

Gradient Descent

$$\theta_0 \leftarrow \text{Random}$$

$$\theta_k \leftarrow \theta_{k-1} - \alpha \underbrace{\nabla_{\theta} J(\theta_{k-1})}_{\text{gradient}} \quad (*)$$

GD

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n \text{cost}(y_i, \hat{y}_i)$$

$$\nabla_{\theta} J(\theta) = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} \text{cost}(y_i, \hat{y}_i)$$

$$\underbrace{\nabla_{\theta} J(\theta)}_{\text{approximation}} = \frac{1}{s} \sum_{i=1}^s \nabla_{\theta} \text{cost}(\underline{y_i}, \underline{\hat{y}_i})$$

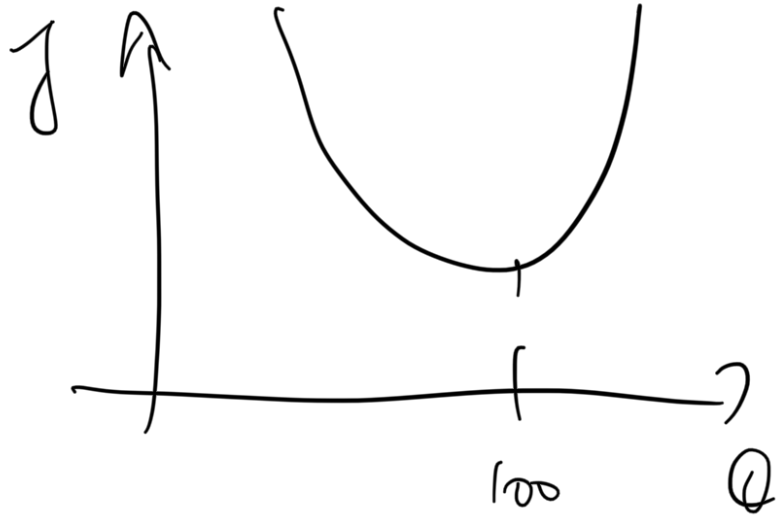
$s \leq n$

SGD

Approximation  
in a faster time (seconds)

Steps are faster  
↳ may need more steps.

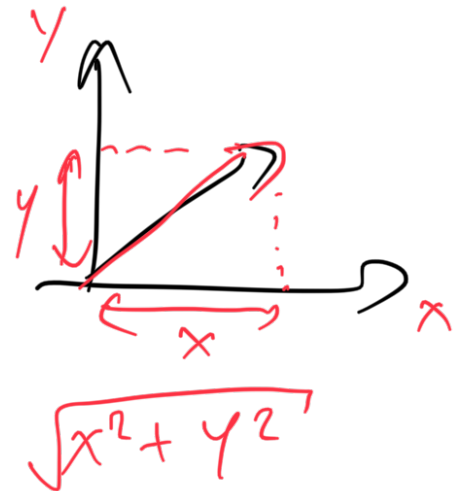
but w- 1



stopping criterion

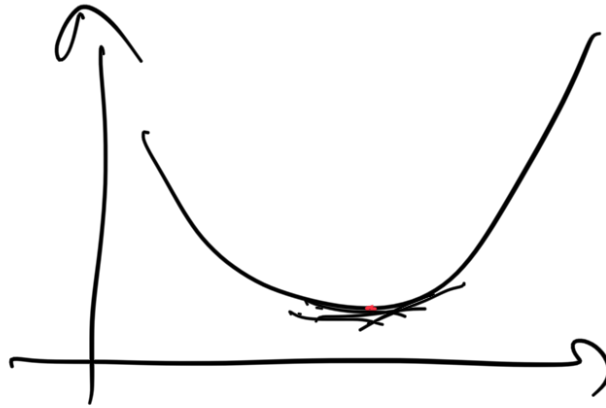
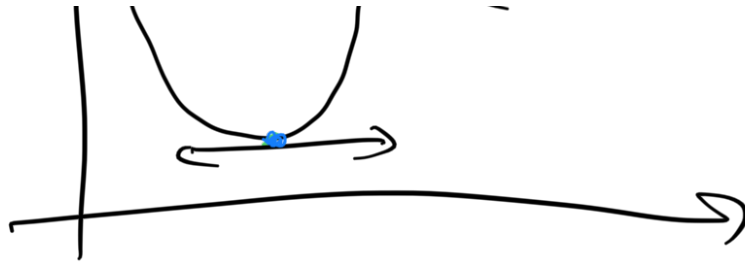
$$\|\nabla_Q J\|_2 \leq \varepsilon$$

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$

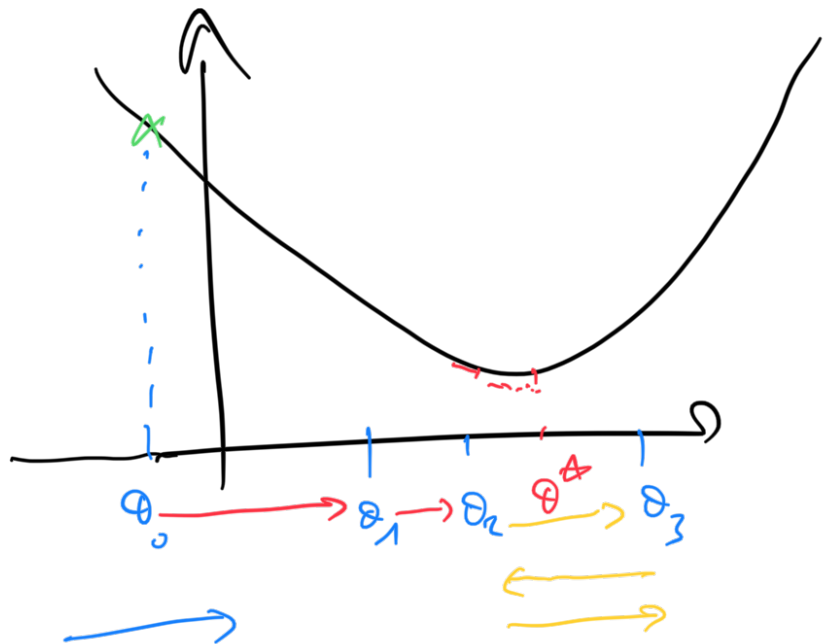


$$\|x\|_2 = \sqrt{x_1^2 + x_2^2 + \dots + x_d^2}$$





$$\|\nabla J\|_2 \leq \varepsilon$$

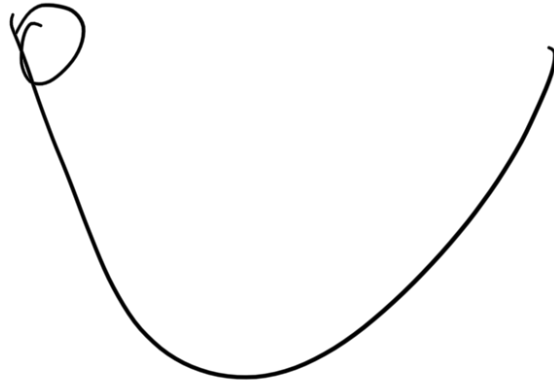


Decrease learning rate over the step

10 steps

eg. multiply by 0.9 every time

---



$$\left\{ \begin{array}{l} V = \mu V - \alpha \nabla_{\theta} J(\theta_{k-1}) \\ \theta_k \leftarrow \theta_{k-1} + V \end{array} \right. \quad \text{with momentum}$$