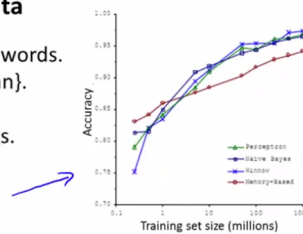


- Gradient Descent with Large Datasets
- Learning With Large Datasets
- Machine Learning and data

Machine learning and data

Classify between confusable words.
E.g., {to, two, too}, {then, than}.

For breakfast I ate two eggs.



→ "It's not who has the best algorithm that wins.
It's who has the most data."

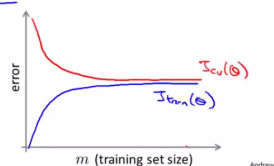
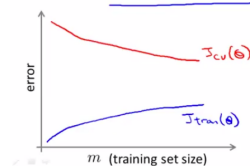
- Learning with large datasets

Learning with large datasets

$m = 100,000,000$ ←

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

$m = 1,000$ ←



- Stochastic gradient descent
- Linear regression with gradient descent

Linear regression with gradient descent

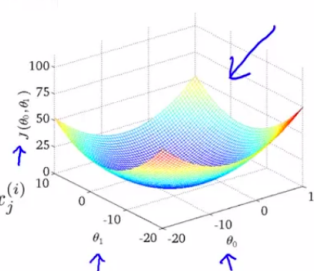
$$h_{\theta}(x) = \sum_{j=0}^n \theta_j x_j$$

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

Repeat {

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

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



- If m is large computing the derivative term can be very expensive

Linear regression with gradient descent

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

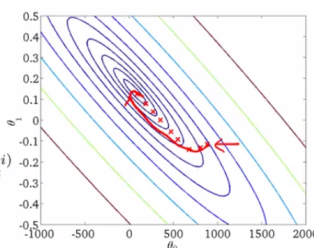
$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Repeat {

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

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

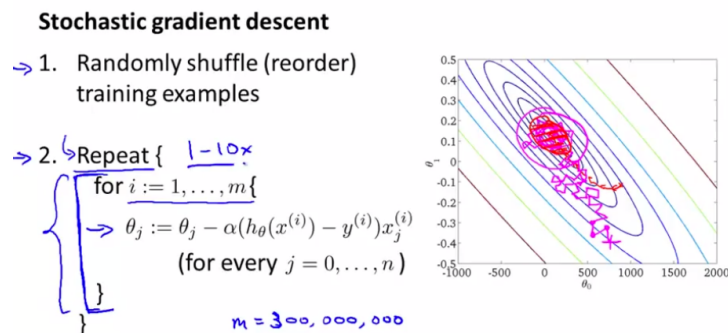
}



- Batch gradient this particular version
 - means that we are looking at all the training examples every time
- Batch gradient descent and Stochastic gradient descent

Batch gradient descent	Stochastic gradient descent
$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$	$cost(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$
$J_{train}(\theta) = \frac{1}{m} \sum_{i=1}^m cost(\theta, (x^{(i)}, y^{(i)}))$	$J_{train}(\theta) = \frac{1}{m} \sum_{i=1}^m cost(\theta, (x^{(i)}, y^{(i)}))$
Repeat { $\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$ $\frac{\partial}{\partial \theta_j} J_{train}(\theta)$ } (for every $j = 0, \dots, n$)	1. Randomly shuffle dataset. 2. Repeat { for $i = 1, \dots, m$ { $\theta_j := \theta_j - \alpha (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$ } $\frac{\partial}{\partial \theta_j} cost(\theta, (x^{(i)}, y^{(i)}))$ $(x^{(i)}, y^{(i)}), (x^{(i)}, y^{(i)}), (x^{(i)}, y^{(i)})$

- stochastic gradient descent
 - doesn't actually find the global minimum stays in a region around the global minimum



- 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 = mini-batch size. $b = 10$ range $b = 2$ to 100 range
 - Get $b = 10$ examples from training set

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

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

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

$i := i + 10$

- Mini-batch gradient descent

Mini-batch gradient descent

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 $j = 0, \dots, n$)

}

}

- Vectorization
 - will output perform stochastic gradient descent if you have a good vectorized implementation
- **Stochastic Gradient Descent Convergence**
- Checking for convergence
 - Batch gradient descent
 - decreasing on every iteration

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:

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

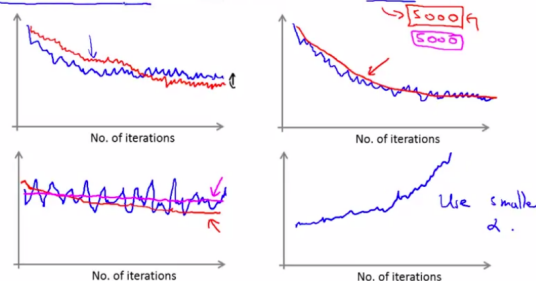
→ During learning, compute $cost(\theta, (x^{(i)}, y^{(i)}))$ before updating θ using $(x^{(i)}, y^{(i)})$.

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

- Checking for convergence

Checking for convergence

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



- Stochastic gradient descent

Stochastic gradient descent

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

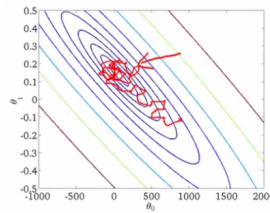
$$J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m \text{cost}(\theta, (x^{(i)}, y^{(i)}))$$

1. Randomly shuffle dataset.

2. Repeat {

$$\text{for } i := 1, \dots, m \quad \left\{ \begin{array}{l} \theta_j := \theta_j - \alpha (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \\ \text{(for } j = 0, \dots, n) \end{array} \right.$$

}

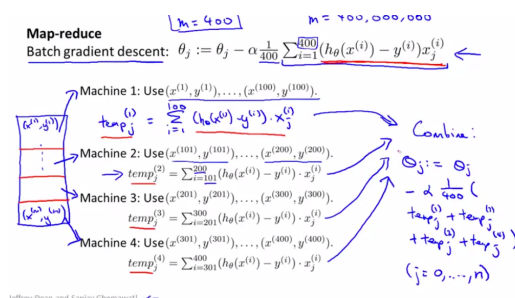


Learning rate α is typically held constant. Can slowly decrease α over time if we want θ to converge. (E.g. $\alpha = \frac{\text{const}}{\text{iterationNumber} + \text{const}}$)

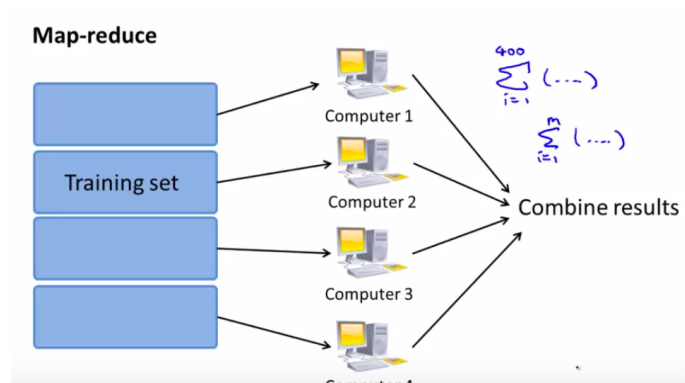
- Online Learning
 - Shipping service website where user comes, specifies origin and destination, you offer to ship their package for some asking price, and users sometimes choose to use your shipping service sometimes not. $Y = 1$ or 0
 - Features x captures properties of user, or origin/ destination and asking price. We want to learn $p(y=1|x;\theta)$ to optimize price.

Repeat forever {
 Get (x, y) corresponding to user.
 Update θ using (x, y) :
 $\theta_j := \theta_j - \alpha (h_{\theta}(x) - y) \cdot x_j \quad (j=0, \dots, n)$
 }
 price logistic regression

- can adapt to changing user preferences
- Other online learning example
 - Product search (learning to search)
 - User search for “Android phone 1080p camera”
 - 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.
 - Learn $p(y = 1 | x; \theta)$ <- predicts CTR
 - User to show user the 10 phones they’re most likely to click on
 - **get example learn from it then discard and move on**
- Map Reduce and Data Parallelism
- Map-reduce



- Map-reduce



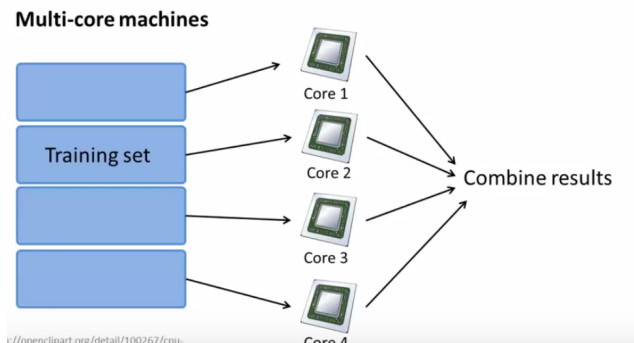
- Map-reduce and summation over the training set
 - many learning algorithms can be expressed as computing sums of functions over the training set

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

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

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

- Multi-core Machines



- don't have to worry about network latency