

# Week : 10

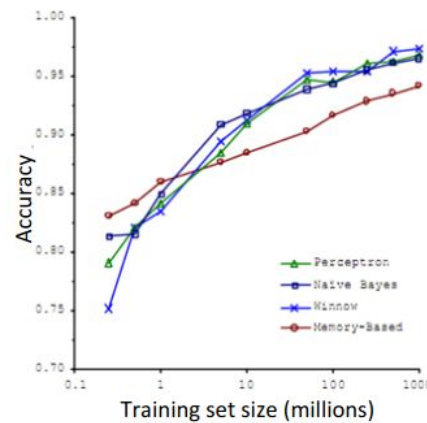
## Learning with large datasets :

If you have a large size of data then it can help you even with a bad algorithm.

### 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.”

[Figure from Banko and Brill, 2001]

Andrew Ng

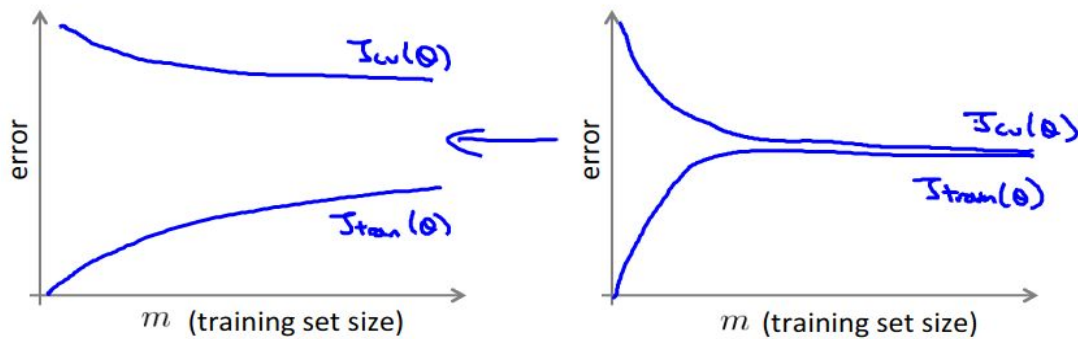
Sometimes you may have hundreds of millions of data, but if you want to use them all to train your algorithm, it will take a huge time. Instead of doing that we can decide how much bigger data we need to train our algorithm by the help of high bias - high variance graph shown in the picture below. If we ended up with a graph like the left one, that means our graph is in high variance and more training sets will help our data to improve result. But if our graph will be like the right one, then it is suffering from high bias. Adding more training sets won't help to improve our algorithm. We can try adding new features to solve high bias.

## Learning with large datasets

$m = 100,000,000$

$m = 1,000?$

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



Andrew Ng

## Stochastic gradient descent :

Original gradient descent where we count all training sets together, if we had a large dataset then it would be so much slower. As we count all training sets together, another name of this gradient descent is batch gradient descent. In order to solve the slowness we have another gradient descent, stochastic gradient descent. It counts a single dataset at a time rather than the whole dataset. As a result the algorithm is much faster than our original one.

## Linear regression with gradient descent

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

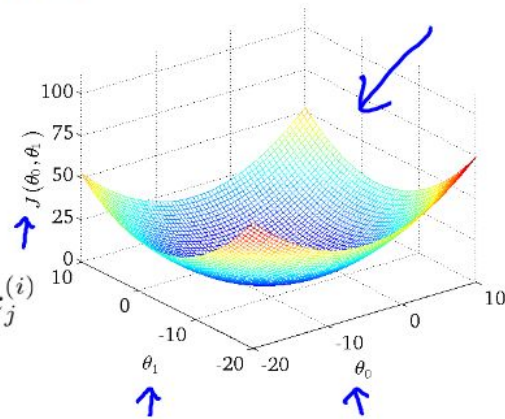
$$\rightarrow J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Repeat {

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

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

}



Andrew Ng

## Linear regression with gradient descent

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

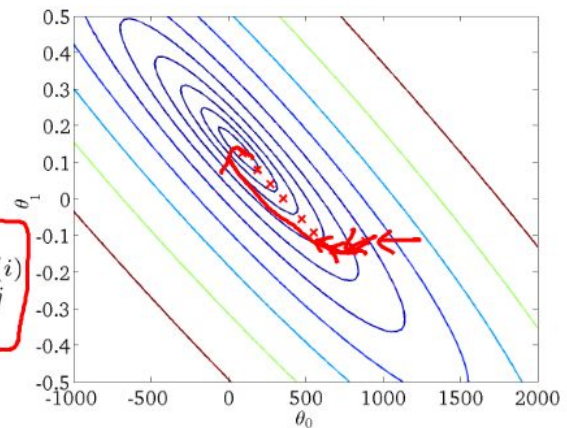
$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Repeat {

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

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

}



$M = 300,000,000$

Batch gradient descent

Andrew Ng

### Batch gradient descent

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

Repeat {

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

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

$m = 300,000,000$

}

### Stochastic gradient descent

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

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

1. Randomly shuffle dataset. ←

2. Repeat {

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

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

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

}

$\rightarrow \frac{\partial}{\partial \theta_j} cost(\theta, (x^{(i)}, y^{(i)}))$

$\rightarrow (x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}), \dots$

Andrew Ng

As we use one single data set at a time, our gradient descent will move in a zigzag mood, and it is normal. But it is faster than the convention one and also can achieve our goal.

### Stochastic gradient descent

→ 1. Randomly shuffle (reorder) training examples

→ 2. Repeat { 1-10x

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

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

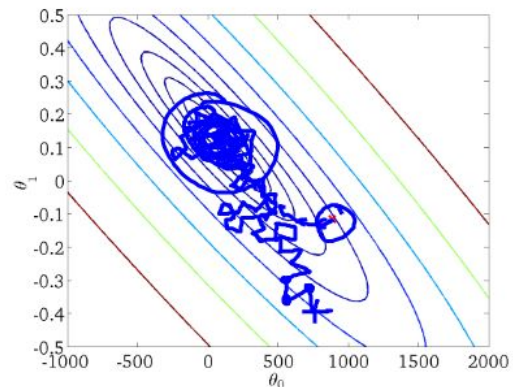
(for  $j = 0, \dots, n$ )

}

every

}

→  $m = 300,000,000$



Andrew Ng

## Mini batch gradient descent :

Besides stochastic gradient descent we have another gradient descent named mini batch gradient descent which is also faster than batch gradient descent. In this gradient descent we will take  $b$  number of test datasets at-a-time and iterate whole datasets by grouping in  $b$  number at a time.

### 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

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

Get  $b = 10$  examples  $(x^{(i)}, y^{(i)}) \dots (x^{(i+9)}, y^{(i+9)})$

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

$i := i + 10$

Andrew Ng

You may be surprised why we should have another gradient descent where stochastic is the fastest. Well, in that case if you can implement a good vectorized implementation then it will be faster than stochastic gradient descent.



## Mini-batch gradient descent

Say  $b = 10$ ,  $m = 1000$ .

Repeat {

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

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

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

}

}

$m = 300,000,000$

→  $b$  examples

→ 1 example

Vectorization

$b = 10$

Andrew Ng

## Stochastic gradient descent convergence :

It is important to check if our algorithm is learning well. To do so, we can plot a graph of convergence for every 1k iterations. The number of iterations may vary but the main motivation of this doing is to check if stochastic gradient descent is working well.

### Checking for convergence

→ Batch gradient descent:

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

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

$m = 300,000,000$

→ Stochastic gradient descent:

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

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

→  $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots$

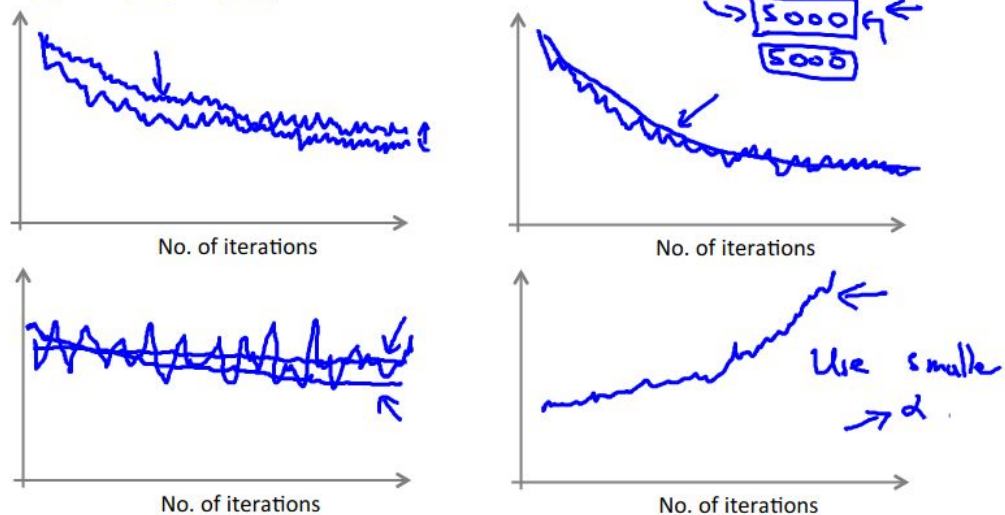
→ Every 1000 iterations (say), plot  $cost(\theta, (x^{(i)}, y^{(i)}))$  averaged over the last 1000 examples processed by algorithm.

Andrew Ng

We may end up with such graphs. If you end up with top like twos, then we can assume that our algorithm is learning well. But if our graph ends up with bottom like twos, then there may be some problem in our algorithm. It may happen because of learning rate alpha. We can set alpha smaller and hope our algorithm learns well.

### Checking for convergence

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



Andrew Ng

Setting such a small learning rate may make our algorithm slower. But if we set a large learning rate then it also misleads our algorithm. To solve both problems we can set a large learning rate at first of our algorithm and decrease it by every iteration. To do so, we have been given a formula in the picture below.

## Stochastic gradient descent

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

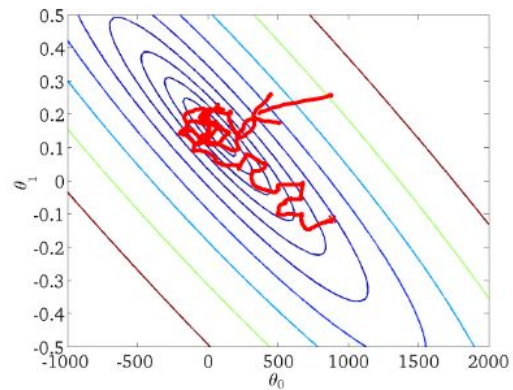
$$J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m \text{cost}(\theta, (x^{(i)}, y^{(i)}))$$

1. Randomly shuffle dataset.
2. Repeat {
 

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.  $\alpha = \frac{\text{const1}}{\text{iterationNumber} + \text{const2}}$ )

Andrew Ng

## Stochastic gradient descent

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

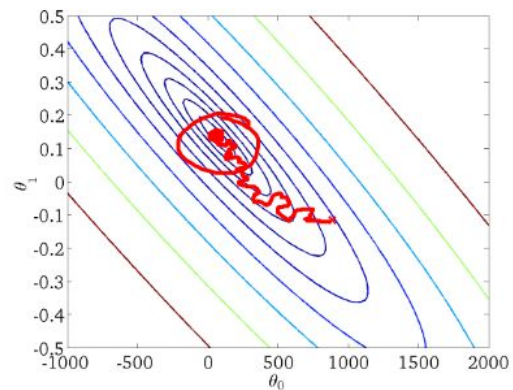
$$J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m \text{cost}(\theta, (x^{(i)}, y^{(i)}))$$

1. Randomly shuffle dataset.
2. Repeat {
 

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.  $\alpha = \frac{\text{const1}}{\text{iterationNumber} + \text{const2}}$ )  $\alpha \rightarrow 0$

Andrew Ng



## Online learning :

Sometimes websites have a large number of users. As a result, they would have a huge number of training set. If they want to train their algorithm in a conventional way they may need a large number of storage as well as a huge time for learning. On the other hand, if they have an algorithm which can learn online and per-user click they can avoid both storage and time problems. In that technique, the online algorithm takes a training set from a user and feeds it to the algorithm and throws it away.

### 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.

Repeat forever {  
  Get  $(x, y)$  corresponding to user.  
  Update  $\theta$  using  $(x, y)$ :  
     $\rightarrow \theta_j := \theta_j - \alpha (h_0(x) - y) \cdot x_j \quad (j = 0, \dots, n)$   
  }  
   $\rightarrow$  Can adapt to changing user preference.

price      logistic regression

Andrew Ng

Suppose a user goes to a mobile phone website and searches mobile phones by writing "android phone 1080p camera". And if the website wants to show 10 most click-likely mobile phones then it needs to pick 10 best mobile phones based on these keywords. And to do so, they can apply logistic regression and learn online.

### Other online learning example:

Product search (learning to search)

User searches 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.  $(x, y)$  ←

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

→ Use to show user the 10 phones they're most likely to click on.

Other examples: Choosing special offers to show user; customized selection of news articles; product recommendation; ...

Andrew Ng

## Map reduce and data parallelism :

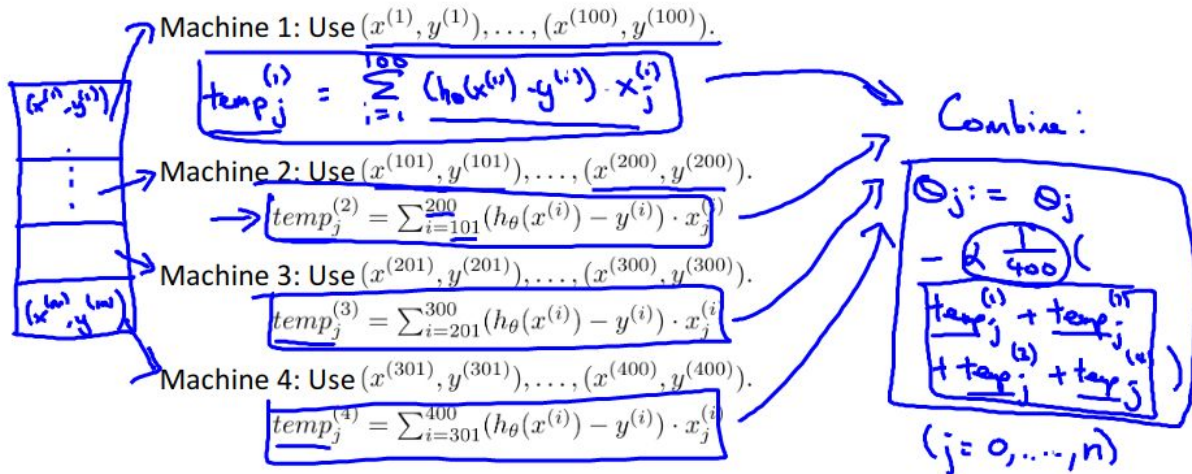
If we have a large number of data that no one machine we have can process all together, then we can use map reduce algorithm to use more than one machine to process data. In that case if we have  $p$  number of machines, we divide our data into  $p$  blocks and send each block to one machine. After computing these data all machines send back the result to the main machine, and the main machine will sum all of the work.

## Map-reduce

Batch gradient descent:

$$\theta_j := \theta_j - \alpha \frac{1}{400} \sum_{i=1}^{400} (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

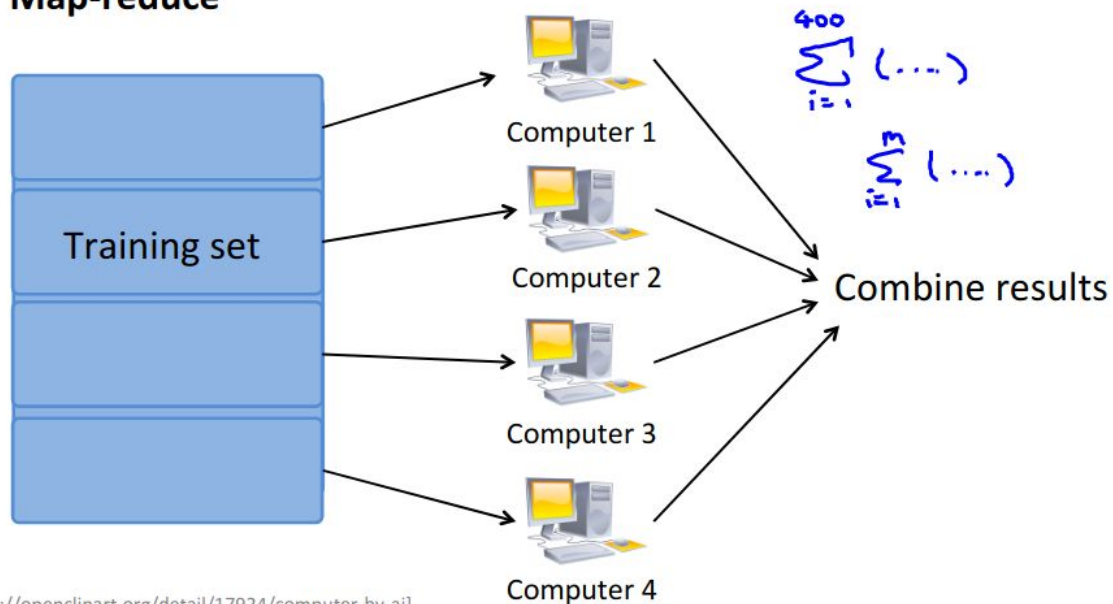
$m = 400$  ←       $m = 400,000,000$



[Jeffrey Dean and Sanjay Ghemawat]

Andrew Ng

## Map-reduce



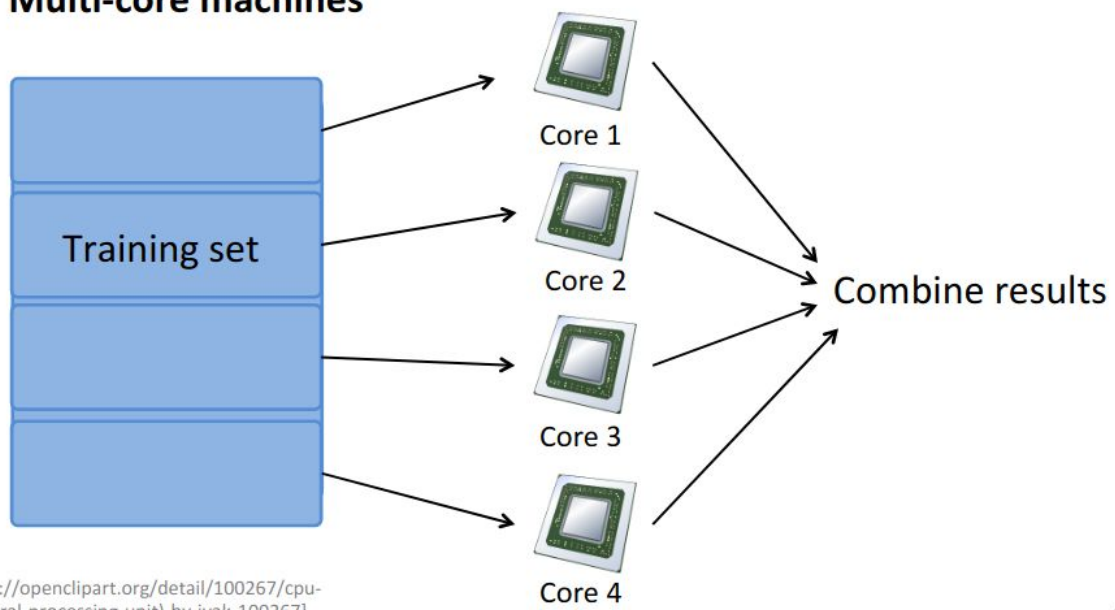
[<http://opencitart.org/detail/17924/computer-by-aj>]

Andrew Ng

We can use map reduce algorithm on all algorithms which have summation over the training set.



## Multi-core machines



Andrew Ng