} / (1 - \beta_2^t)$$

$$W := W - \alpha \frac{V_{dw}^{\text{corrected}}}{\sqrt{S_{dw}^{\text{corrected}} + \epsilon}}$$

$$b := b - \alpha \frac{V_{db}^{\text{corrected}}}{\sqrt{S_{db}^{\text{corrected}} + \epsilon}}$$

Hyperparameters choice:

→ α : needs to be tune

→ β_1 : 0.9 → (dw)

→ β_2 : 0.999 → (dw^2)

→ ϵ : 10^{-8}

Adam: Adaptive moment estimation



Adam Coates

Adam

So, combining the effect of momentum expression as well as learning rate adjustment, the weight update in the i^{th} iteration of Adam algorithm is formulated as:

$$W_i = W_{i-1} - \alpha_i m_{wi}$$

Adam also does a bias correction. When we start, m and s have a small value close to 0. Down the iterations, their values may grow. To give a dampening effect to the growing values of m and s , the following equations are used to adjust the values of m and s .

$$m_{wi}^{\text{corrected}} = \frac{m_{wi}}{1 - \beta_1^i}$$

$$s_{wi}^{\text{corrected}} = \frac{s_{wi}}{1 - \beta_2^i}$$

$$\therefore \alpha_i^{\text{corrected}} = \frac{\alpha}{\sqrt{\delta + s_{wi}^{\text{corrected}}}}$$

Adam

Applying bias correction, the weight update for the i^{th} iteration can be done as:

$$W_i = W_{i-1} - \alpha_i^{\text{corrected}} m_{wi}^{\text{corrected}}$$

Similarly, the bias update for the i^{th} iteration can be done as:

$$m_{bi} = \beta_1 m_{bi-1} + (1 - \beta_1) \partial b$$

$$s_{bi} = \beta_2 s_{bi-1} + (1 - \beta_2) \partial b_i^2$$

$$\alpha_i = \frac{\alpha}{\sqrt{\delta + s_{bi}}}$$

Adam

$$b_i = b_{i-1} - \alpha_i m_{bi}$$

$$m_{bi}^{\text{corrected}} = \frac{m_{bi}}{1 - \beta_1^i}$$

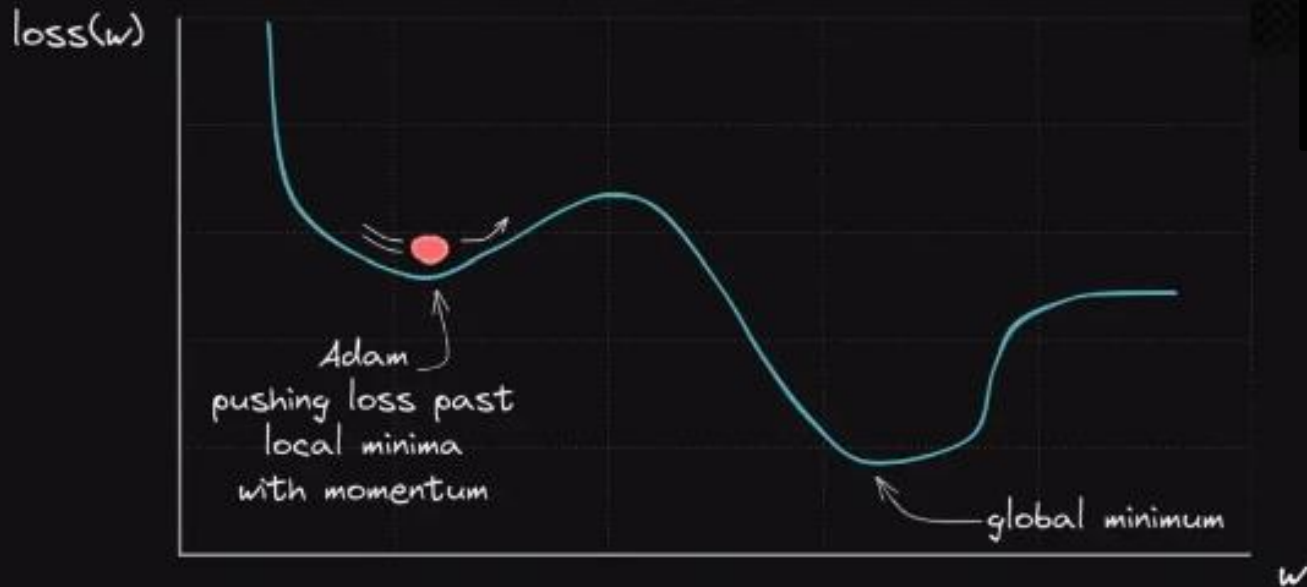
$$s_{bi}^{\text{corrected}} = \frac{s_{bi}}{1 - \beta_2^i}$$

$$\therefore \alpha_i^{\text{corrected}} = \frac{\alpha}{\sqrt{\delta + s_{bi}^{\text{corrected}}}}$$

$$b_i = b_{i-1} - \alpha_i^{\text{corrected}} m_{bi}^{\text{corrected}}$$

ADAM OPTIMIZER

Network Loss
(with Adam optimizer)



Popular optimizers:

- Adagrad
- AdaDelta
- RMSProp
- Adam

Momentum term allows Adam to accelerate towards the minimized loss with little oscillation.

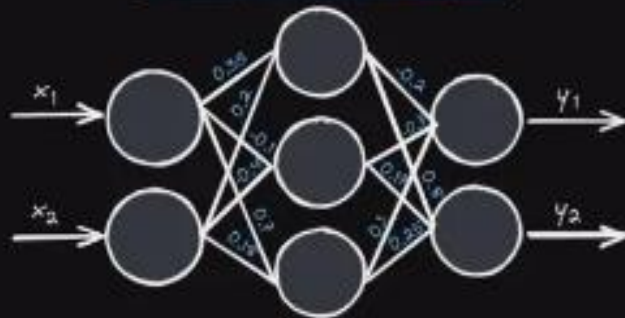
Adaptive learning rates allow Adam to accelerate along flatter regions of the loss curve and slow down along steeper regions.

Optimizer choices....

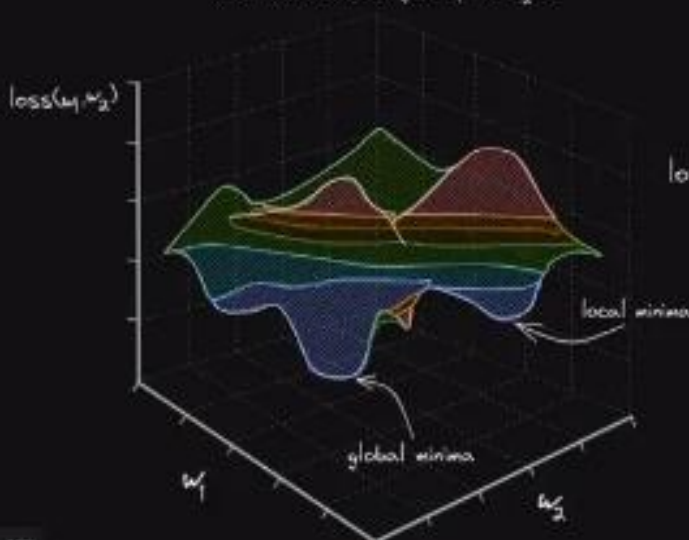
- For straight forward problems, normal gradient descent (mini-batch incase if the volume of data is huge) can be used.
- If it doesn't yields satisfactory result, RMSProp can be used or to speed up the learning process, gradient descent with momentum can be used.
- For large and sparse data like text data, Adagrad or Adadelata can be used.
- Adam is most popularly adopted optimizer by DL community.

OPTIMIZATION ALGORITHMS

Randomly Initialized Weights



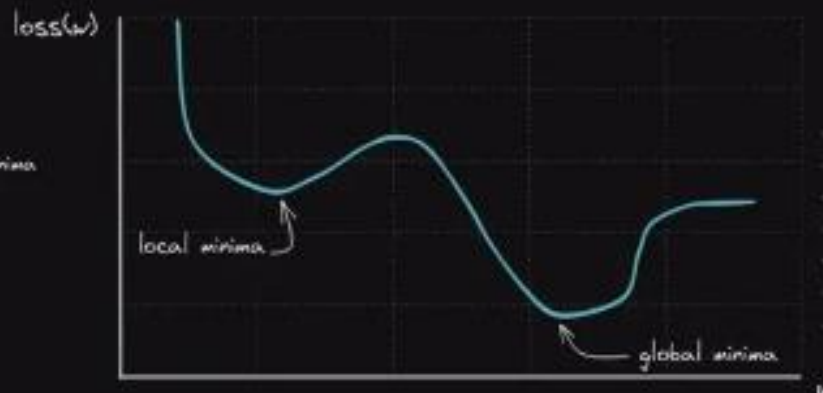
Network Loss
(as a function of weights w_1 and w_2)



Training Steps

1. Forward pass input data through the network.
2. Obtain output from the network.
3. Measure loss between network output and true labels of data.
4. Complete backwards pass to obtain the gradients (backpropagation).
5. Update network weights with the gradients to step towards minimized loss.

Network Loss
(as a function of weight w)



Optimizers

- Batch gradient descent
- Stochastic gradient descent (SGD)
- Mini-batch gradient descent
- Adagrad
- AdaDelta
- RMSProp
- Adam

Thank You!