

## Difference b/w Gradient Descent vs Regularization

↓

(model learns by minimizing a cost function)

(we use this when model overfits)

- ① Regularization methods have a pre-defined cost function, unlike Gradient Descent in which the cost function is the gradient of the given function.
- ② Regularization is used when a model overfits, unlike Gradient Descent which is used regardless of overfitting (we may find the model overfits after its validation but we do not use Gradient Descent to prevent overfitting)

## 6 Gradient Descent

It is an optimization algorithm used for minimizing the cost function in various ML algos.

It is basically used for updating the parameters of the learning model.

### Types of gradient descent:

#### ① Batch Gradient Descent:

(All in one go)

→ All training examples are processed in single iteration.

→ If the no. of training examples is large, then batch gradient descent is computationally expensive & not preferred.

#### ② Stochastic Gradient Descent (One by one)

→ It processes 1 training example each iteration.

Hence, the parameters are updated even after one iteration. Hence, this is quite faster than batch gradient.

→ But, when training examples are high even then it processes only one example which can be additional overhead of the system so the no. of iterations will be quite large.

(in batches) Best

#### ③ Mini-Batch G.D.: → It works faster than both batch & stochastic

→ Here, b. example where b < m are processed per iteration.

→ So, even if the no. of iteration training example is large, it is processed in batches of b training examples in one go.

Conclusion: Mini-Batch G.D. works for larger training examples & that too with lesser no. of iterations.

## 6 Gradient Descent

It is an optimization algorithm used for minimizing the cost function in various ML algos.

It is basically used for updating the parameters of the learning model.

### Types of gradient descent:

① Batch Gradient Descent: (All in one go)

- All training examples are processed in single iteration.
- If the no. of training examples is large, then batch gradient descent is computationally expensive & not preferred.

② Stochastic Gradient Descent. (One by one)

- It processes 1 training example each iteration.
- Hence, the parameters is updated even after one iteration. Hence, this is quite faster than batch gradient.
- But, when training examples are high even then it processes only one example which can be additional overhead of the system & the no. of iterations will be quite large.

③ Mini-Batch Gr.D.: → It works faster than both batch & stochastic

→ Here, b. example where b < m are processed per iteration.

→ So, even if the no. of iteration training example is large, it is processed in batches of b training examples in one go.

Conclusion: Mini Batch G.D. works for larger training examples & that too with lesser no. of iterations.

Convergence trends in diff. variants of Gradient Descents:

- ① In case of Batch Gradient Descent, the algorithm follows
- a straight path towards the minimum. If the cost function is convex, then it converges to global min. &
  - if the C.F. is not convex, then it converges to a local min. Here; the learning rate is typically held constant.
- ② In case of SGD & mini-gradient descent, the algo. does not converge but keeps on fluctuating around the global minimum. Therefore, in order to make it converge, we have to slowly change the learning rate.  
(However, the convergence of SGD is much noisier as in one iteration, it processes only one training example.)

## Optimization techniques for Gradient Descent (class notes)

- ① Gr.D. is an iterative optimization algo; used to find the min. value for a function.
- ② The general idea is to initialize the parameters to random values, and then take small steps in the direction of slope at each iteration.
- ③ Gr.D. is highly used in supervised learning to minimize the error function & find the optimal value of parameters.
- Various extensions have been designed for the gradient descent algo. :

Optimization Techniques for Gr.D.

- 1. Momentum method
- 2. RMSprop
- 3. Adam optimization

- ① Momentum method: this method is used to accelerate the Gr.D. algo. by taking into consideration the exponentially weighted average of the gradients.

Using averages makes the algo. converges towards the minima in a faster way, as the gradients towards the uncommon dirn's are cancelled out.

Pseudocode :

$$v = 0$$

for each iteration  $i$ :

compute  $\partial w$

$$v = \beta v + (1 - \beta) \partial w$$

$$w = w - \alpha v$$

→  $v$  &  $\partial w$  are analogous to  
velocity & acceleration.

→  $\alpha$  is learning rate

→  $\beta$  is momentum.

generally kept 0.9.

RMSprop: The intuition of this method is to apply  
an exponentially weighted average method to the  
[Second moment of the gradient  $\partial w^2$ ].

Pseudocode :

$$s = 0$$

for each iteration  $i$ :

compute  $\partial w$

$$s = \beta s + (1 - \beta) \partial w^2$$

$$w = w - \alpha \left( \frac{\partial w}{\sqrt{s} + \epsilon} \right)$$

Adam optimization : this algo incorporates the momentum  
method & RMSprop along with bias correction.

Pseudo:

$$v = 0$$

$s = 0$

for each iteration  $i$ :

compute  $\partial w$

$$v = \beta_1 v + (1 - \beta_1) \partial w$$

$$s = \beta_2 s + (1 - \beta_2) \partial w^2$$

$$v = v / \{ 1 - \beta_1^{t+1} \}$$

$$s = s / \{ 1 - \beta_2^{t+1} \}$$

$$w = w - \alpha \frac{v}{\sqrt{s} + \epsilon}$$

where;  
 $\beta_1 = 0.9$

## Mode Limitation of G.D.J.

