

Learning with large datasets

Sheer amount of data \rightarrow Learning algorithm works better

Large datasets \rightarrow low-bias learning algorithm

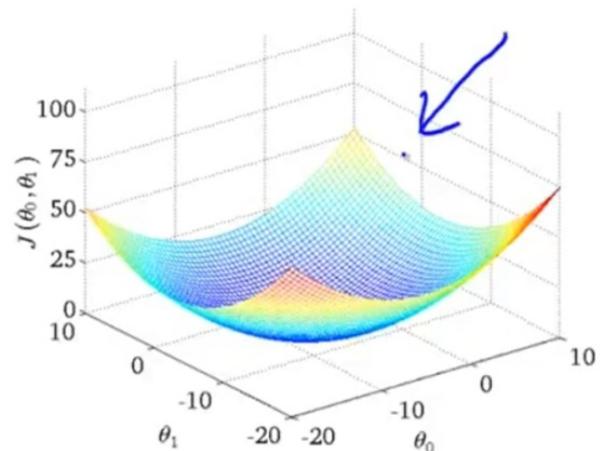
Own unique problem - Computational prob.

Stochastic Gradient Descent

Linear regression with gradient descent

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

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



Repeat {

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

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

}

if m is large, computing \sum can be very expensive.

Batch gradient descent

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

Repeat {

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

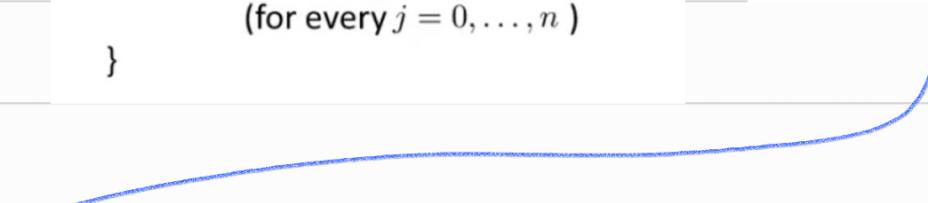
$\underbrace{\frac{\partial}{\partial \theta_j} J_{train}(\theta)}_{\text{for every } j = 0, \dots, n}$

}

Stochastic gradient descent

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

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



→ 1. Randomly shuffle dataset

2. Repeat {

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

$$\theta_j := \theta_j - \alpha \underbrace{(h_\theta(x^{(i)}) - y^{(i)})}_{\frac{\partial}{\partial \theta_j} cost(\theta, (x^{(i)}, y^{(i)}))} \cdot x_j^{(i)}$$

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

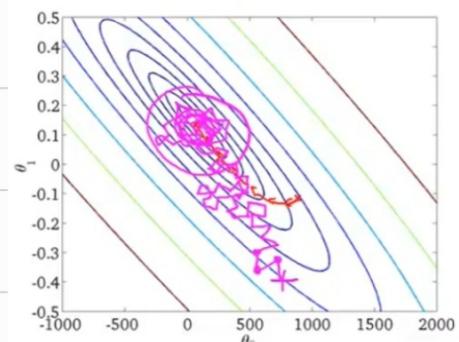
}

}

Batch gradient descent will tend to take a reasonably straight line trajectory to get to the global minimum



Stochastic gradient descent: every iteration is just trying to fit single training ex. better



Mini-batch gradient descent

Batch GD : Use all m examples in each iteration

SGD : Use 1 example in each iteration

Mini-batch GD : Use b examples in each iteration

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

Say $b = 10, m = 1000.$

Repeat { ↗

→ for $i = 1, 11, 21, 31, \dots, 991\{$

→ $\theta_j := \theta_j - \alpha \frac{1}{10} \sum_{k=i}^{i+9} (h_\theta(x^{(k)}) - y^{(k)}) x_j^{(k)}$

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

}

}

첫 10 examples로 cost를 α 만큼 고정,

다음 10 examples로 고정 ... m 개 솔루션까지 반복

good vectorized implementation : better than SGD

Stochastic Gradient Descent Convergence

Stochastic gradient descent:

→ $cost(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2}(h_\theta(x^{(i)}) - y^{(i)})^2 \Rightarrow (x^{(i)}, y^{(i)}), (x^{(i+1)}, y^{(i+1)}), \dots$

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

→ Every 1000 iterations (say), plot $cost(\theta, (x^{(i)}, y^{(i)}))$ averaged over the last 1000 examples processed by algorithm.

