~ **Gradient Descent**

- It is an OPTIMIZATION ALGORITHM
- Optimization = Search, FINDING BEST SOLUTION FROM ALL THE FEASIBLE ONES!
- Minimizing or Maximizing AN OBJECTIVE FUNCTION WRT SOME VARS.

○ **Analytical methods** → FINDING A CLOSED-FORM SOLUTION!
  - GIVES EXACT SOLUTIONS
  - time - CONSUMING
  - CLOSED-FORM SOLUTIONS DON'T ALWAYS EXIST!

○ **Numerical Methods:**
  - APPROX. SOLUTIONS WITH ALLOWABLE TOLERANCE
  - MIGHT BE FASTER!
  - POSSIBLE IN MOST CASES
  - BETTER SUITED WHERE THERE ARE A LOT OF FEATURES.

- It is a GENERAL OPT. ALGO. to FIND minimum.

**Resume:**
- OPTIMIZATION METHOD
- GENERAL PURPOSE
- ITERATIVE
- USES FIRST DERIVATIVES

**Gradient Vector:** VECTOR OF PARTIAL DERIVATIVES OF A FUNCTION WRT ITS PARAMS:

GRADIENT $\nabla f = \begin{bmatrix} \dfrac{\partial f}{\partial \theta_1} \\ ----- \\ \dfrac{\partial f}{\partial \theta_m} \end{bmatrix}$

IT POINTS IN THE DIRECTION OF GREATEST ASCENT OF $f$.

IN ORDER TO DESCENT THE $\nabla$ WE USE NEGATIVE OF $\nabla f$

**3 STEPS:**
- SET A LEARNING RATE: $\alpha$
- RANDOMLY INITIALIZE PARAMETERS
- REPEAT UNTIL CONVERGENCE

SUBSTITUTE WITH VALUE $\theta^t$

$$\theta^{(t+2)} = \theta^{(t)} - \alpha \dfrac{\partial f}{\partial \theta}\Big|_{\theta^t}$$

LEARNING RATE

PARTIAL DERIVATIVE WRT to $\theta$ (VAR i.e. x)

**NOTE!** THE - IS DUE TO THE FACT THAT WE WANT OPPOSITE OF GRADIENT!!!

**Example!**

$f(\theta) = (\theta - 5)^2$

$\dfrac{\partial f}{\partial \theta} = 2(\theta - 5)$

$\theta^{(0)} = 1$ RANDOMLY CHOSEN INITIAL VAL

$\theta^{(1)} = 1 - 0.1 \cdot 2(1 - 5) = 1 + 0.8 = 1.8$

$\alpha = 0.1$ CHOSEN BY US!

$\theta^{(2)} = 1.8 - 0.1 \cdot 2(1.8 - 5) = 1.8 + 0.64 = 2.44$

⋮

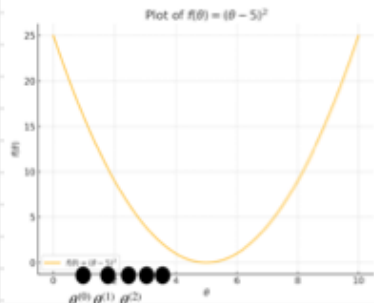WHEN to **STOP**? WHEN WE RICH Convergence!!

(NON MANDATORY GRADIENT = 0)
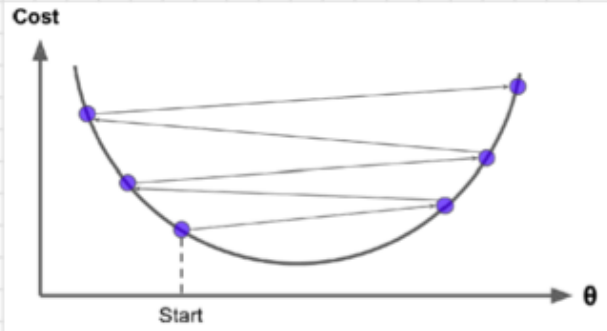
Plot of $f(\theta) = (\theta - 5)^2$

- GD IS SENSITIVE TO THE LEARNING RATE!!

Cost

Start
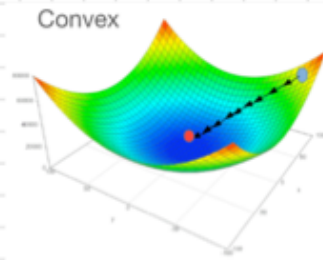
Too small!

Cost

Start

Too large!

- THE GD FINDS THE GLOBAL OPTIMUM IF THE FUNCTION IS:

  - DIFFERENTIABLE

  - CONVEX

Convex

## MULTIVARIATED FUNCTION

* IN THE EXAMPLE ABOVE WE SAW AN APPLICATION ON JUST ONE VARIABLE $\theta$

* THE GRADIENT $= \nabla$ OF A MULTIVARIATE $f$ IS:

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial \theta_1} \\ ------ \\ \frac{\partial f}{\partial \theta_m} \end{bmatrix} \longrightarrow \text{VECTOR OF ITS PARTIAL DERIVATIVES}$$

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \nabla f \big|_{\theta^t}$$

## TRAIN MODELS

- DIFFERENTIABLE
- CONVEX

LINEAR REG. EXAMPLE

* USE GD TO TRAIN ML MODELS.

* FIND OPTIMAL PARAMETERS OF THE MODEL.

$f(\theta) = \text{COST/LOSS FUNCTION}.$

$\theta = \text{VECTOR OF MODEL'S PARAMETERS}.$

$\hat{y} = \beta_0 + \beta_1 x$

$\beta = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}$

$$J(\beta) = \frac{1}{n} \sum_{i=1}^{n} (\hat{y}_i - y_i)^2 \longrightarrow \text{COST FUNCTION}$$

$$\min_{\beta} [J(\beta)] \rightarrow \text{OPTIMIZATION OBJECTIVE!}$$

EXAMPLE FOR LINEAR REGRESSION:

$\alpha = 0.01 \quad \beta_0^{(0)} = 0, \ \beta_1^{(0)} = 0$

$$J(\beta) = \frac{1}{n} \sum_{i=1}^{n} (\hat{y}_i - y_i)^2 = \frac{1}{n} \sum_{i=1}^{n} [(\beta_0 + \beta_1 x_i) - y_i]^2$$

$$\frac{\partial J}{\partial \beta_0} = \frac{1}{n} \sum_{i=1}^{n} 2(\beta_0 + \beta_1 x_i - y_i) + 0 = \frac{2}{n} \sum_{i=1}^{n} (\hat{y}_i - y_i)$$

$$\frac{\partial J}{\partial \beta_1} = \frac{1}{n} \sum_{i=1}^{n} 2(\beta_0 + \beta_1 x_i - y_i) + x_i = \frac{2}{n} \sum_{i=1}^{n} x_i (\hat{y}_i - y_i)$$

- Repeat until convergence (e.g., J is not changing much):
  ○ Update $\beta$:

$$\beta_0^{(t+1)} = \beta_0^{(t)} - \alpha \frac{2}{n} \sum_{i=1}^{n} (\hat{y}_i^{(t)} - y_i)$$

$$\beta_1^{(t+1)} = \beta_1^{(t)} - \alpha \frac{2}{n} \sum_{i=1}^{n} x_i (\hat{y}_i^{(t)} - y_i)$$

# Variations of Gradient Descent

## Batch

- **All training data is used to compute gradient:**
  - Gradient of **cost function**
  - Used **in every step**
- **Very slow** in training
- Can cause memory problem
- Directly move towards optimum
- Works well **only for convex**

$$\beta^{(t+1)} = \beta^{(t)} - \alpha \frac{1}{n} \sum_{i=1}^{n} \nabla L_{BGD}(\hat{y}_i^{(t)}, y_i)$$

↳ n = # of all points in DF.

## Stochastic

- Considers **just one** datapoint **at a time**
- **Repeat** until **convergence:**
  - **Repeat** for a given **# of iter:**
    - **Take a random sample**
    - **Feed it** to formula to update params.
- By convention repeat n = size of DF
- Much **easier**
- **Escapes local minima**
- May never reach min

$$\beta^{(t+1)} = \beta^{(t)} - \alpha \left( \nabla L_{BGD}(\hat{y}_i^{(t)}, y_i) \right)$$

only 1 point

## Mini-Batch

- Computes gradients on small random sets
- May escape local minima
- **Less fluctuation**

$$\beta^{(t+1)} = \beta^{(t)} - \alpha \frac{1}{n} \sum_{i \in S} \nabla L_{BGD}(\hat{y}_i^{(t)}, y_i)$$

S is a random set of samples