

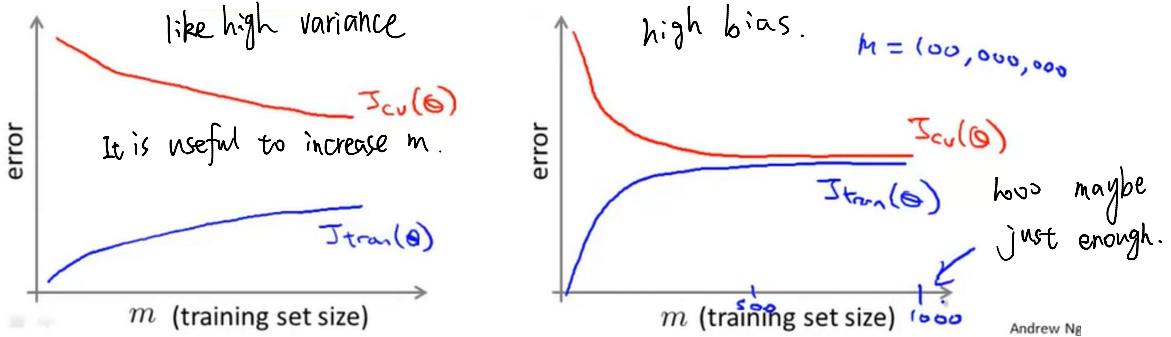
# Week 10

## • Gradient Descent with Large Datasets:

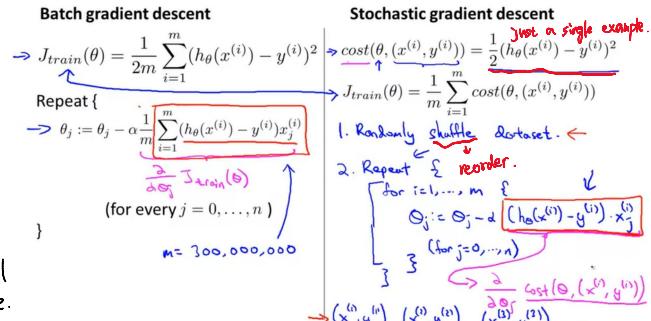
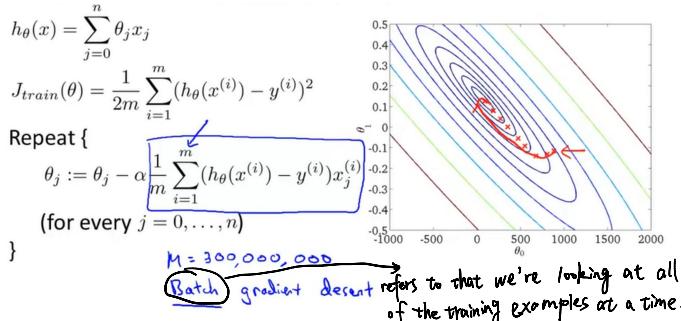
### a. Learning With Large Datasets:

Before we use massive data sets, we can check if one thousand is ok.

The way to check of using a much smaller training set: plotting the learning curve.



### b. Stochastic Gradient Descent:



### c. Mini-Batch Gradient Descent:

→ Batch gradient descent: Use all  $m$  examples in each iteration

→ Stochastic gradient descent: Use 1 example in each iteration

Mini-batch gradient descent: Use  $b$  examples in each iteration

$$b = \text{mini-batch size. } b = 10. \quad \frac{2-100}{2}$$

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

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

$$i := i + 10$$

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  $j = 0, \dots, n$ )

}

}

### d. Stochastic Gradient Descent Convergence:

How do you make sure that the algorithm is completely debugged and is converging okay? And how to tune the learning rate.

#### Checking for convergence

→ Batch gradient descent:

→ Plot  $J_{train}(\theta)$  as a function of the number of iterations of gradient descent.

$$\rightarrow J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \quad M = 300,000,000$$

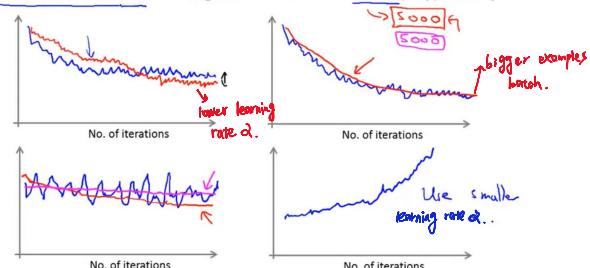
→ Stochastic gradient descent:

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

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

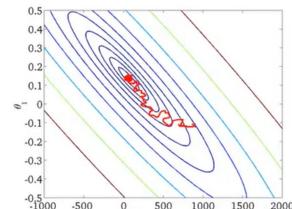
Plot  $cost(\theta, (x^{(i)}, y^{(i)}))$ , averaged over the last 1000 (say) examples



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

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

1. Randomly shuffle dataset.  
 2. Repeat {  
 for  $i := 1, \dots, m$  {  
 $\theta_j := \theta_j - \alpha(h_\theta(x^{(i)}) - y^{(i)})x_j^{(i)}$   
 (for  $j = 0, \dots, n$ )  
 }  
 }



Learning rate  $\alpha$  is typically held constant. Can slowly decrease  $\alpha$  over time if we want  $\theta$  to converge. (E.g.  $\alpha = \frac{\text{const1}}{\text{iterationNumber} + \text{const2}}$ )  $\rightarrow 0$

## • Advanced Topics:

### a. Online Learning:

#### Online learning

Shipping service website where user comes, specifies origin and destination, you offer to ship their package for some asking price, and users sometimes choose to use your shipping service ( $y = 1$ ), sometimes not ( $y = 0$ ).

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

Report forever {  
 Get  $(x, y)$  corresponding to user.  
 Update  $\theta$  using  $(x, y)$ .  
 $\theta_j := \theta_j - \alpha(h_\theta(x) - y) \cdot x_j$   
} . Can adapt to changing user preference.

#### Other online learning example:

Product search (learning to search)

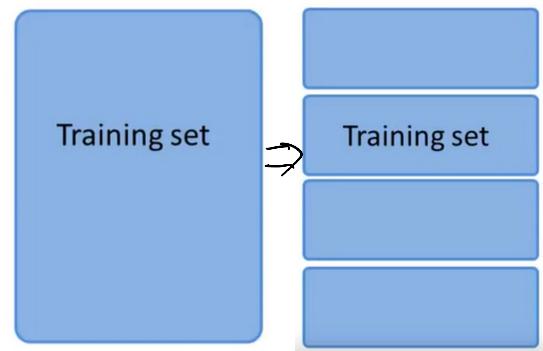
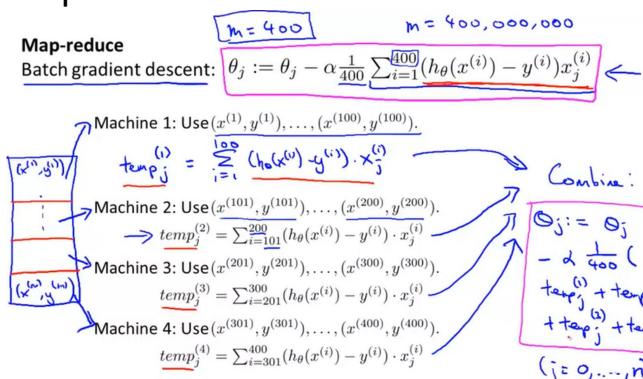
User searches for "Android phone 1080p camera"  $\leftarrow$   
 Have 100 phones in store. Will return 10 results.

$x$  = features of phone, how many words in user query match name of phone, how many words in query match description of phone, etc.

$y = 1$  if user clicks on link.  $y = 0$  otherwise.

$\rightarrow$  Learn  $p(y=1|x; \theta)$ .  $\leftarrow$  predicted CTR

### b. Map Reduce and Data Parallelism:



Many learning algorithms can be expressed as computing sums of functions over the training set.

E.g. for advanced optimization, with logistic regression, need:

$$\begin{aligned} \Rightarrow J_{\text{train}}(\theta) &= -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)}) - (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)})) \\ \Rightarrow \frac{\partial}{\partial \theta_j} J_{\text{train}}(\theta) &= \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} \end{aligned}$$

$\text{temp}_j^{(1)} \quad \text{temp}_j^{(2)} \leftarrow$

