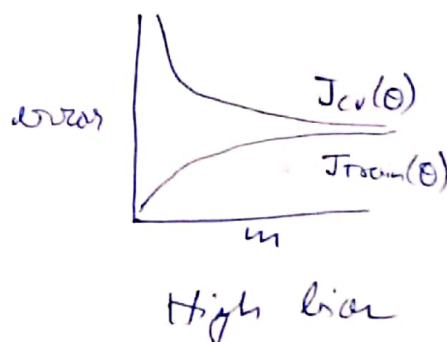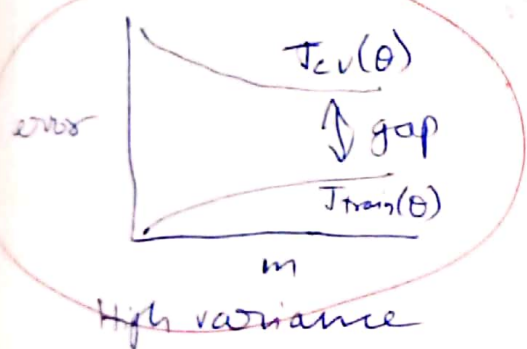# Weak 10

## Gradient Descent with large datasets

- Large datasets can improve accuracy

Lets say we have m = 100 million training examples. Inorder to check if more data can help, we need to check if it has high variance.

We do this by first taking a smaller dataset m = 1000, plotting the learning curve, if its high variance - only then will more data (100 million) result in better accuracy



High variance



High bias

## Stochastic Gradient Descent

Batch gradient descent is very computationally expensive, since if we have M = 300 million training examples, it'll take 300 million iterations to get $\theta$ values
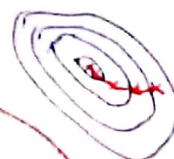
Linear Regression with gradient descent (batch)

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

Repeat {

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

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

}



Happens 300 million times

Stochastic gradient descent:

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

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

1. Randomly shuffle dataset
2. Repeat {

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

  <span style="color:red">This is nothing but</span>
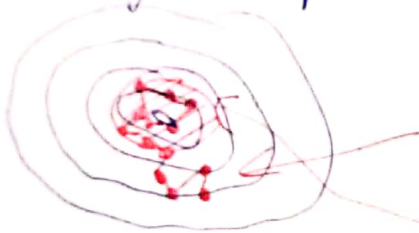  $$\color{red}\frac{d}{d\theta_j} cost(\theta, (x^{(i)}, y^{(i)}))$$

  $$\theta_1 := \theta_j - \alpha \boxed{(h_\theta(x^{(i)}) - y^{(i)})x_j^{(i)}}$$

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

  }
}

Here for every iteration, it tries to fit one training example better

<span style="color:red">It doesn't need to reduce every step, it can even increase. But after a few iterations it reaches an area that is close the global minima</span>

∴ <span style="color:red">we can repeat 1 ~ 10 times</span>

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 6 examples in each iteration

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

Repeat {

    for $c = 1, 11, 21, 31, \ldots 991$ {   Depends on

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

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

    }

}

- If $M = 300$ million, we can start making progress after only looking at $b = 10$ examples rather than all 300 million for every iteration

- Mini-batch gradient descent is likely to outperform stochastic gradient descent only if you have a good vectorised implementation. So you can use the numerical algebra libraries and parallelize your gradient computations over $b$ examples. You can use $b = 10$

## Stochastic Gradient Descent Convergence
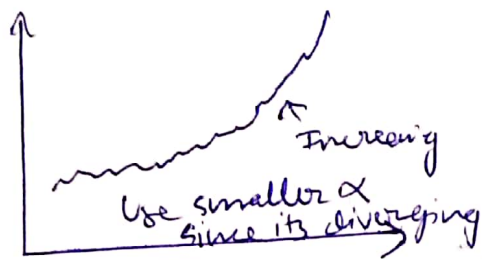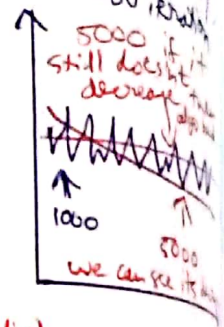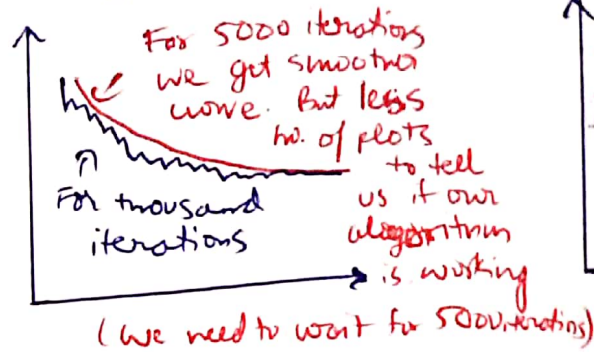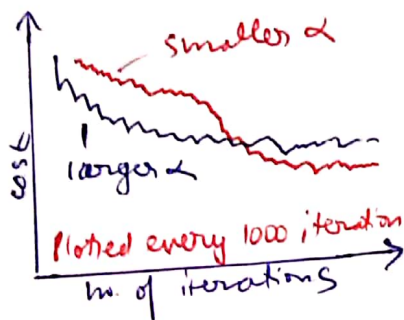
Checking for convergence:

Before in batch gradient descent we would check if cost is decreasing every iteration. We can't do this for stochastic since a) It may not always decrease b) For large no. of examples ($m$), it can slow down the process.

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

Different cases:

Plotting cost $(\theta, (x^{(i)}, y^{(i)}))$ over no. of iterations
averaged over different no. of ~~that~~ plots (eg: every 1000 iterations)



smaller $\alpha$

larger $\alpha$

Plotted every 1000 iteration

m. of iterations

For 5000 iterations we get smoother curve. But less no. of plots to tell us if our algorithm is working

For thousand iterations

is working

(we need to wait for 5000 iterations)

5000 it still doesn't decrease

1000

5000 we can see it



Increasing

Use smaller $\alpha$ since its diverging

Before we found out that stochastic gradient descent doesn't converge to the global minima but just oscillated near it. However one way to avoid this is by slowly decreasing $\alpha = \dfrac{const 1}{iterationNumber + const 2}$ and because of the smaller steps, it oscillate closer to the minima. However people don't use this since configuring const1 & const 2 is extra work and the estimated minima given by stochastic gradient descent is good enough.

# Online learning

If we have a stream of data coming we use this algorithm.

Repeat forever {

    Get $(x,y)$ corresponding to user

    Update $\theta$ using $(x,y)$

        $\rightarrow \theta_j := \theta_j - \alpha (h_\theta(x) - y) \cdot x_j \quad (j = 0, \ldots, n)$

}

It can adapt to changing user performance

Ex:

~~other plan~~
## Product Search

$\rightarrow$ User searches for "Android phone 1080p camera"

$\rightarrow$ we have 100 phones in store. Will return 10 results.

$\Rightarrow x =$ features of phone, how many words in user query match name of phone, how many words in query match description of phone, etc

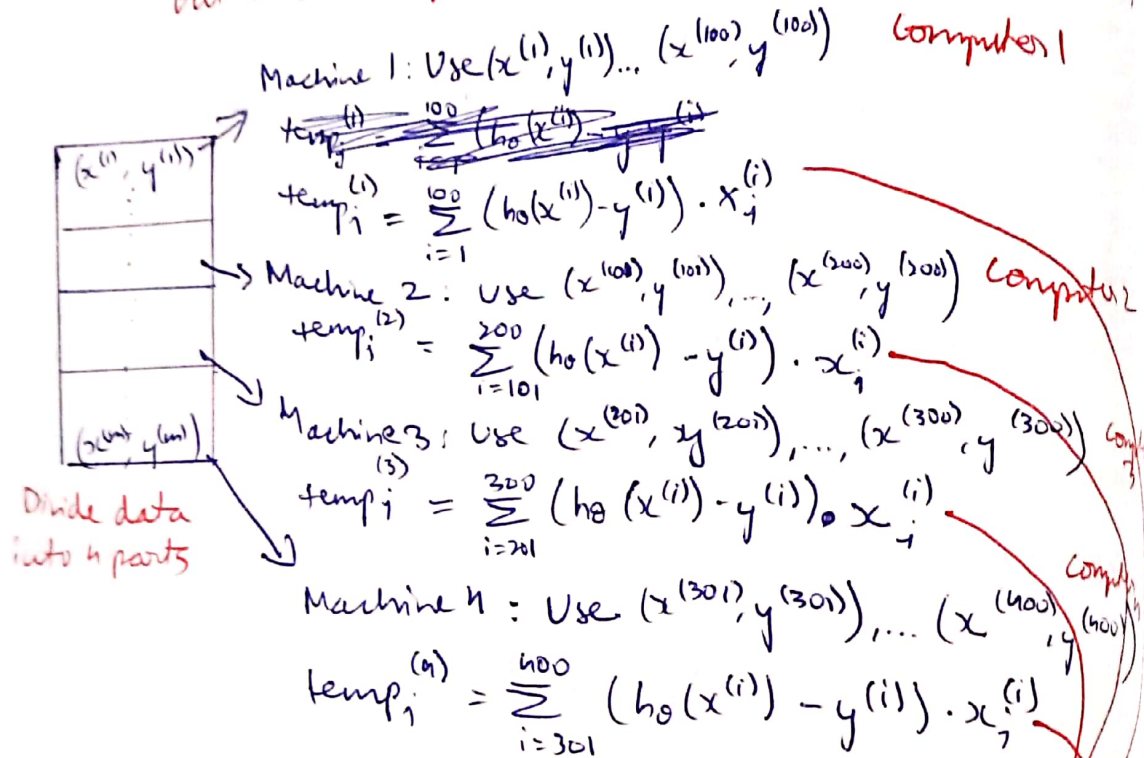$\rightarrow y = 1$ if user clicks on link, $y = 0$ otherwise

$\rightarrow$ Learn $p(y = 1 | x; \theta) \leftarrow$ predicted CTR (click through rate)

$\rightarrow$ 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; ..

# Map reduce & data parallelism

If we want to run the algo on multiple computers

## Map reduce

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

Here we are taking $m = 400$ as an example but in reality this can be big like $m = 400$ million

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

$temp_j^{(1)} = \sum \frac{100}{i=1}(h_\theta(x^{(i)}) - y^{(i)})$

$temp_j^{(1)} = \sum_{i=1}^{100} (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$

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

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

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

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

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

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

(left margin)
$(x^{(1)}, y^{(1)})$

$(x^{(400)}, y^{(400)})$

Divide data into 4 parts

---

**Centralised server**            combine

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

$(j = 0, ... 4)$

∴ You get an almost 4x speed (ignoring network latency, etc).

$\rightarrow$ You can use map reduce where summation (summation is the computationally intensive)

$\rightarrow$ You can use it in a multi-core computer (some numerical algebra libraries use this)