

My
Wing
know
our score.

Week 10

large datasets
have computational
problems

Date: _____

Learning with large Data Sets

ML 8 data:

To get high performance ML algo:

Take low bias algo & train on lots of data

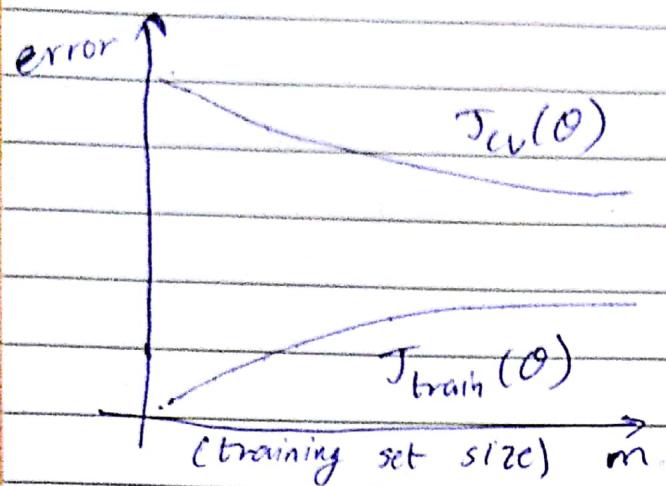
"It's not who has the best algorithm that wins.
It's who has the most data."

Suppose $m = 100,000,000$ & you are applying linear regression via ~~the~~ gradient descent:

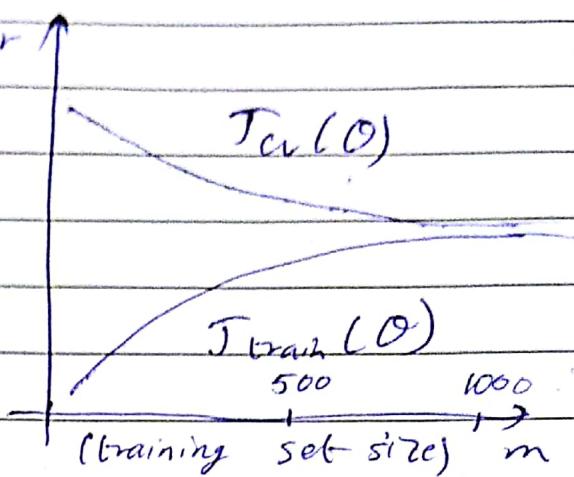
$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

↳ to take 1 step of gradient descent
you need to sum over 100,000,000

Why not just randomly pick 1,000 examples from data set & use those for ~~linear~~
training linear regression?



High variance issue



High bias issue

If the problem is high variance then stick to $m = 100,000,000$ as the use of more examples will decrease $J_{cv}(\theta)$

If the problem is high bias then it is possible to keep $m = 1000$ as there is no point in using more data.

Naturally if we have a scenario that resembles the situation to the right we will try to change the scenario as such to make it look like the situation on the left (high variance)

when dealing with large data sets we need to find computationally efficient ways of training ML algos

modification to gradient descent that scales to
bigger training sets
→ applies to NN, logistic regression &
other functions just as Date: _____ gradient
descent

Stochastic Gradient Descent

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

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

Repeat {

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

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

}

Normal gradient descent is also called batch
gradient descent because all training examples
are considered ~~at~~ in a gradient descent step

Suppose $m = 300,000,000$ then you can't
bring all records into memory but instead
instead you will have to read from disk
↳ slows down process even further

Stochastic Gradient Descent

Date: _____

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

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

Random shuffling ensures that each record scanned is visited in randomly sorted order
 ↳ speeds up stochastic gradient descent by a bit

Sees how well hypothesis is doing on single example $x^{(i)}, y^{(i)}$

Sequence: ✓

1) Randomly shuffle dataset

2) Repeat { ←

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

(External loop)
 causes program to run over entire training set multiple times

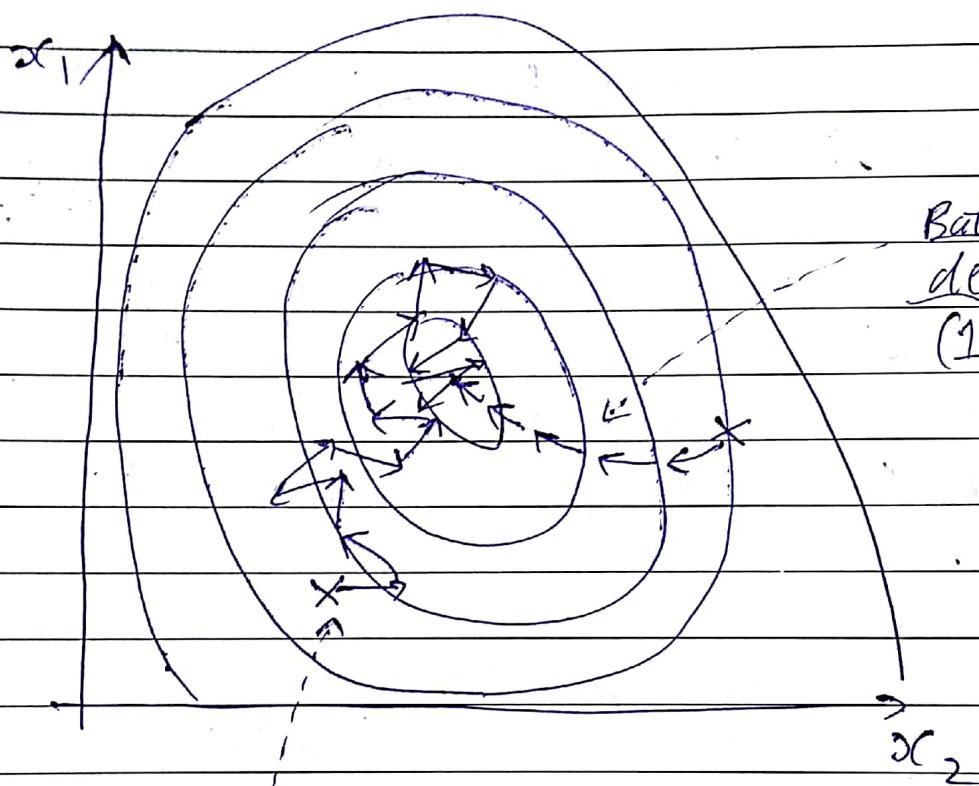
$$\theta_j := \theta_j - \alpha (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

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

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

Stochastic gradient descent will look at the first example only & take a step based on only that first example. Then modify $\theta_1, \dots, \theta_n$ based on only 2nd example & etc until it goes through whole data set

Stochastic gradient descent - instead of waiting for all values to be summed in data set we start improving the parameters early after iterating over a single data point.



Stochastic gradient descent steps

Stochastic gradient descent will generally move towards global minimum but not always, more specifically in that direction at a point in time. You can get unlucky based on the randomly chosen record that is used to retrieve cost.

Date: _____

Stochastic gradient descent does not converge as is the case with batch gradient descent, instead it wanders around in ~~the~~ some region that is close to global minimum - does not get to global minimum & stay there

↳ gives θ values close to global minimum & that is enough

In stochastic gradient descent, it depends on ~~how~~ m , how many times to go through the outer loop i.e. how many times to go through entire training set.

$m \uparrow$ then less outer loop iteration

1-10 passes is fairly common

Stochastic gradient descent is faster because outer loop is repeated less times as parameter update happens m times when going through entire training set while batch gradient descent updates only once

Date: _____

Mini-batch gradient Descent

↳ can sometimes work even faster
than stochastic 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 is in between batch &
stochastic. Batch ~~is~~ (b) usually
ranges from 2 - 100

$$b = 10$$

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

$$\rightarrow \theta_j := \theta_j - \alpha_{10} \sum_{k=i}^9 (\hat{y}_k - y^{(k)}) \delta_j^{(k)}$$

$$i := i + 10$$

→ update parameters after b examples

Date: _____

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 $j=0, \dots, n$)

}

}

Minibatch can outperform stochastic gradient descent only if implementation of minibatch is well vectorized

→ the sum over b examples is calculated in a vectorised manner to parallelize gradient computations

Stochastic Gradient Descent

Convergence

→ how do you know that stochastic gradient descent is completely debugged & is converging?

→ how do you tune α in stochastic gradient descent

Checking for convergence

Batch gradient descent:

- Plot $J_{\text{train}}(\theta)$ as a function of # of iterations of gradient descent
- $J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

Suppose $m = 300,000,000$

When updating parameters while stochastic grad. descent you don't want to pause periodically to calculate cost because calculating cost requires summation over all training examples hence becomes costly when **victory** m is large

Stochastic gradient descent:

$$-\text{cost}(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2} (\theta_0(x^{(i)}) - y^{(i)})^2$$

- during learning, compute

$\text{cost}(\theta, (x^{(i)}, y^{(i)}))$ before updating θ
using $(x^{(i)}, y^{(i)})$

{ every 1000 iterations, plot

$\text{cost}(\theta, (x^{(i)}, y^{(i)}))$ averaged over the
last 1000 examples processed by algo

cost computed before θ update

so that cost is computed for data

set examples that it has not seen yet

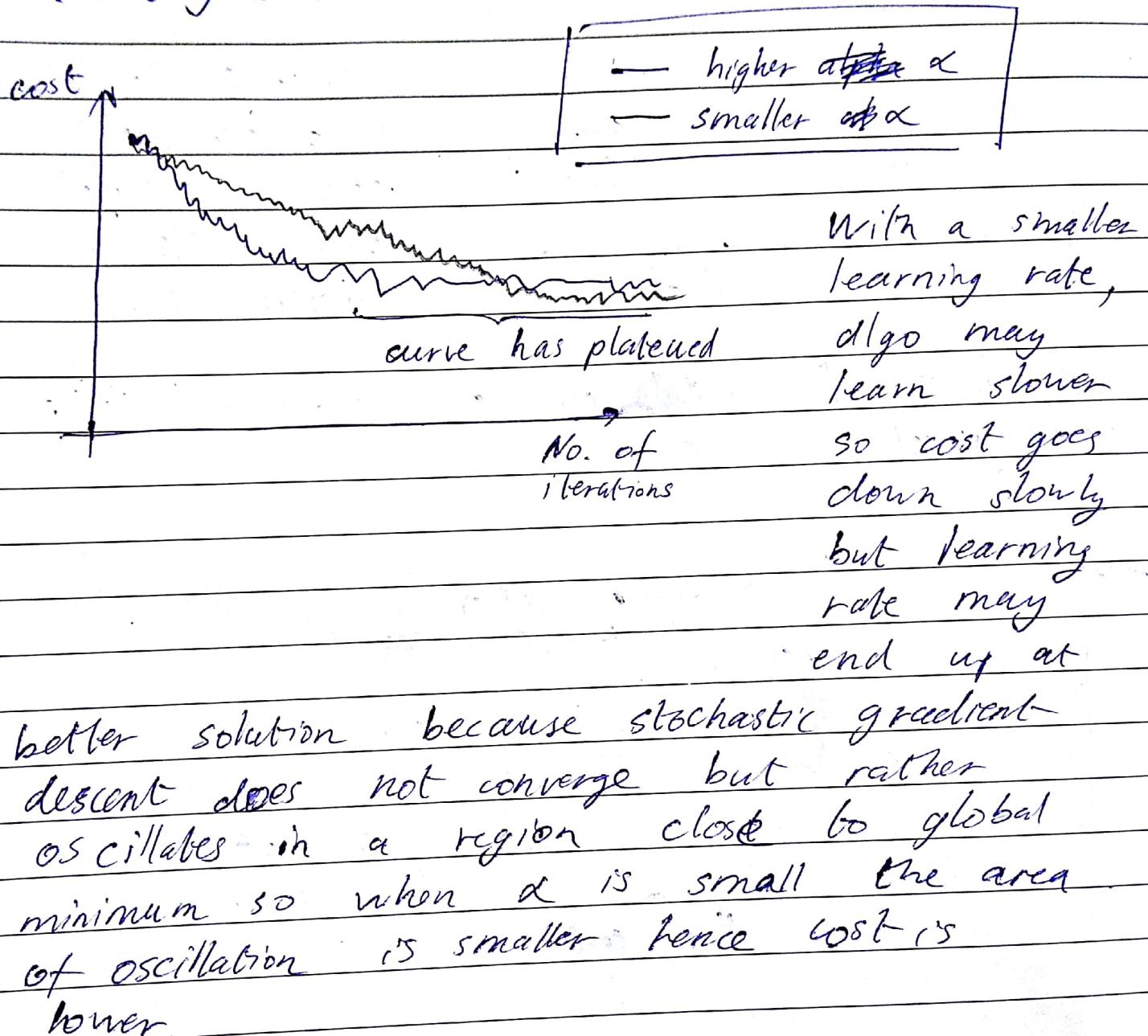
estimates how

well algo is doing

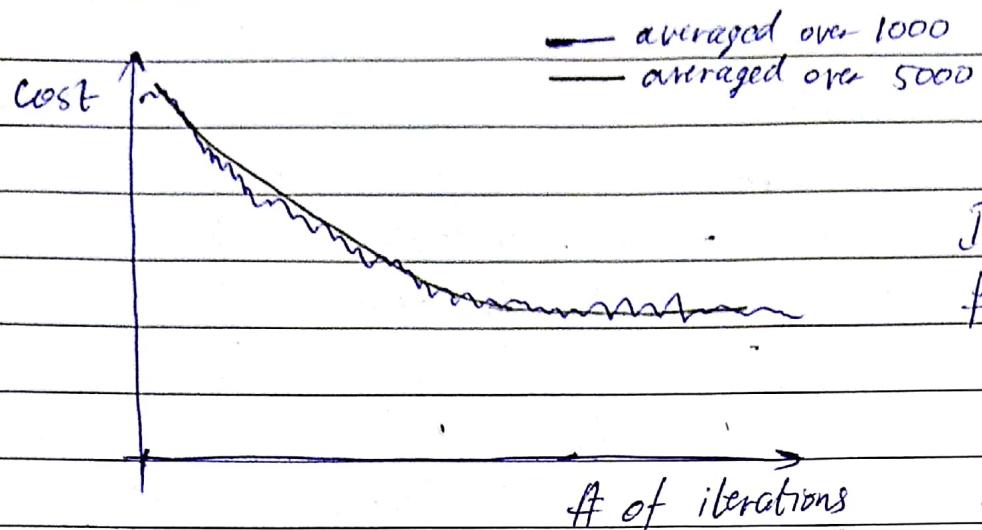
you don't need to iterate over entire
data set now

noise is introduced, will not
necessarily decrease on every
iteration

Checking for convergence

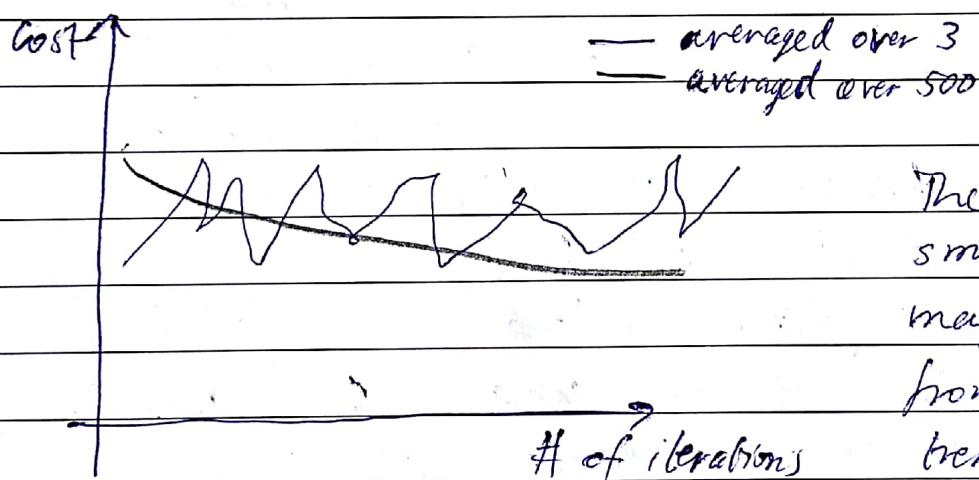


Date: _____

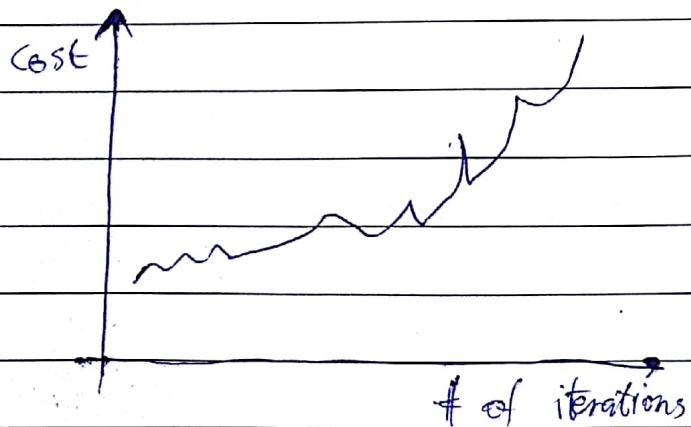


Increasing the # of examples cost is averaged over creates a smoother curve

Disadvantage - fewer data points over entire training set



The noise with small value averages may prohibit us from seeing the trend



Shows that algorithm is diverging - use smaller value of α

Stochastic gradient descent does not converge but rather wanders around the global minimum, we can make it converge by decreasing α & get a better value

Date:

Learning rate α is typically held constant.

Can slowly decrease α over time if we want α to converge

$$[\alpha \rightarrow 0]$$

→ e.g. $\alpha = \frac{\text{constant 1}}{\text{iteration Number} + \text{constant 2}}$

→ The reason why people don't do this sometimes is because then you have additional parameters to tune algo that you need to tinker with

Online Learning

→ large-scale ML setting that allows to model problems where algorithm must learn from a continuous stream of data

In online learning, we discard the notion of there being a fixed training set. You get an example, learn using that example & then discard it.

→ this is only practically implemented when you have a high-width, high-velocity data stream available

Scenario

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

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

Use logistic regression.

Repeat forever {

Get (x, y) corresponding to user

Update θ using (x, y) :

$$\theta_j := \theta_j - \alpha (\text{logit}(x) - y) x_j \quad (j = 0, \dots, n)$$

}

Online learning platforms can adapt to changing user preferences

↳ θ will reflect the change in preference

Online learning is very similar to stochastic gradient descent, only instead of scanning a fixed training set, we're getting one example from a user.

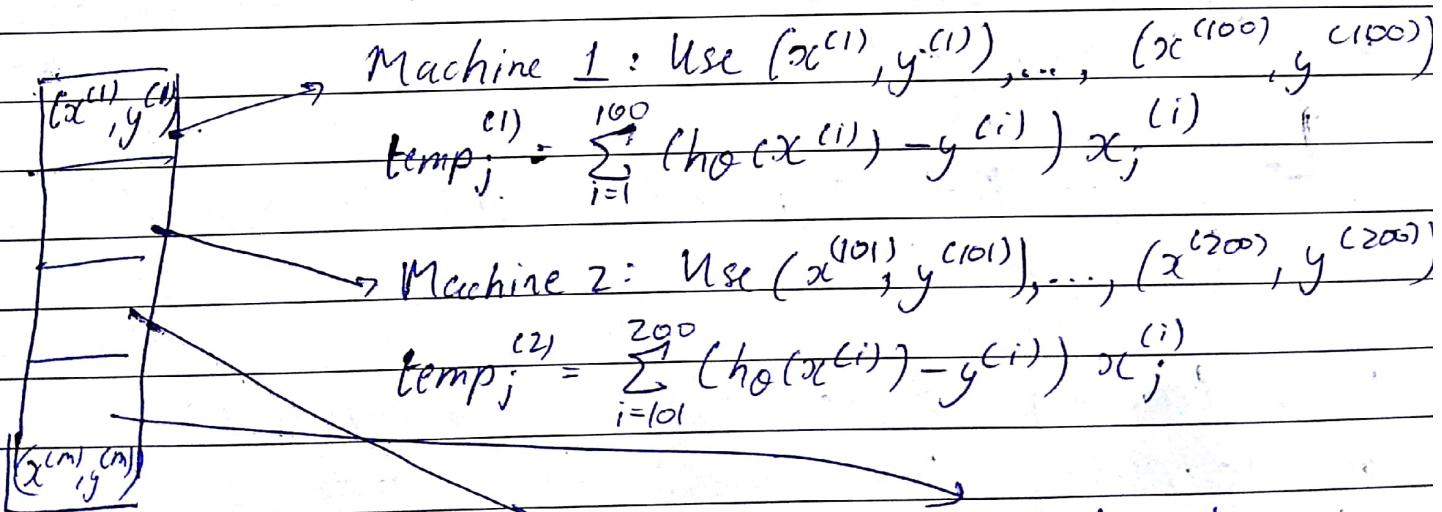
Map-reduce & Data Parallelism

→ run ML algorithms on multiple machines

Batch gradient descent: ($m = 400$)

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

↳ if m is large then this becomes computationally expensive



Split training set into subsets

$$\text{temp}_j^{(3)}$$

$$\text{temp}_j^{(4)}$$

Date: _____

The architecture used works 4 times as fast now, considering no network latencies.

Master Server will combine all the $\text{temp}_j^{(i)}$ values:

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

$$[j = 0, \dots, n]$$

Map-reduce & summation over the training set

If you wish to apply Map reduce to some learning algo, you must find out if learning algo can be expressed as summation over training set

↳ if yes then Map reduce can be applied to scale learning algo

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

E.g:

$$J_{\text{train}}(0) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_0(x^{(i)}) - (1-y^{(i)}) \log (1-h_0(x^{(i)}))$$

↳ The ~~compu~~ training set can be divided amongst PCs & computation executed then recombined on central server

~~Map~~ Until now map reduce has been discussed as a way of parallelizing over multiple computers in a computer cluster but it is also applicable on CPUs with multiple cores

- ↳ each cores works on a share of the training set
- ↳ problem of network latency goes away when working with cores

At times some really good linear algebra libraries ~~can't~~ automatically parallelize the work over multiple cores so all you need to do is implement the vectorized code & the libraries will automatically parallelize it for you