



Machine Learning

Large scale  
machine  
~~Learning~~  
with large  
datasets

# Machine learning and data

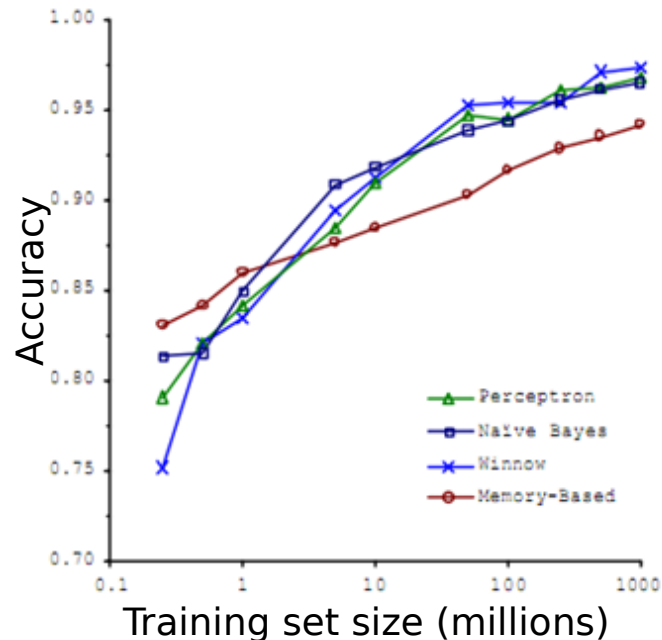
Classify between confusable words.

E.g., {to, two, too}, {then, than}.

For breakfast I ate \_\_\_\_\_  
eggs.

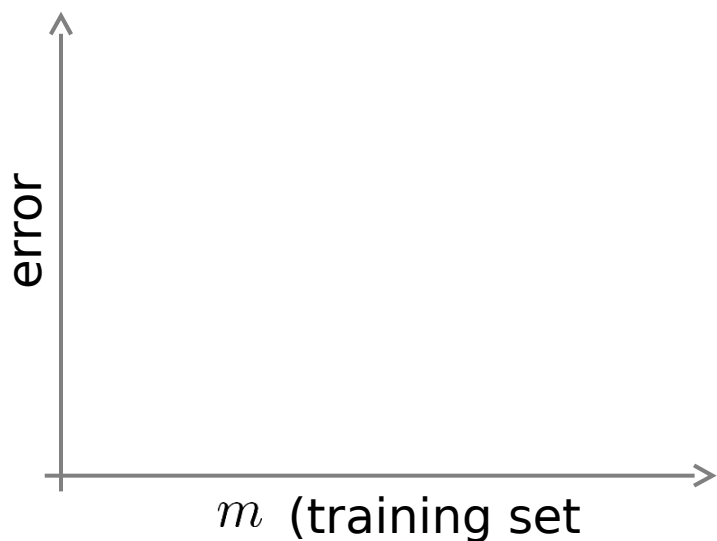
“It’s not who has the best algorithm that wins.

It’s who has the most data.”

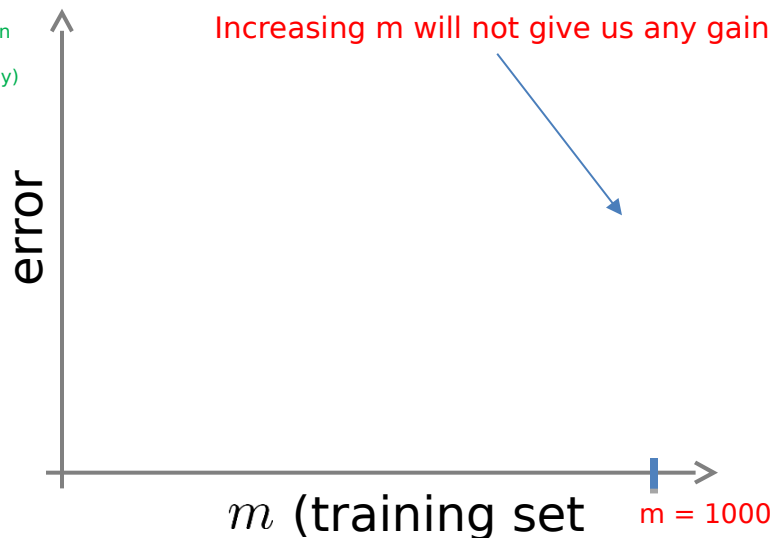


# Learning with large datasets

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



Adding extra features, hidden unit in neural networks, .. (say) and then increasing m



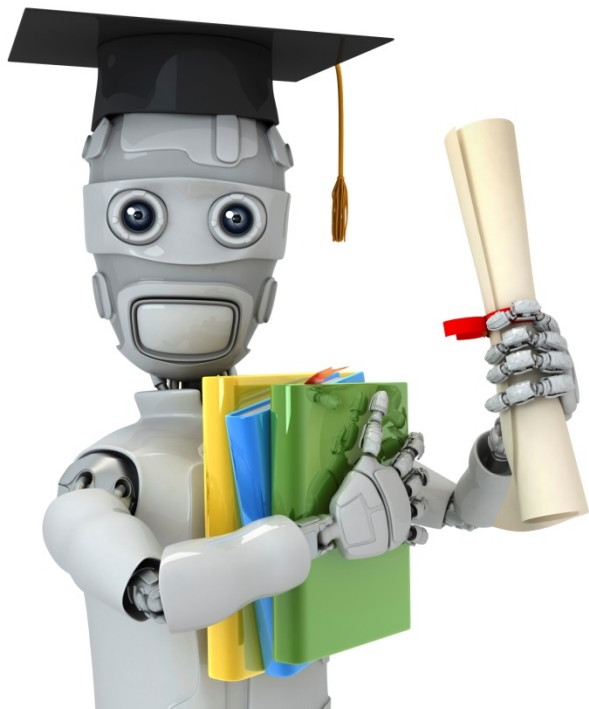
Increasing m will not give us any gain

Suppose you are facing a supervised learning problem and have a very large dataset ( $m = 100,000,000$ ). How can you tell if using all of the data is likely to perform much better than using a small subset of the data (say  $m = 1,000$ )?

---

- ☐ There is no need to verify this; using a larger dataset always gives much better performance.
- ☐ Plot  $J_{\text{train}}(\theta)$  as a function of the number of iterations of the optimization algorithm (such as gradient descent).
- ☐ Plot a learning curve ( $J_{\text{train}}(\theta)$  and  $J_{\text{CV}}(\theta)$ , plotted as a function of  $m$ ) for some range of values of  $m$  (say up to  $m = 1,000$ ) and verify that the algorithm has bias when  $m$  is small.
- ☒ Plot a learning curve for a range of values of  $m$  and verify that the algorithm has high variance when  $m$  is small.

**Correct Response**



Machine Learning

Large scale  
machine  
~~learning~~  
gradient  
descent

# Linear regression with gradient descent

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

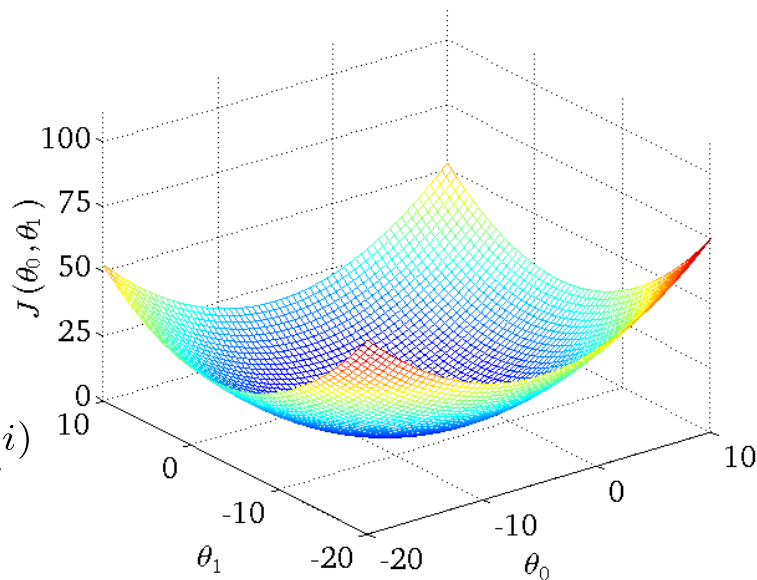
$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Repeat {

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

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

}



# Linear regression with gradient descent

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

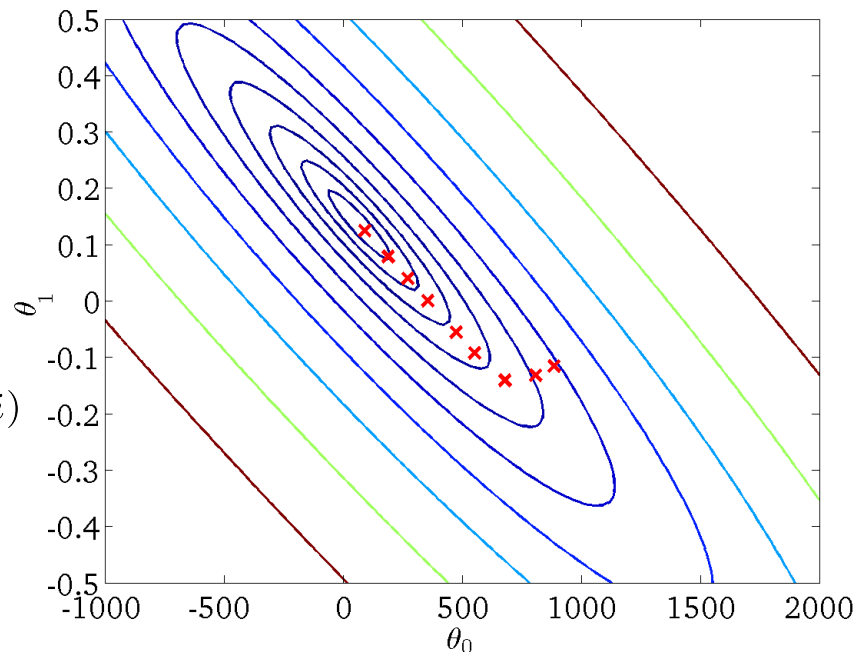
$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Repeat {

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

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

}



When  $m$  is large then computing the derivative becomes very difficult

## Batch gradient descent

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

Repeat {

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

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

}

Step 1: we modify parameters and try to fit 1<sup>st</sup> training example.  
Step 1: we modify parameters and try to fit 2<sup>nd</sup> training example.  
... on

## Stochastic gradient descent

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

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

(Pre-processing step)



# Stochastic gradient descent

1. Randomly shuffle  
(reorder) training  
examples

How much to repeat the for loop?  
Generally depends on the size of the training set but normally 1-10 repetition works fine

2. Repeat  $\{ \dots, m$

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

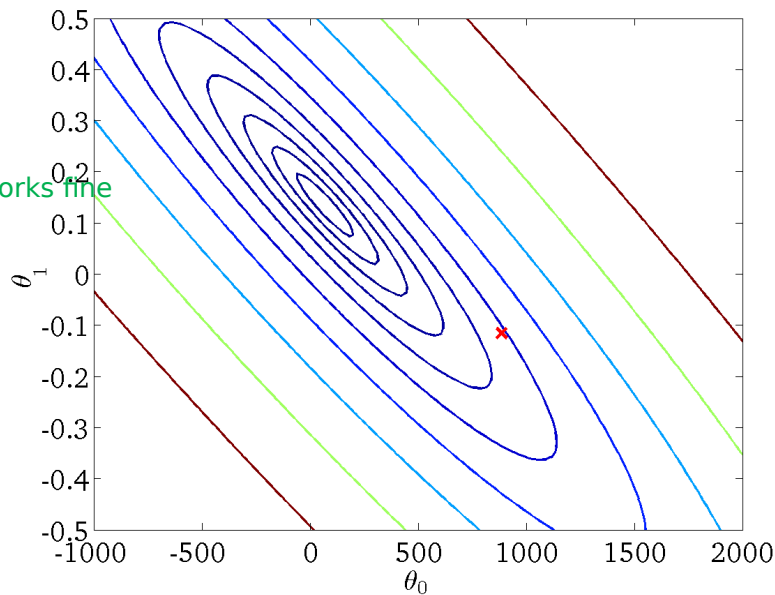
$j = 0, \dots, n$

(for every

)

}

}



Which of the following statements about stochastic gradient descent are true? Check all that apply.

- ☒ When the training set size  $m$  is very large, stochastic gradient descent can be much faster than gradient descent.

Correct Response

- ☒ The cost function  $J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$  should go down with every iteration of batch gradient descent (assuming a well-tuned learning rate  $\alpha$ ) but not necessarily with stochastic gradient descent.

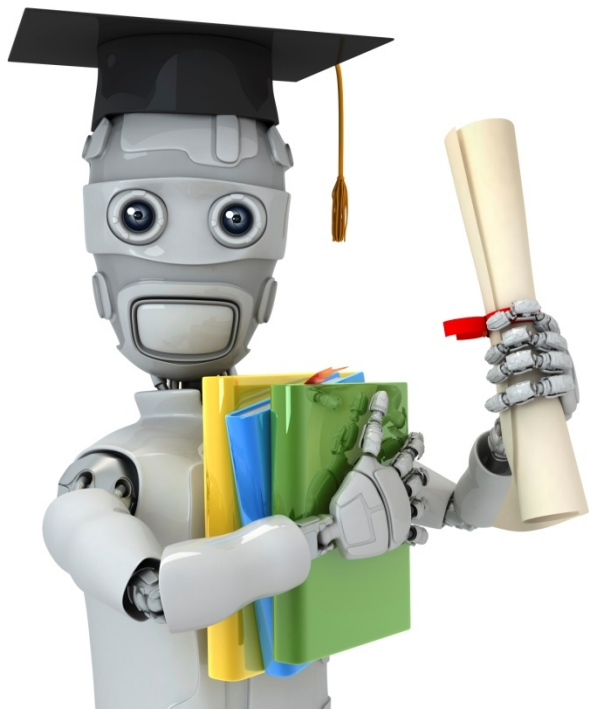
Correct Response

- ☐ Stochastic gradient descent is applicable only to linear regression but not to other models (such as logistic regression or neural networks).

Correct Response

- ☒ Before beginning the main loop of stochastic gradient descent, it is a good idea to "shuffle" your training data into a random order.

Correct Response



Machine Learning

Large scale  
machine  
~~learning~~  
gradient  
descent

## 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

# 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  $j = 0, \dots, n$ )

}

}

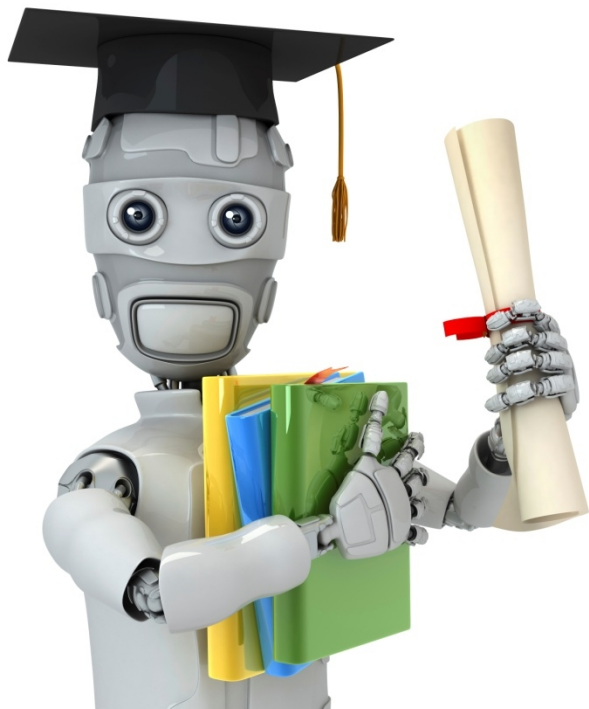
Suppose you use mini-batch gradient descent on a training set of size  $m$ , and you use a mini-batch size of  $b$ . The algorithm becomes the same as batch gradient descent if:

---

- ☐  $b = 1$
- ☐  $b = m / 2$
- ☒  $b = m$

**Correct Response**

- ☐ None of the above



Machine Learning

Large scale  
machine  
~~learning~~  
gradient  
descent  
convergence

# Checking for convergence

Batch gradient descent:

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

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

Stochastic gradient descent:

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

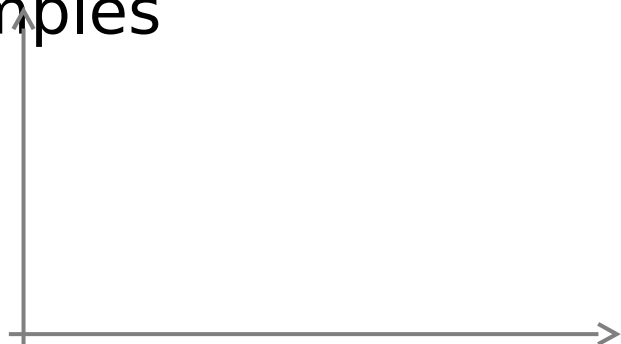
During learning, compute  $cost(\theta, (x^{(i)}, y^{(i)}))$  before updating  $\theta$

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



# Checking for convergence

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



No. of iterations



No. of iterations



No. of iterations



No. of iterations

Diverging

# 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  $j = 1, \dots, m$

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

    for  $j = 1, \dots, n$

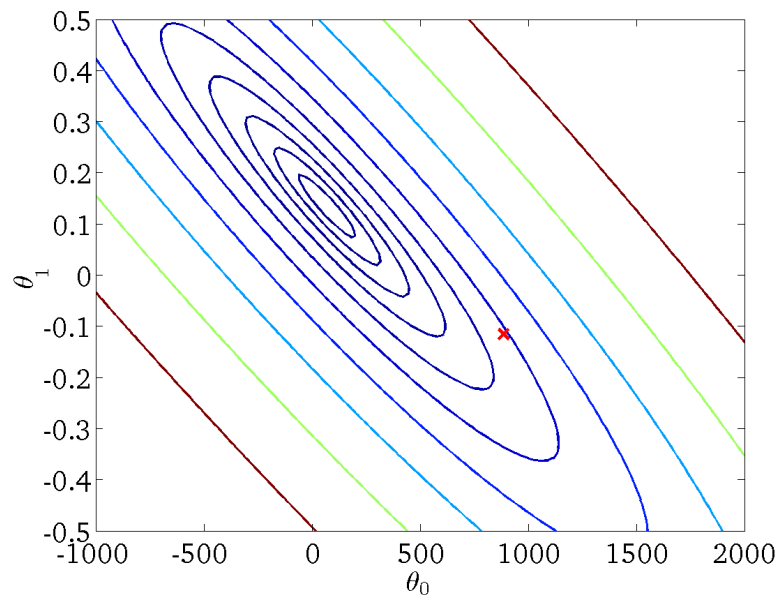
)

}

Learning rate  $\alpha$  is typically held constant. Can

slowly decrease  $\alpha$

over time if we want to converge. (E.g.



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

# 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  $j = 1, \dots, m$

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

    for  $j = 1, \dots, n$

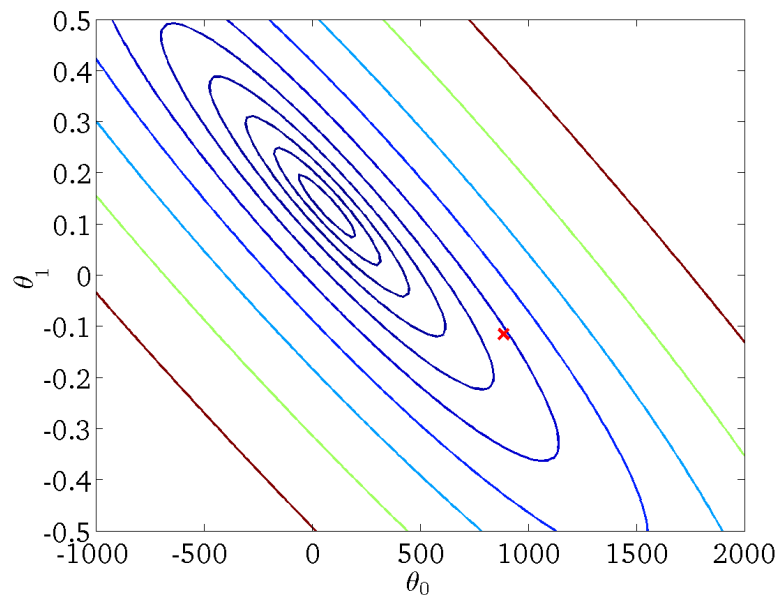
)

}

Learning rate  $\alpha$  is typically held constant. Can

slowly decrease  $\alpha$

over time if we want to converge. (E.g.



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

Which of the following statements about stochastic gradient descent are true? Check all that apply.

- ☐ Picking a learning rate  $\alpha$  that is very small has no disadvantage and can only speed up learning.

Correct Response

- ☒ If we reduce the learning rate  $\alpha$  (and run stochastic gradient descent long enough), it's possible that we may find a set of better parameters than with larger  $\alpha$ .

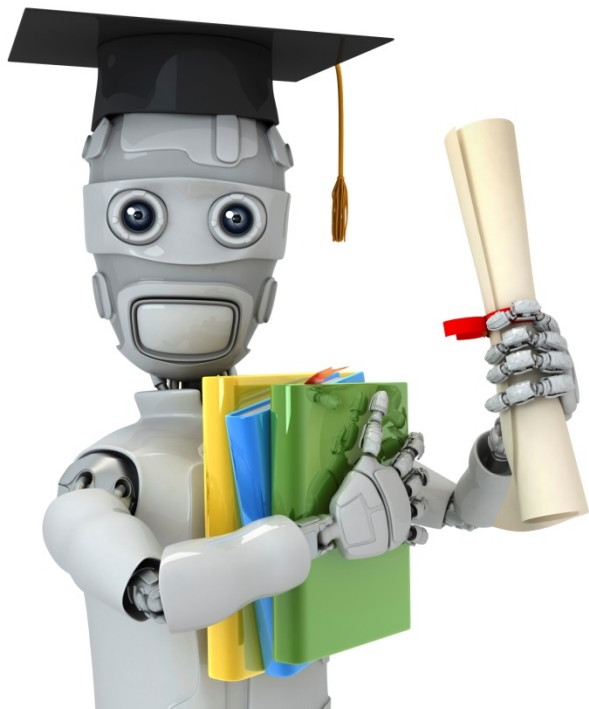
Correct Response

- ☐ If we want stochastic gradient descent to converge to a (local) minimum rather than wander of "oscillate" around it, we should slowly increase  $\alpha$  over time.

Correct Response

- ☒ If we plot  $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$  (averaged over the last 1000 examples) and stochastic gradient descent does not seem to be reducing the cost, one possible problem may be that the learning rate  $\alpha$  is poorly tuned.

Correct Response



Machine Learning

Large scale  
machine  
~~learning~~  
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  $y = 1$  sometimes choose to use your shipping service (  $y = 0$  ) sometimes capture properties of user, of origin/destination and asking price. We want to learn to optimize price.

## 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, 0 otherwise.

Learn  $p(y = 1|x; \theta)$  . Click through rate

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; ...

Some of the advantages of using an online learning algorithm are:

---

- ☒ It can adapt to changing user tastes (i.e., if  $p(y|x; \theta)$  changes over time).

Correct Response

- ☐ There is no need to pick a learning rate  $\alpha$ .

Correct Response

- ☒ It allows us to learn from a continuous stream of data, since we use each example once then no longer need to process it again.

Correct Response

- ☐ It does not require that good features be chosen for the learning task.

Correct Response





Machine Learning

Large scale  
machine  
~~learning~~  
Map-reduce  
and data  
parallelism

# Map-reduce

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

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

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

$$temp_j^{(2)} = \sum_{i=101}^{200} (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

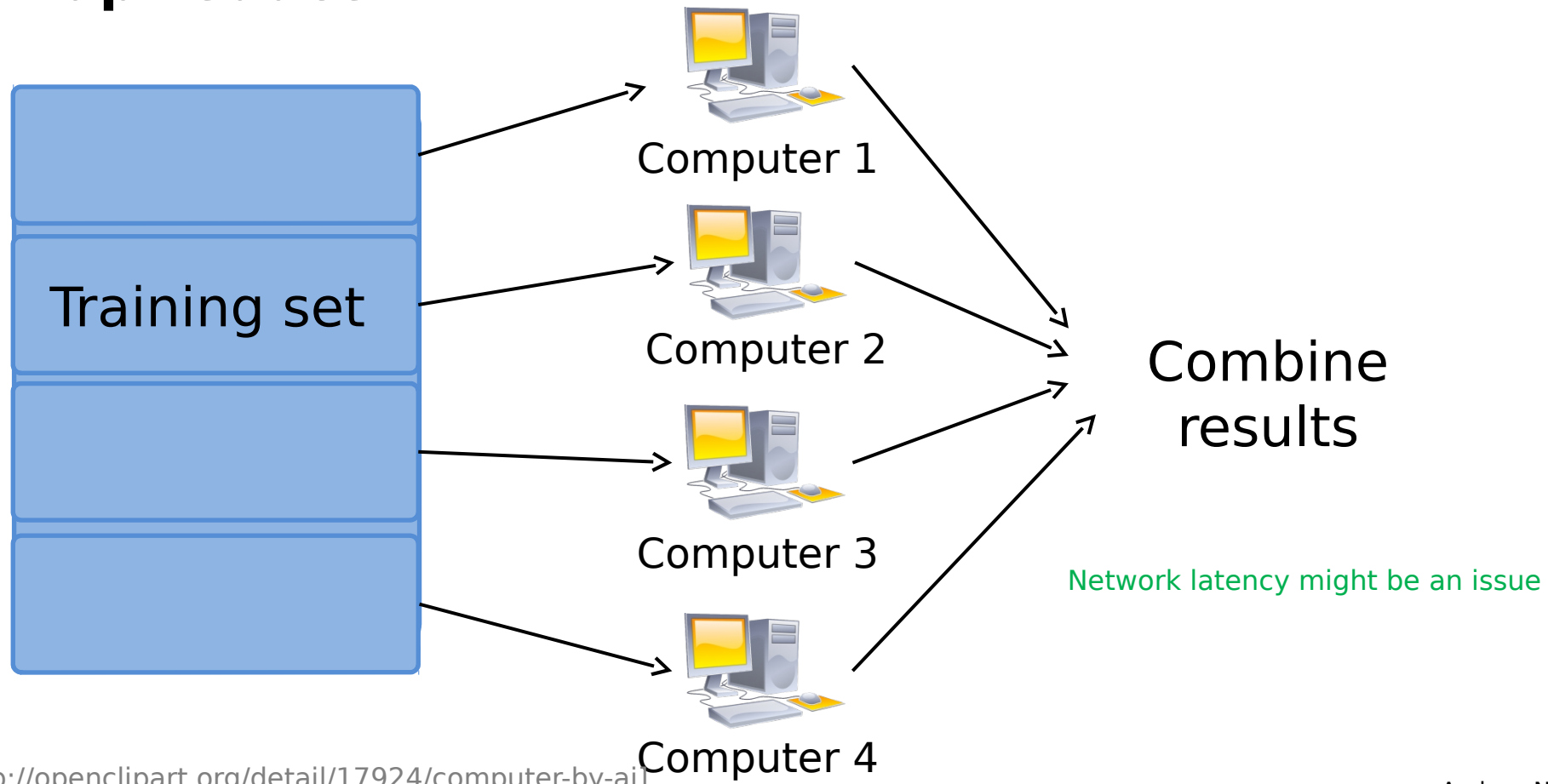
Machine 3: Use  $(x^{(201)}, y^{(201)}), \dots, (x^{(300)}, y^{(300)})$ .

$$temp_j^{(3)} = \sum_{i=201}^{300} (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

Machine 4: Use  $(x^{(301)}, y^{(301)}), \dots, (x^{(400)}, y^{(400)})$ .

$$temp_j^{(4)} = \sum_{i=301}^{400} (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

# 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.

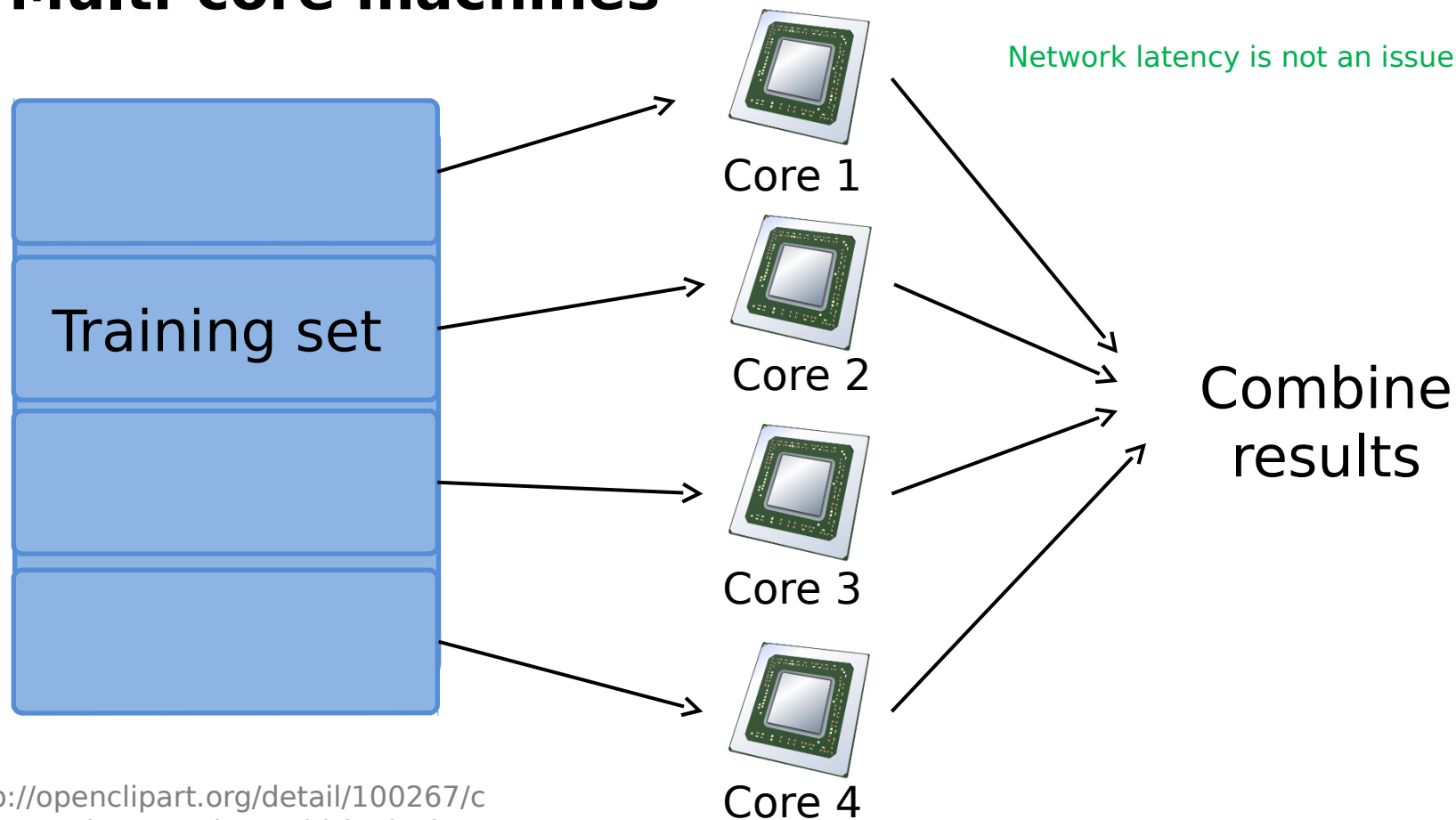
E.g. for advanced optimization, with logistic

regression, need:

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

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

# Multi-core machines



Suppose you apply the map-reduce method to train a neural network on ten machines. In each iteration, what will each of the machines do?

---

- ☐ Computer either forward propagation or back propagation on 1/5 of the data.
- ☒ Compute forward propagation and back propagation on 1/10 of the data to compute the derivative with respect to that 1/10 of the data.

**Correct Response**

- ☐ Compute only forward propagation on 1/10 of the data. (The centralized machine then performs back propagation on all the data).
- ☐ Compute back propagation on 1/10 of the data (after the centralized machine has computed forward propagation on all of the data).