

## Week 10

-- @ learning with large datasets  
 #. large scale machine learning  $\rightarrow$  Algo but viewing with big data sets.

- Dealing with massive datasets.

Best ways to get a high performance machine learning  $\rightarrow$   
 Take low bias learning algo, and train on lots of data.

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

It has computational problems.

# learning with large datasets.

$$n = 100,000,000$$

$$\hat{Q}_j := Q_j - \alpha \sum_{i=1}^m (h(x^{(i)})) f_j^{(i)} x_j^{(i)}$$

A good sanity check  $\Rightarrow$  Randomly pick 1000 e.g. of 100,000,000 and train our algo on 1000 e.g. might do just as well

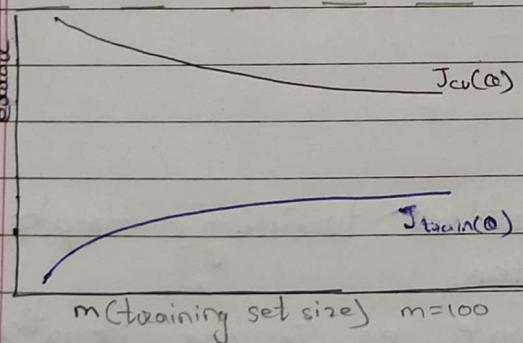
(a) Suppose you are facing a supervised learning problem and have a very large dataset ( $m = 100,000,000$ ). How can you tell if using all the data is likely to perform much better than using a small subset of data (say  $m=1,000$ )?

→ • There is no need to verify this; using large dataset always gives much better performance.

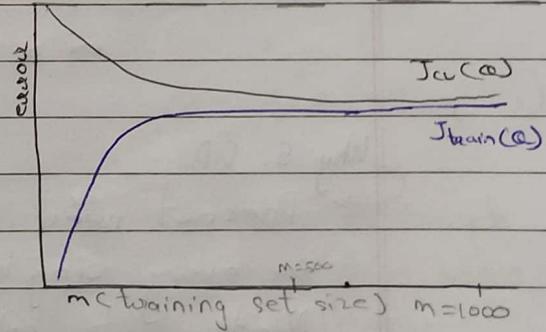
• Plot  $J_{\text{train}}(\theta)$  as a function of the number of iterations of optimization algo. (gradient descent)

• 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 testing range of values of  $m$  and verify that the algo has high variance when  $m$  is small.



High variance learning algo.  
↓ adding extra terms  
↓ Improve performance.



High bias learning algo  
will not increase performance.

stop

To go → Add extra features/hidden layers  
Situation on left

And use more examples.

## (Batch Gradient Descent)

## --@ Stochastic Gradient Descent -- --

# modification to basic gradient descent. (batch gradient descent)  
 applies to logistic reg., neural network, linear reg. (training gradient descent)  
 focus of video

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

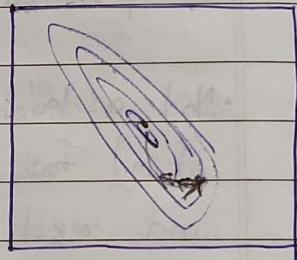
$$J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Repeat {

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

}

m = 300,000,000



sum of 300000000 and single by step,

epoch

Why S.G.D.

→ Does not need to look at all training examples  
in every single iteration.

Needs to look at single eg. in one iteration.

## Batch Gradient Descent

$$J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i - y_i)^2$$

Repeat {

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i) x_j$$

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

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

1. Randomly shuffle dataset (good practice)

2. Repeat {

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

$$\theta_j := \theta_j - \alpha (\hat{y}_i - y_i) x_j$$

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

}

$$\text{cost}(\theta, (x^{(i)}, y^{(i)}))$$

$\frac{\partial}{\partial \theta_j}$

i) looking at only this first example  $(x^{(1)}, y^{(1)})$ , it's gonna take like a basically little gradient step.

what cost of just first e.g.

To fit 1st e.g. a little bit better update  $\theta_0$ .

ii)  $(x^{(2)}, y^{(2)})$

iii)  $(x^{(3)}, y^{(3)})$

Repeat  $\Rightarrow$  takes multiple passes over entire training set

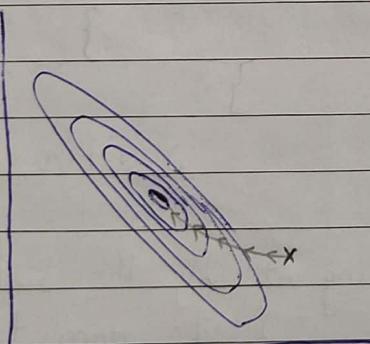
# Stochastic gd :- like

# lot like gd but,

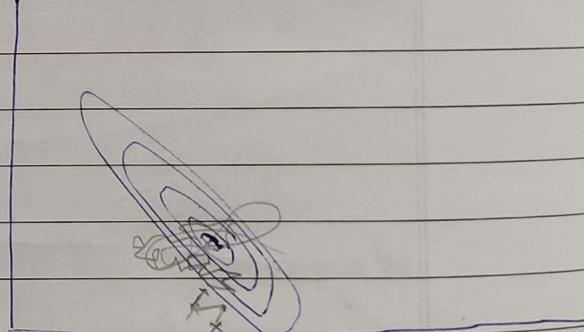
Rather than wait to sum up those gradient terms over all m training e.g., we are taking this gradient term using just one single training e.g., And we are starting to make progress in improving parameters already.

Rather than scan through all m training e.g., before this, we can modify the parameters a little bit and make progress towards global minimum.

We just need to look at single e.g. we're already starting to make progress. Moving the parameters towards global minimum.



Batch GD.



Stochastic GD.

2<sup>nd</sup>, it's unlucky, head in bad direction.  
3<sup>rd</sup>, better

Stochastic GD :

- We find, generally move in direction of global minimum, but not always.
- Doesn't converge in same sense as Batch GD does, wandering around continuously in some region close to global minimum, but it doesn't get to global minimum & stay there.  $\rightarrow$  No Problem

We get pretty good hypothesis.

- We get parameters near global minimum & that's good enough.

How many times to run Repeat { loop }

Depending on size of dataset, doing this loop just single time may be enough!!! (1-10 times)

### Batch GD:

After taking pass through entire dataset

Just one single gradient descent step. SGD > BGD  
(one of baby step)

statements about Stochastic GD?

When the training set size m is very large, SGD can be much faster than BGD.

The cost function  $J_{\text{train}}(\theta) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$  should go down with every iteration of BGD, but not necessarily with SGD.

One e.g. may ↑, & other may ↓, other ↑ in cost in SGD.

# SGD is applicable to LR, Logistic R, NN.

Before beginning of main loop of SGD, its a good idea to shuffle dataset.

## --- Mini Batch Gradient Descent - - - - -

# sometime may run faster than SGD.

Batch GD: Use all  $m$  examples in one each iteration.

Stochastic GD: Use 1 e.g. in each iteration.

Mini batch GD: Use  $b$  examples in each iteration.

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

$b = 10$

usually

$b = (2-10)$

Get  $b=10$  e.g.'s  $(x^{(i)}, y^{(i)}), \dots, (x^{(i+9)}, y^{(i+9)})$

$$\boxed{\begin{aligned} \theta_j := \theta_j - \alpha \frac{1}{10} \sum_{k=i}^{i+9} (\hat{y}_k - y^{(k)}) x_j^{(k)} \\ i = i + 10; \end{aligned}}$$

## \* Mini batch Gradient Descent

Say  $b=10, m=1000$ .

Repeat {

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

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

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

}

$$m = 300,000,000$$

$$b = 10$$

looking at just first 10 examples we can start to make progress. in improving parameter  $\alpha$

Second 10 e.g.'s.

MbGD can be faster than SGD.

# MbGD vs SGD. MbGD likely to outperform SGD only if we have a good vectorized implementation.

b e.g's.

1 e.g.

- Parallelization. (Sum over 10 e.g's can be performed in more vectorized way)

Use libraries to parallelize gradient computations over b examples.

- MbGD disadvantage  $\rightarrow$  Extra parameter, b.

Q Suppose you use mini-Gradient Descent on training set of size m, & use MbGD size of b. The algo becomes

SGD if.  $b=m$

## --@ Stochastic Gradient Descent Convergence --

# making sure  $\rightarrow$  converging.

- Picking  $\alpha$
- Debugging program

### \* Checking for convergence.

# Batch gradient descent:

- Plot  $J_{\text{train}}(\theta)$  as a fun of no. of iterations of gradient descent.
- $J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (\theta_0(x^{(i)}) - y^{(i)})^2$

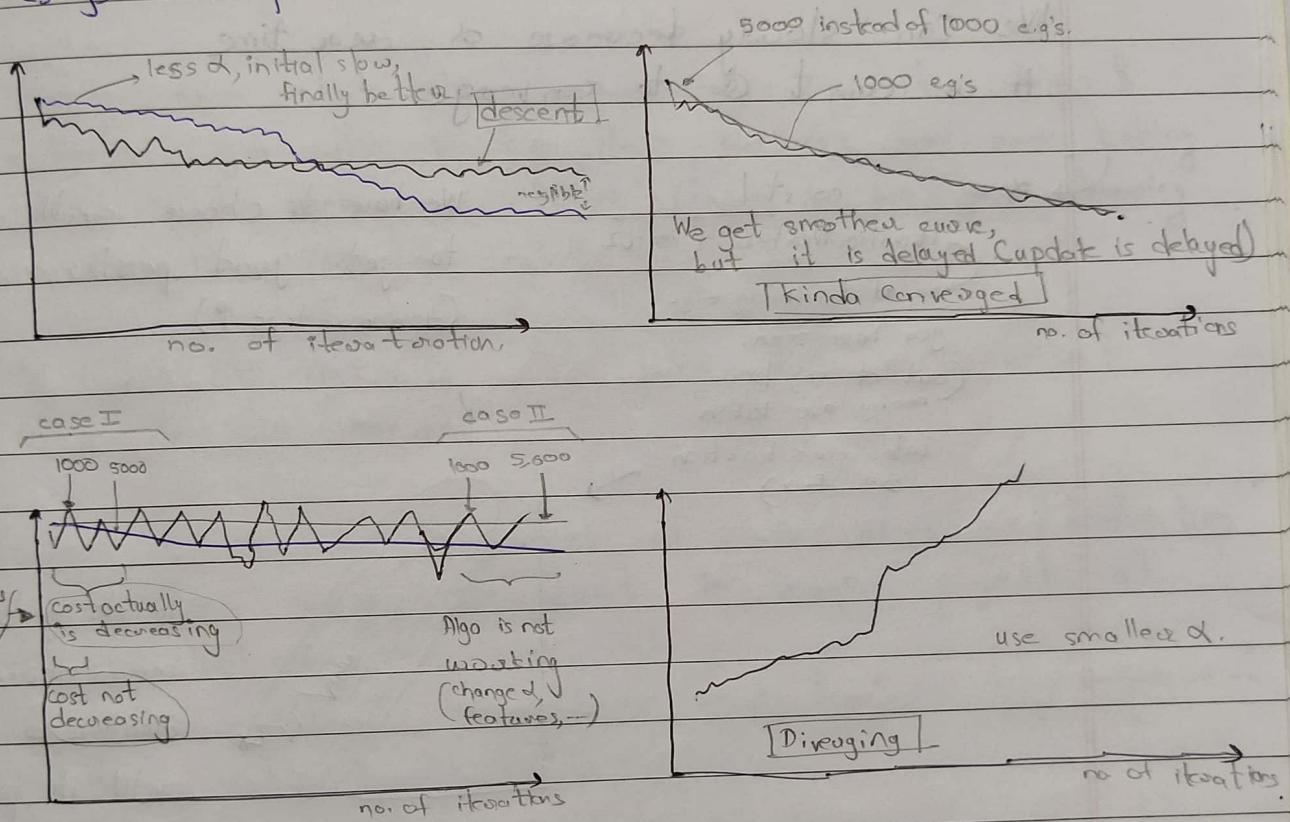
$J(\theta)$  ↓ on  
every  
iteration.

# Stochastic gradient descent :

- $\text{cost}(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2} (\theta_0(x^{(i)}) - y^{(i)})^2 \rightarrow (x^{(i)}, y^{(i)}), (\theta^{(i)}, y^{(i)})$
- During learning, compute  $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$  before updating  $\theta$  using  $(x^{(i)}, y^{(i)})$ .
- Every 1000 iterations (say), plot  $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$  averaged over the last 1000 examples processed by algorithm.

#Checking for convergence in Stochastic GD.

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



Stochastic GD doesn't just converge to global minimum.

Parameters oscillate around global minimum.

small alpha, you'll end up with small oscillations.

### 1) Learning rate $\alpha$ :-

We get close to global minimum, but we don't get exactly global minimum.

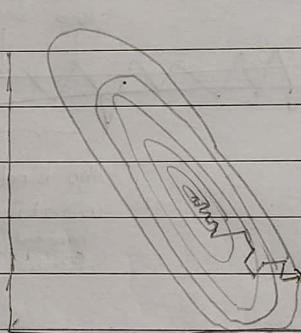
Learning rate  $\alpha$  is typically held constant. (We never change  $\alpha$  while algo. is running)

We can slowly decrease  $\alpha$  over time if we want  $\alpha$  to converge.

$$\alpha = \frac{\text{const1}}{\text{iteration Number} + \text{const2}}$$

↳ no. of training e.g. you have seen.  
(out of m, how many are taken into consideration so far)

We have to choose const1, const2 to get good performance.  
(#parameters ↑)



As it gets closer mean deviating gets smaller & smaller until it's pretty much close to global minimum.

Alpha get smaller, so step size  $\beta$ 's

But still keep  $\alpha$  const, is more common (const1, const extra param)

### 2) Choose about stochastic GD.

- Picking a  $\alpha$  very small has no disadvantage & can speed up learning.
- ✓ If we reduce  $\alpha$  (can step long enough). It's possible that we may find a set of better parameters  $\omega$  than with larger  $\alpha$ .
- If we want SGD to converge to a local minimum rather than wander or "oscillate" around it, we should slowly increase  $\alpha$  over time.
- ✓ If we plot  $\text{cost}(\alpha, \mathbf{x}^{(i)}, \mathbf{y}^{(i)})$  (averaged over the last 1000 examples) and SGD doesn't seem to be reducing the cost, one possible problem may be that the  $\alpha$  is poorly tuned.

## -@ Online learning - - - -

#large scale ml setting. Use only if data coming is very large.  
 • continuous stream of data is coming. (We learn & discard  
 (if data set is big) assumption.)

#online learning.

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

We want to optimize asking price that we want to offer to our users

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

We will use logistic Regress. price ~~we're asking for~~

• Algo flow.

Repeat forever {  
 get  $(x_i, y_i)$  corresponding to user.  
 update  $\alpha$  using  $j$ th  $(x_i, y_i)$   
 $\alpha_j := \alpha_j - \alpha(h_\alpha(x) - y)x_j$  }  $(j=0, \dots, n)$

• Can Adapt to changing user preference.

e.g. 2) Other online learning example.

Product search (learning to search)

User searches for "Android phone 100 pcam" → type search  
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.)

$y=1$  if user clicks on link.

$y=0$  otherwise.

Learn  $p(y=1 | x; \theta)$  Estimate probability that user will click on link of mobile phone.

If  $y=1$ , show in search result.

This is called "The problem of learning the predicted click-through rate, the predicted CTR"

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

Every time user searches we show 10 results,

We will get  $10 (x, y)$  pairs & use online learning algo to update parameters using 10 steps of GD on these  $10 (x, y)$  pairs e.g.'s & then throw the data away.

Other eg.) Choosing special offer to show user; customized selection of news article; product recommendation; -

If you have collaborative filtering, you'd give you additional features to feed into logistic regression classifier to try to predict CTR for diff products.

Any of these problems can be formulated as standard ml problem where you have fixed dataset.

- Some of advantages of using an online learning algorithm
- ✓ It can adapt to changing user taste  
( $p(\theta|x)$  changes over time)

- There is no need to pick a d.
- ✓ 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.
- It doesn't require that good features be chosen for learning task.

# Similar algo to SGD (training set doesn't change)

400,000,000 → 400 (we are using)

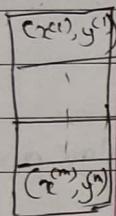
## --① Map Reduce and Data Parallelism -- --

## ★ Map Reduce:

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

Machine 1: Use  $(x^{(1)}, y^{(1)}), \dots, (x^{(400)}, y^{(400)})$ .

$$\text{tmp}_j^{(1)} = \sum_{i=1}^{400} (\text{h}_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$



Machine 2: Use  $(x^{(1)}, y^{(1)}), \dots, (x^{(200)}, y^{(200)})$ .

$$\text{tmp}_j^{(2)} = \sum_{i=1}^{200} (\text{h}_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

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

$$\text{tmp}_j^{(3)} = \sum_{i=201}^{300} (\text{h}_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

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

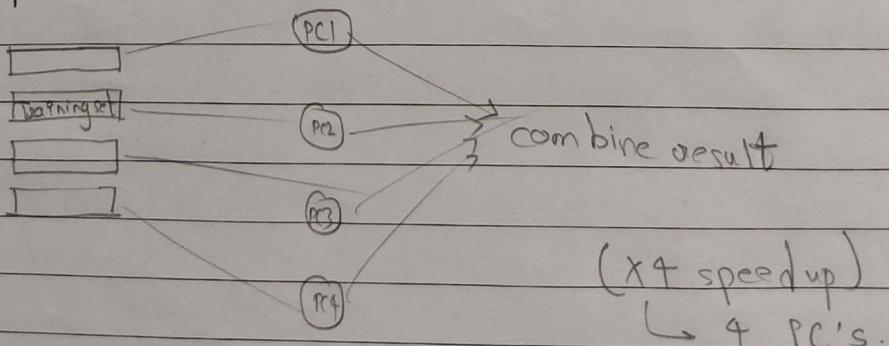
$$\text{tmp}_j^{(4)} = \sum_{i=301}^{400} (\text{h}_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Combine.

$$\theta_j := \theta_j - \alpha \frac{1}{400} (\text{tmp}_j^{(1)} + \text{tmp}_j^{(2)} + \text{tmp}_j^{(3)} + \text{tmp}_j^{(4)})$$

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

## # Map Reduce.



# Map reduce and summation over the training set.

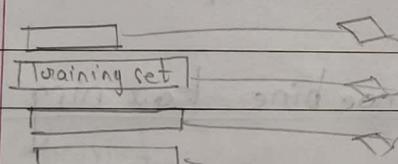
Many learning algorithms can be expressed as computing sums of functions over the training set.

e.g. for advanced optimization, with logistic regn, need:

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

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

# Multi-core machines.



combine & until ready

Q Suppose you apply mapreduce method to train a neural network on ten machines. In each iteration, what will each of machines do?

→ Compute forward propagation & back propagation on  $\frac{1}{10}$  of the data to compute the derivative with respect to that  $\frac{1}{10}$  of the data.

• Compute either forward prop / back prop on  $\frac{1}{5}$  of data.

• Compute only FP on  $\frac{1}{10}$  of data & BP on all data → centralized

• compute BP on  $\frac{1}{10}$  of data & FP