

Week 10

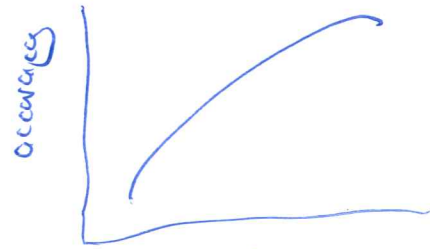
Large Scale machine learning - learning with large datasets

Machine learning and data

classify between two 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"



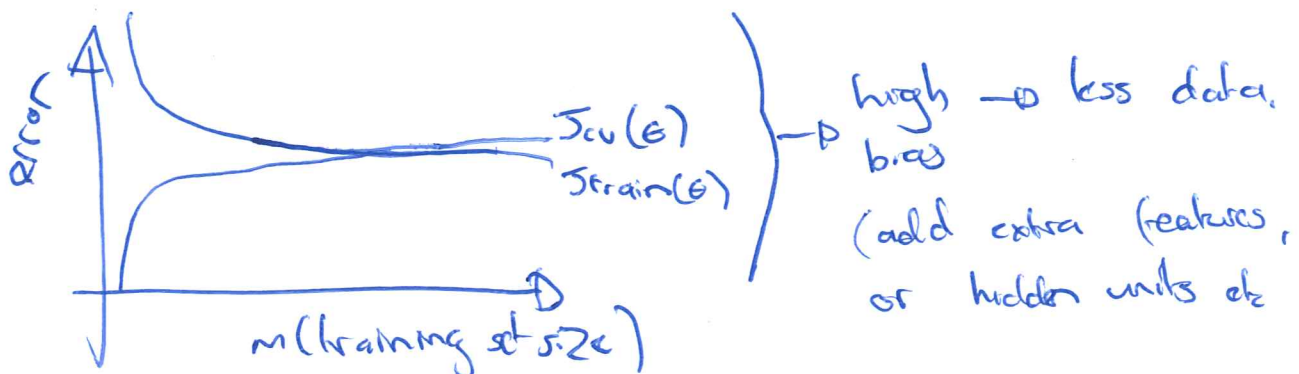
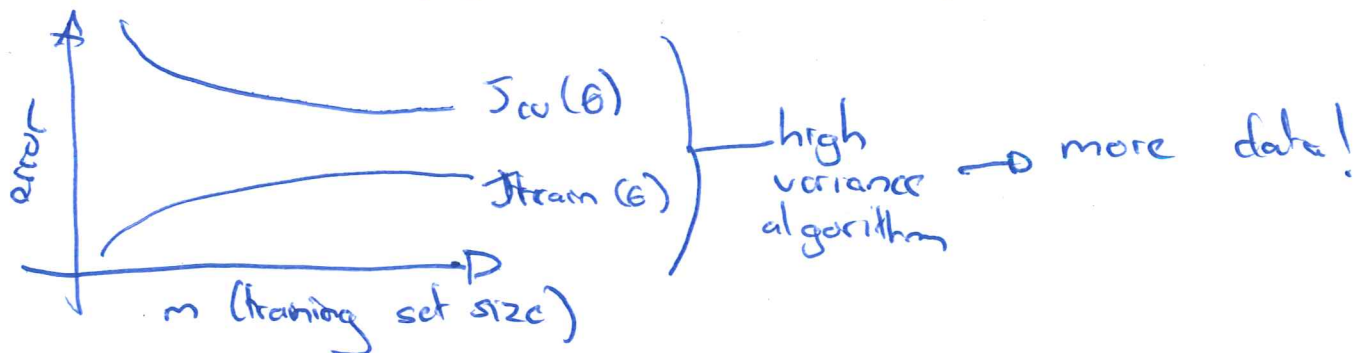
learning with large datasets

$m = 100,000,000$ (lets say website)

and train gradient descent $100,000,000$

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

why not use $m=1000$?



Week 10

Large scale machine learning - stochastic gradient descent

Linear regression with gradient descent

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

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

↗ cost function

Gradient Descent: ("Batch or all")

Repeat {

$$\theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \right] \quad \left\{ \frac{\partial}{\partial \theta_j} J_{\text{train}}(\theta) \right.$$

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

↗ If m is large, gradient descent is computationally expensive.

$M = 300,000,000$ population of U.S.

Batch gradient descent "All"

Our Alternative is: stochastic gradient descent →

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

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

1) Randomly shuffle dataset. ↗ interest of safety

2) Repeat {

for $i=1, \dots, m$ ↗ repeat inner loop 1 to m times

$$\theta_j := \theta_j - \alpha \left[(h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} \right] \quad \left\{ \frac{\partial}{\partial \theta_j} \text{Cost}(\theta, (x^{(i)}, y^{(i)})) \right.$$

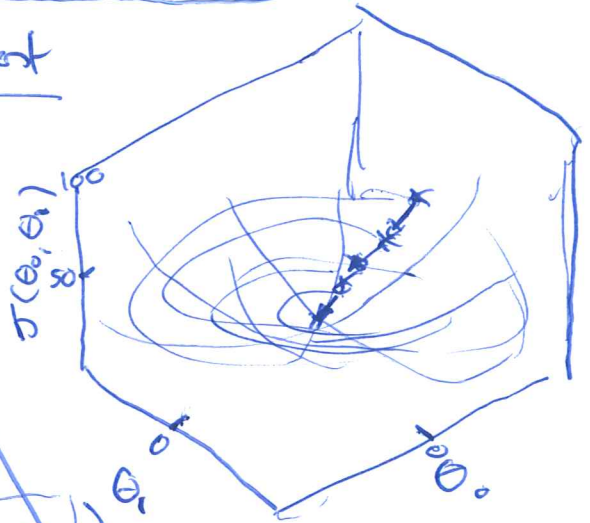
(for $j=0, \dots, n$)

↗ so rather than waiting for the sum to modify parameters, we modify with each training example

①

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

↗ data looks at the next training example & takes another step in parameter space to fit data



takes a more random path to the global minimum

↳ does not converge on the same sense as batch gradient descent

Week 10

Large scale machine learning - mini batch 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.

b = mini-batch size. $b = 10$ $2 - 100$

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

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

$i := i + 10$

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 $j = 0, \dots, n$)

}

}

can parallelize the "mini-batch" computations
 b examples
1 example

Mini batch will outperform stochastic gradient descent only if you have a good vectorized implementation.

Week 10

Large scale machine learning - stochastic gradient descent convergence

Checking for convergence:

Batch gradient descent:

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

$$J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \quad \text{B = 1000}$$

stochastic gradient descent:

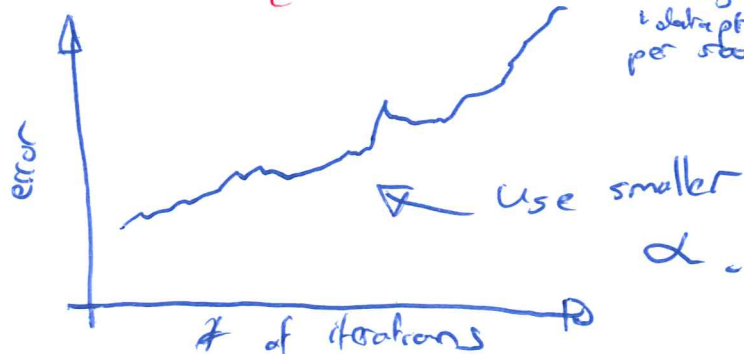
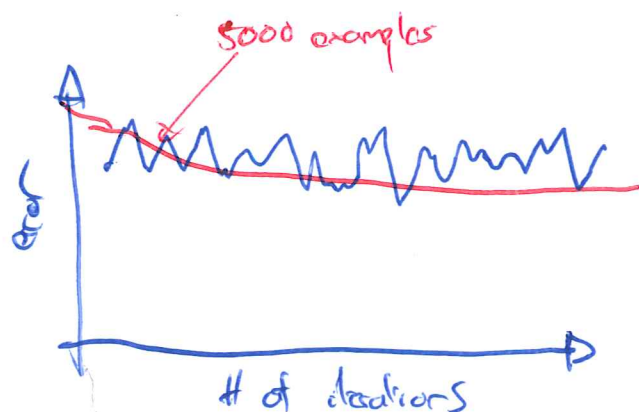
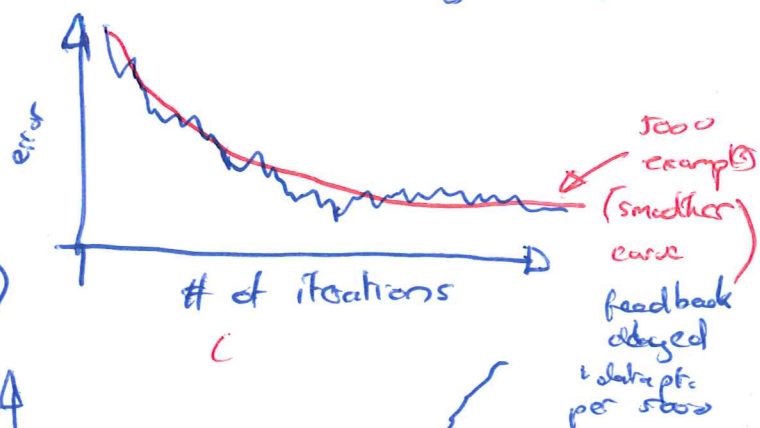
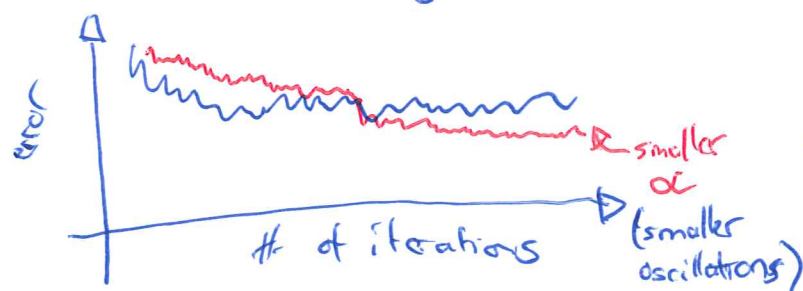
$$\text{cost}(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

during learning, compute $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$ before updating θ using $(x^{(i)}, y^{(i)})$. Right before we update our parameters θ ; compute cost on current example $(x^{(i)}, y^{(i)})$

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

Checking for convergence

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



Week 10

Large scale machine learning - stochastic gradient descent convergence

Learning rate α is typically held constant. Can slowly decrease α over time if we want to converge e.g. $\alpha = \frac{\text{const } 1}{(\text{iteration} + 1) + \text{const } 2}$

$\alpha \rightarrow 0$.

Week 10

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 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 θ using (x, y) :

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

} ~~Disco~~ we are discarding the idea that there is a trading set i.e. only single (x, y)

Can adapt to changes in user preferences

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.

- Learn $p(y=1|x; \theta)$ "problem of learning CTR"
- Use to show user the 10 phones they're most likely to click on.

Other examples: choosing special offer to show user; ~~ask~~ customized selection of news articles; product recommendation

Week 10

Large scale machine learning - Map reduce and data parallelism

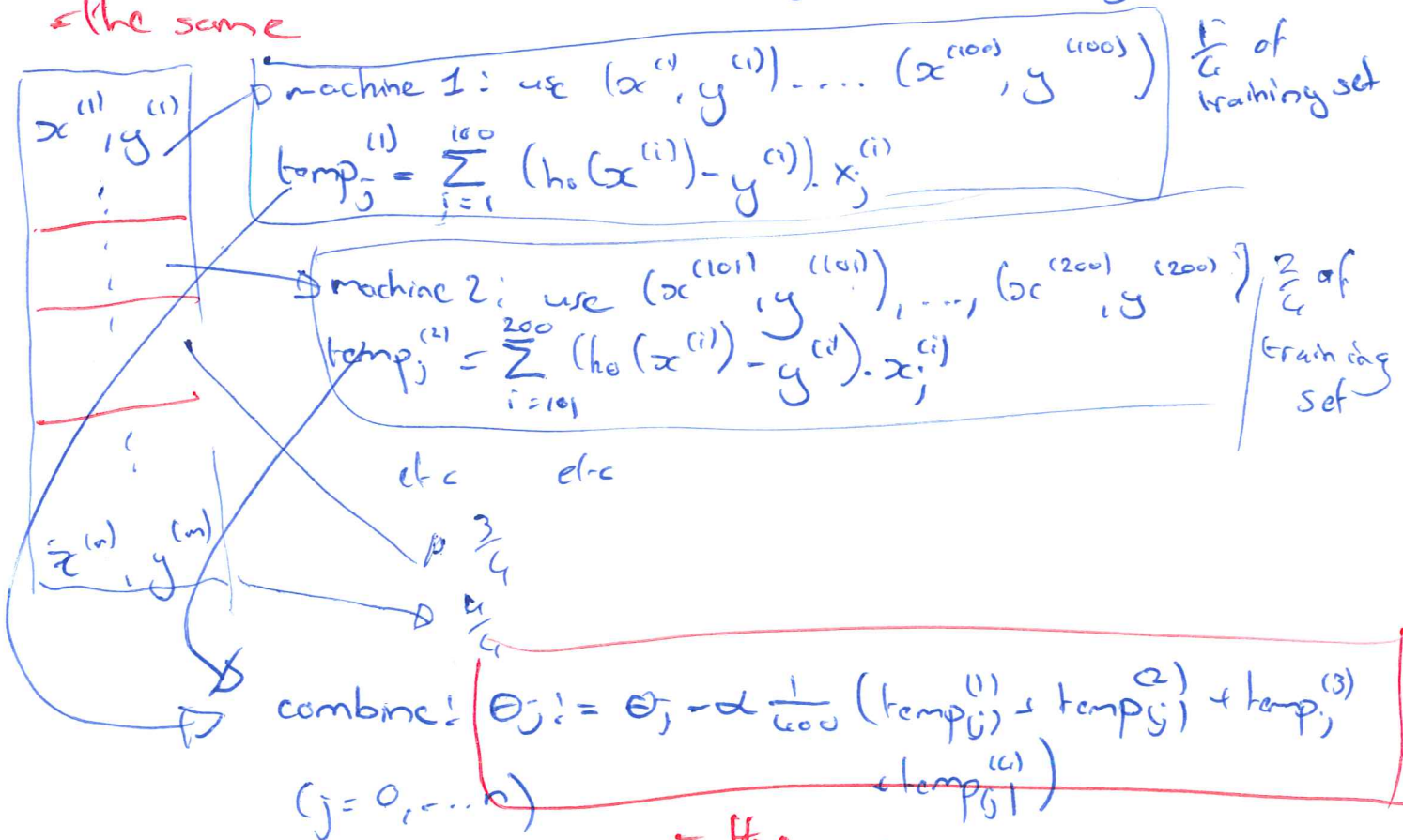
Map reduce

$m = 400,000,000 \rightarrow$ pretend $n = 400$

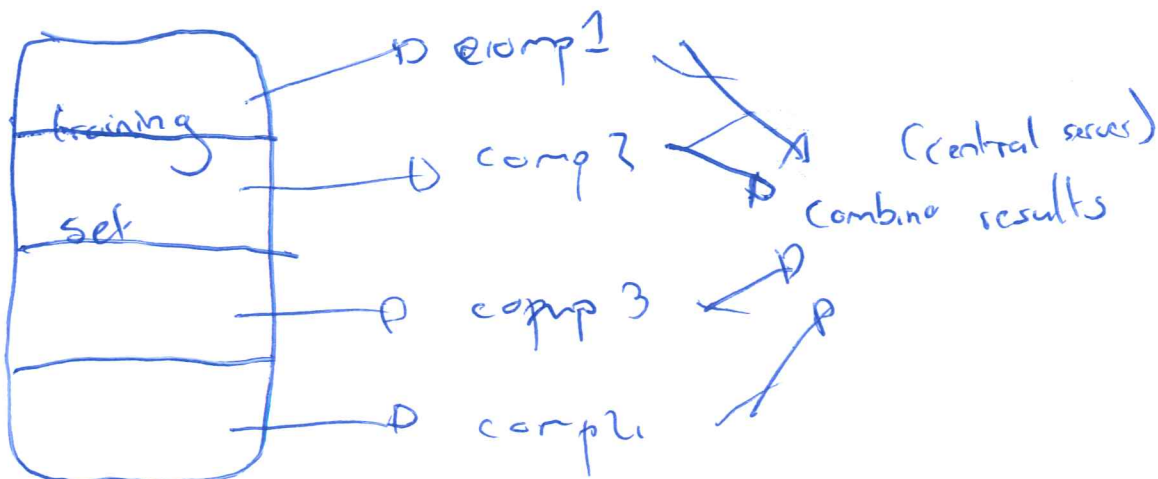
Batch gradient descent $\theta_j := \theta_j - \alpha \frac{1}{400} \sum_{i=1}^{400} (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$

machine 1: use $(x^{(1)}, y^{(1)}), \dots, (x^{(100)}, y^{(100)})$

← the same



Map-reduce



Week 10

Large scale machine learning - map reduce & data parallelism

map-reduce & 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: ^{e.g. train}

Multi-core machines

