# Optimization Algorithms
in machine learning

Optimization algorithms play a crucial role in training machine learning models by iteratively **updating the model parameters to minimize a certain objective function**, often referred to as the loss or cost function.

Here are some key optimization algorithms used in machine learning:

# Gradient Descent

**Basic Idea:**
- Update model parameters in the opposite direction of the gradient of the loss function with respect to the parameters.

**Variants:**
- Stochastic Gradient Descent (SGD): Updates parameters using a single random data point at a time.
- Mini-Batch Gradient Descent: Updates parameters using a small batch of randomly selected data points.

# Momentum

**Idea:**
- Introduces a momentum term to prevent oscillations and speed up convergence.

**Update Rule:**
- $v = \beta \cdot v - \alpha \cdot \nabla J(\theta)$
- $\theta = \theta + v$

**Hyperparameter $\beta$:**
- Controls the contribution of the previous update to the current update.

# Adagrad

**_Idea:_**
- Adjusts the learning rates of each parameter individually based on the historical gradient information.

**_Update Rule:_**
- $\theta = \theta - (\alpha/\text{root}(G+\epsilon)) \cdot \nabla J(\theta)$

**_Accumulated Squared Gradients G:_**
- Keeps track of the sum of the squares of the gradients.

# RMSprop

**_Idea:_**
- Addresses the diminishing learning rates problem of Adagrad by using a moving average of squared gradients.

**_Update Rule:_**
- $\theta = \theta - (\alpha/\text{root}(E[G]+\epsilon)) \cdot \nabla J(\theta)$

**_Exponential Moving Average E[G]:_**
  - A moving average of squared gradients.

# Adam

**_Combines Momentum and RMSprop:_**
- _Utilizes both momentum-based updates and adaptive learning rates._

**_Update Rule:_**
- $m = \beta_1 \cdot m + (1 - \beta_1) \cdot \nabla J(\theta)$
- $v = \beta_2 \cdot v + (1 - \beta_2) \cdot (\nabla J(\theta))2$
- $\theta = \theta - (\alpha / root(v + \epsilon)) \cdot m$

**_Hyperparameters $\beta_1$ and $\beta_2$:_**
- _Control the decay rates of the moment estimates._

# AdaDelta

**_Idea:_**
- _Addresses the accumulation of squared gradients issue in RMSprop by using a more sophisticated update rule._

**Update Rule:**
- $\theta = \theta - (root(E[G] + \epsilon) / root(E[\Delta\theta] + \epsilon)) \cdot \nabla J(\theta)$

**_Exponential Moving Average of Squared Parameter Updates $E[\Delta\theta]$:_**
- _A moving average of squared parameter updates._

# Nadam

**_Combines Nesterov Accelerated Gradient (NAG) with Adam:_**

- *Incorporates NAG's momentum term into Adam.*

**_Update Rule:_**

- *Similar to Adam, but with Nesterov momentum.*

# Adapting Learning Rate Schedules

**_Learning Rate Schedulers:_**

- *Techniques that dynamically adjust the learning rate during training.*

**_Examples:_**

- *Step Decay: Reduce the learning rate after a fixed number of epochs.*
- *Exponential Decay: Reduce the learning rate exponentially over time.*