# 1. Gradient descent

Basic Concept for single variable function:

- Function : $f(x)$ that we want to minimize

- Gradient : $\nabla f(x)$ = derivative/slope at point x

- Update Rule : $x_{new} = x_{old} - \alpha.\nabla f(x)$

- Learning Rate : $\alpha$ (step size parameter)

Extending to multiple variable function

- Function : $f(x_1, x_2, \cdots, x_n)$

- Gradient : $\nabla f = [\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \cdots, \frac{\partial f}{\partial x_n}]$

---

**Algorithm 1** Gradient Descent Algorithm

---

**Require:** Training data, cost function $Q(w)$
**Ensure:** Optimized parameters (weights) $w^*$
1: Initialize parameters $w$ (randomly or zeros)
2: Set learning rate $\alpha$
3: **repeat**
4:     Calculate gradient: $\nabla Q(w)$
5:     Update parameters: $w \leftarrow w - \alpha.\nabla Q(w)$
6:     Check convergence criteria
7: **until** convergence
8: **return** optimized parameters $w$

---

Traditional gradient descent (batch gradient descent) uses average gradient computed over the entire dataset for each update. This is a very simplified version (and not the entire picture) of what happens in gradient descent, but good enough to understand.

- Pros : Stable convergence, guaranteed to converge for convex functions

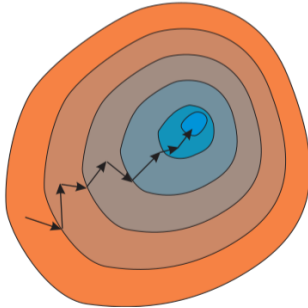- Cons : Slow for large datasets, memory intensive

## 2. STOCHASTIC GRADIENT DESCENT

Minimization of average loss over the training data

$$Q(w) = \frac{1}{l} \sum_{i=1}^{l} \mathcal{L}_i(w) \to \min_{w}$$

**Input**: dataset $X^l$, learning rate $\mu$, parameter $\lambda$
**Output**: weights $w = \left( w_{jh}, w_{hm} \right)$

**Initialization**
Set all the weights $w$ to small random numbers
Evaluate the objective function $Q(w)$

**do**

     select $x_i$ from $X^l$
     compute the loss function $\mathcal{L}_i(w)$
     gradient step $w := w - \mu \mathcal{L}_i'(w)$
     update the objective function $Q := \lambda \mathcal{L}_i + (1 - \lambda)Q$
**until** $Q$ and/or $w$ converges

## 3. MINI-BATCH GRADIENT DESCENT

Mini-batch gradient descent uses average of gradients computed over a mini-batch (small set) of dataset for each update.

- Pros : Balance between batch-GD and SGD

- Cons : Need to tune optimal batch size

## 4. GRADIENT DESCENT OPTIMIZERS

In all the below approaches $\nabla f(\theta)$ is the gradient w.r.t the parameter vector $\theta$.

**DEFINITION 1: Momentum.**
Key idea: As the descent direction is fixed add historical vector to accelerate convergence $v_t = v_{t-1} + \epsilon \nabla f(\theta_{t-1})$

**DEFINITION 2: Nesterov nomentum (NAG).**
Key idea: In addition to the descent direction "Look ahead" by computing gradient at anticipated future position. $v_t = v_{t-1} + \epsilon \nabla f(\theta_{t-1} - \mu v_{t-1})$

**DEFINITION 3: AdaGrad.**
Key idea: Adaptive learning rates - larger updates for infrequent parameters, smaller updates for frequent ones.

$$v_t = v_{t-1} + g^2, \text{where } g = \nabla f(\theta_{t-1})$$
$$\theta_{t+1} = \theta_t - \frac{\mu}{\sqrt{v_t} + \epsilon} g$$

- Pros

  - No manual learning rate tuning per parameter

– Great for sparse data

– Larger updates for rare features

- Cons

  – Learning rate decreases too aggressively, may stop learning

**DEFINITION 4: RMSprop.**

Key Idea : Fix AdaGrad's vanishing learning rate using exponential moving average.

$$v_t = \beta v_{t-1} + (1 - \beta)g^2, \beta \approx 1$$
$$\theta_{t+1} = \theta_t - \frac{\mu}{\sqrt{v_t} + \epsilon}g$$

- Pros

  – Maintains adaptive learning rates without vanishing

  – Works well for non-stationary objectives

- Cons

  – Sensitive to choice of $\mu$

**DEFINITION 5: Adam (Adaptive Moment Estimation).**

Key Idea : Combine momentum + RMSprop with bias correction.

$$v_t = \beta_1 v_{t-1} + (1 - \beta_1)g; \hat{v}_t = \frac{v_t}{1 - \beta_1^t}$$
$$s_t = \beta_2 s_{t-1} + (1 - \beta_2)g^2; \hat{s}_t = \frac{s_t}{1 - \beta_2^t}$$
$$\theta_{t+1} = \theta_t - \frac{\mu\hat{v}_t}{\sqrt{\hat{s}_t} + \epsilon}g$$

## 5. PRACTICAL ISSUES IN NEURAL NETWORKS TRAINING

- Overfitting

  – Regularization

  – Dropout

- Vanishing or exploding gradients

  – Adaptive learning rate

  – Batch normalization

- Local optima

  – Pre-training