

Large Scale Machine Learning

Learning with Large Data Sets

$m = 100,000,000$

Why not $m = 1,000$? Sanity check!

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Sanity check $m = 1000$

high variance learning curves



Stochastic Gradient Descent

linear regression w/ grad descent

$$h_{\theta}(x) = \sum_{j=0}^n \theta_j x_j$$

$$J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Repeat {

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

(for every $j = 0, \dots, n$) $\underbrace{\frac{\partial}{\partial \theta_j} J_{\text{train}}(\theta)}$

}

$m = 300,000,000$
Batch gradient descent SLOW

Stochastic GD

$$\text{cost}(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$J_{\text{train}}(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(\theta, (x^{(i)}, y^{(i)}))$$

1. Randomly shuffle dataset

2. Repeat { 1-10 times

for $i = 1, \dots, m$ {

$$\theta_j := \theta_j - \alpha (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$$\left. \begin{array}{l} \text{(for } j=0, \dots, n) \quad \frac{\partial}{\partial \theta_j} \text{Cost}(\theta, (x^{(i)}, y^{(i)})) \\ \} \end{array} \right\}$$

$$(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)})$$

1st \rightarrow 2nd \dots

Mini-Batch Gradient Descent (sometimes faster than stochastic)

Batch: Use all m examples in each iteration

Stochastic gradient descent: Use 1 example "

Minibatch: Use b examples "

$$b = \text{mini-batch size} \quad b=10$$

$$\text{Get } b=10 \text{ examples } (x^{(1)}, y^{(1)}), \dots, (x^{(b)}, y^{(b)})$$

$$\text{for } i=1, 2, \dots \left\{ \theta_j := \theta_j - \alpha \frac{1}{b} \sum_{k=1}^b (h_{\theta}(x^{(k)}) - y^{(k)}) x_j \right.$$

(for every $j=0, \dots, n$)

}

b examples vs. 1 example

Vectorization

Disadvantage new parameter b

Stochastic Gradient Descent Convergence

$$\text{Cost}(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

During learning, compute $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$ before updating θ using $(x^{(i)}, y^{(i)})$

Every 1000 iterations (say), plot $\text{cost}(\Theta, (x^{(i)}, y^{(i)}))$ averaged over the last 1,000 examples processed by algo
 ↑ 5,000 for flatter lines / if too noisy

Learning rate α is typically held constant. (AND decrease α over time if we want Θ to converge
 (E.g. $\alpha = \frac{\text{const1}}{\text{iteration} + \text{const2}}$)

Online Learning (streaming)

Features x capture properties of user, of origin/destination and asking price. We want to learn $p(y=1|x; \Theta)$

Repeat forever {

Get (x, y) corresponding to user.

Update Θ using (x, y) :

$$\Theta_j := \Theta_j - \alpha (h_{\Theta}(x) - y) x_j \quad (j = 0, \dots, n)$$

} \rightarrow Can adapt to changing user preference.

Predicted Click-through-rate (CTR)

Map Reduce & Data Parallelism

M/R $m=400$ $m=400,000,000$

Batch : $\Theta_j := \Theta_j - \alpha \frac{1}{n} \sum_{i=1}^n (h_{\Theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$

Machine 1: $U_x(x^{(1)}, y^{(1)}), \dots, (x^{(100)}, y^{(100)})$

$$\text{temp}_j^{(1)} = \sum_{i=1}^{100} (h_{\Theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \longrightarrow$$

Machine 2: $U_x(x^{(101)}, y^{(101)}), \dots, (x^{(200)}, y^{(200)})$

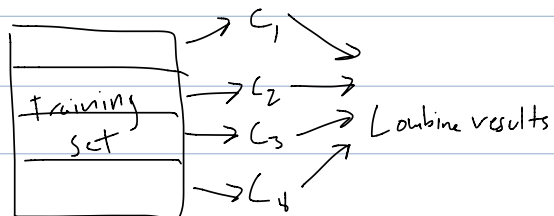
$$\text{temp}_j^{(2)} = \sum_{i=1}^{100} (h_{\Theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \longrightarrow$$

4 Machines



Master Server Combine

$$\Theta_j := \Theta_j - \alpha \frac{1}{400} (t_{\text{tmp},j}^{(1)} + t_{\text{tmp},j}^{(2)} + t_{\text{tmp},j}^{(3)} + t_{\text{tmp},j}^{(4)}) \quad (j=0, \dots, n)$$



Multi-core Machines