

CS229a - Week 10

Written by hackerwins on July 29, 2019

ML

CS229a

Gradient Descent with Large Datasets

Learning with Large Datasets

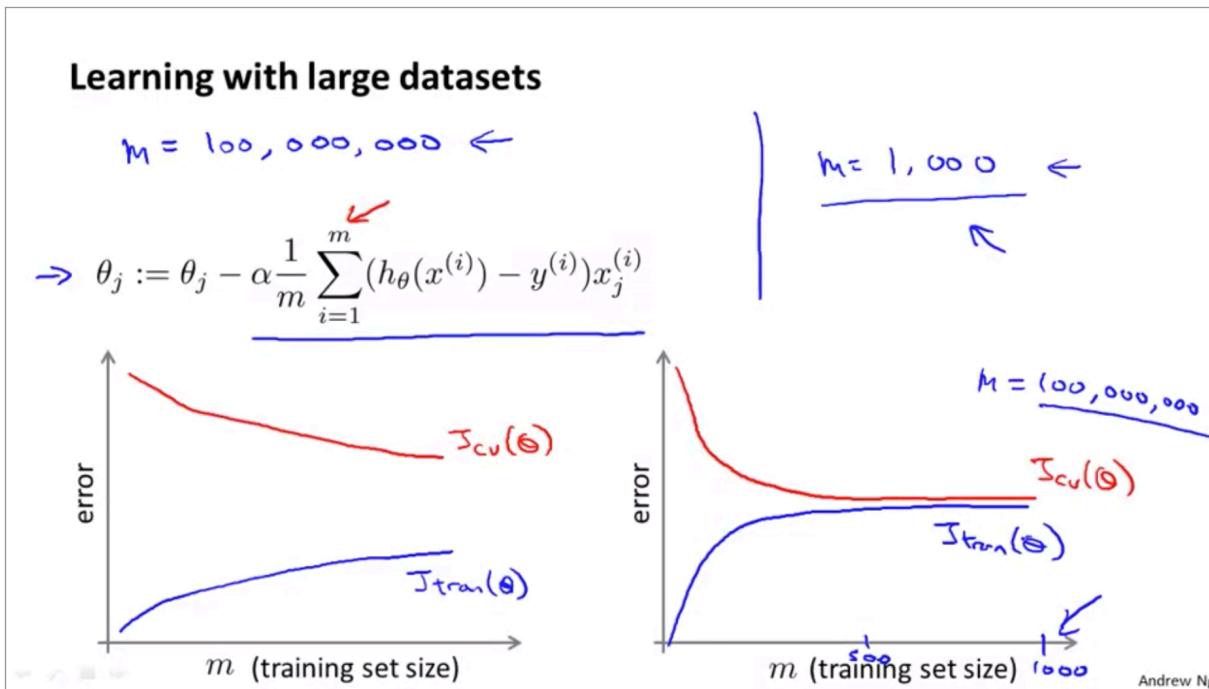
One of the reasons that learning algorithms work so much better now than even say, 5-years ago, is just the amount of data that we have now and that we can train our algorithms on. So why do we want to use such large data sets? We've already seen that one of the best ways to get a high performance machine learning system, is if you take a low-bias learning algorithm, and train that on a lot of data.

But learning with large data sets comes with its own unique problems, specifically, computational problems. Let's say your training set size is M equals 100,000,000.

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

You need to carry out a summation over a hundred million terms, in order to compute these derivatives terms and to perform a single step of decent.

Of course, before we put in the effort into training a model with a hundred million examples, We should also ask ourselves, well, why not use just a thousand examples. 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 $m = 1,000$)? Plot a learning curve for a range of values of m and verify that the algorithm has high variance when m is small.

Stochastic gradient descent

When we have a very large training set, gradient descent becomes a computationally very expensive procedure.

Batch gradient descent

Hypothesis:

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

Cost function:

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

Gradient descent Repeat(for every $j = 0, \dots, n$):

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

If we have hundred million examples ,then computing this derivative term can be very expensive, because the surprise, summing over all m examples. To give the algorithm a name, this particular version of gradient descent is also called **Batch gradient descent**.

Stochastic gradient descent

Cost function:

$$\begin{aligned} cost(\theta, (x^{(i)}, y^{(i)})) &= \frac{1}{2} (h_{\Theta}(x^{(i)}) - y^{(i)})^2 \\ J_{train}(\Theta) &= \frac{1}{2m} \sum_{i=1}^m cost(\theta, (x^{(i)}, y^{(i)})) \end{aligned}$$

Gradient descent:

1. Randomly shuffle training examples
2. Repeat {

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

$$\begin{aligned} \theta_j &:= \theta_j - \alpha (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \\ &\quad (\text{for } j = 0, \dots, n) \end{aligned}$$

}

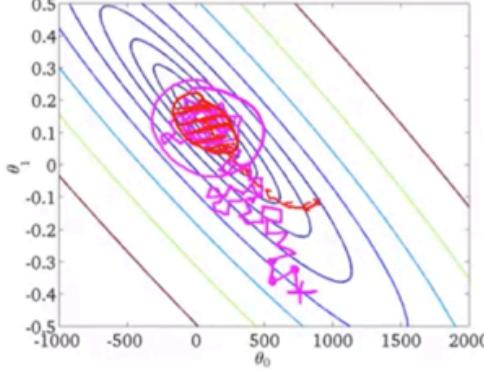
}

Stochastic gradient descent is a lot like descent but rather than wait to sum up these gradient terms over all m training examples, what we're doing is we're taking this gradient term using just one single training example and we're starting to make progress in improving the parameters already.

Stochastic gradient descent

```
→ 1. Randomly shuffle (reorder) training examples
```

```
→ 2. Repeat {
    for i := 1, ..., m{
        →  $\theta_j := \theta_j - \alpha(h_\theta(x^{(i)}) - y^{(i)})x_j^{(i)}$ 
        (for every  $j = 0, \dots, n$ )
    }
}
```



Andrew Ng

As you run Stochastic gradient descent, what you find is that it will generally move the parameters in the direction of the global minimum, but not always. And so take some more random-looking, circuitous path to watch the global minimum. And in fact as you run Stochastic gradient descent it doesn't actually converge in the same sense as Batch gradient descent does and what it ends up doing is wandering around continuously in some region that's in some region close to the global minimum, but it doesn't just get to the global minimum and stay there. But in practice this isn't a problem because, you know, so long as the parameters end up in some region there maybe it is pretty close to the global minimum.

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 (typical range: 2~100)

So we have again data with 300 million training examples, then what we're saying is after looking at just the first 10 examples we can start to make progress in improving the parameters theta so we don't need to scan through the entire training set. We just need to look at the first 10 examples and this will start letting us make progress and then we can look at the second ten examples and modify the parameters a little bit again and so on. So, that is why Mini-batch gradient descent can be faster than batch gradient descent.

So, why do we want to look at b examples at a time rather than look at just a single example at a time as the Stochastic gradient descent? The answer is in vectorization. In particular, Mini-batch gradient descent is likely to outperform Stochastic gradient descent only if you have a good vectorized implementation.

Stochastic Gradient Descent Convergence

Back when we were using batch gradient descent, our standard way for making sure that gradient descent was converging was we would plot the optimization cost function as a function of the number of iterations. So that was the cost function and we would make sure that this cost function is decreasing on every iteration.

Batch gradient descent:

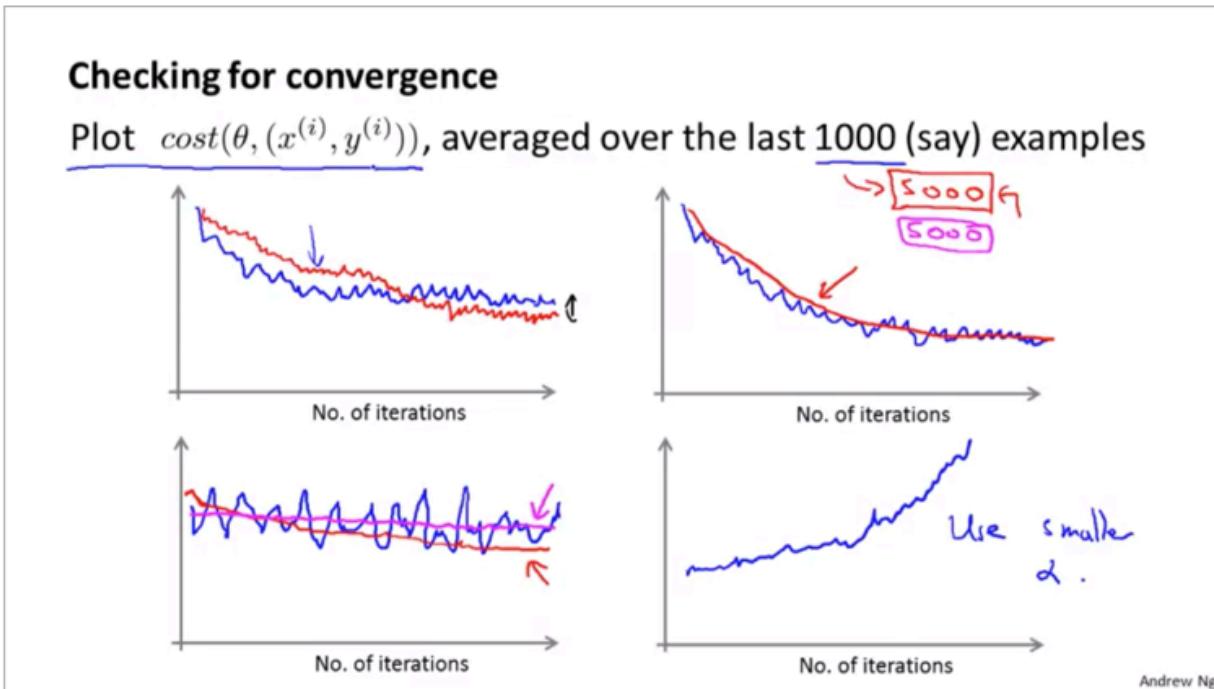
- Plot $J_{train}(\theta)$ as function of the number of iterations of gradient decent

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

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



Typically this process of decreasing alpha slowly is usually not done and keeping the learning rate alpha constant is the more common application of stochastic gradient descent although you will see people use either version.

Stochastic gradient descent

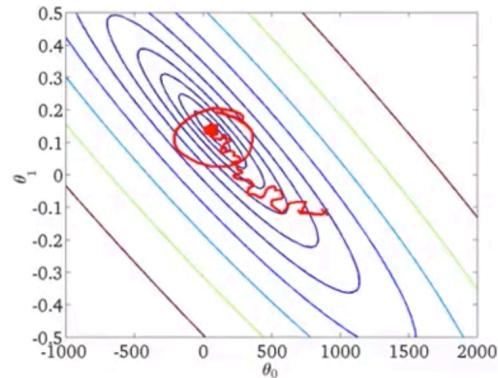
$$\text{cost}(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2}(h_{\theta}(x^{(i)}) - y^{(i)})^2$$

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

1. Randomly shuffle dataset.

2. Repeat {

```
for i := 1, ..., m {
    theta_j := theta_j - alpha * (h_theta(x^{(i)}) - y^{(i)}) * x_j^{(i)}
    (for j = 0, ..., n)
}
```



Learning rate α is typically held constant. Can slowly decrease α over time if we want θ to converge. (E.g. $\alpha = \frac{\text{const1}}{\text{iterationNumber} + \text{const2}}$) $\alpha \rightarrow 0$

Andrew Ng

Advanced Topic

Online Learning

The online learning setting allows us to model problems where we have a continuous flood or a continuous stream of data coming in and we would like an algorithm to learn from that.

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  $\theta$  using  $(x, y)$ .
     $\rightarrow \theta_j := \theta_j - \alpha (h_\theta(x) - y) \cdot x_j \quad (j=0, \dots, n)$ 
}
Can adapt to changing user preference.
  
```

Andrew Ng

Other online learning example:

Product search (learning to search)

User searches for "Android phone 1080p camera" \leftarrow

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

 $(x, y) \leftarrow$

\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; ...

Andrew Ng

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:

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

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

- Machine 2:

$$(x^{(101)}, y^{(101)}), \dots, (x^{(200)}, y^{(200)})$$

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

- Machine 3:

$$(x^{(201)}, y^{(201)}), \dots, (x^{(300)}, y^{(300)})$$

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

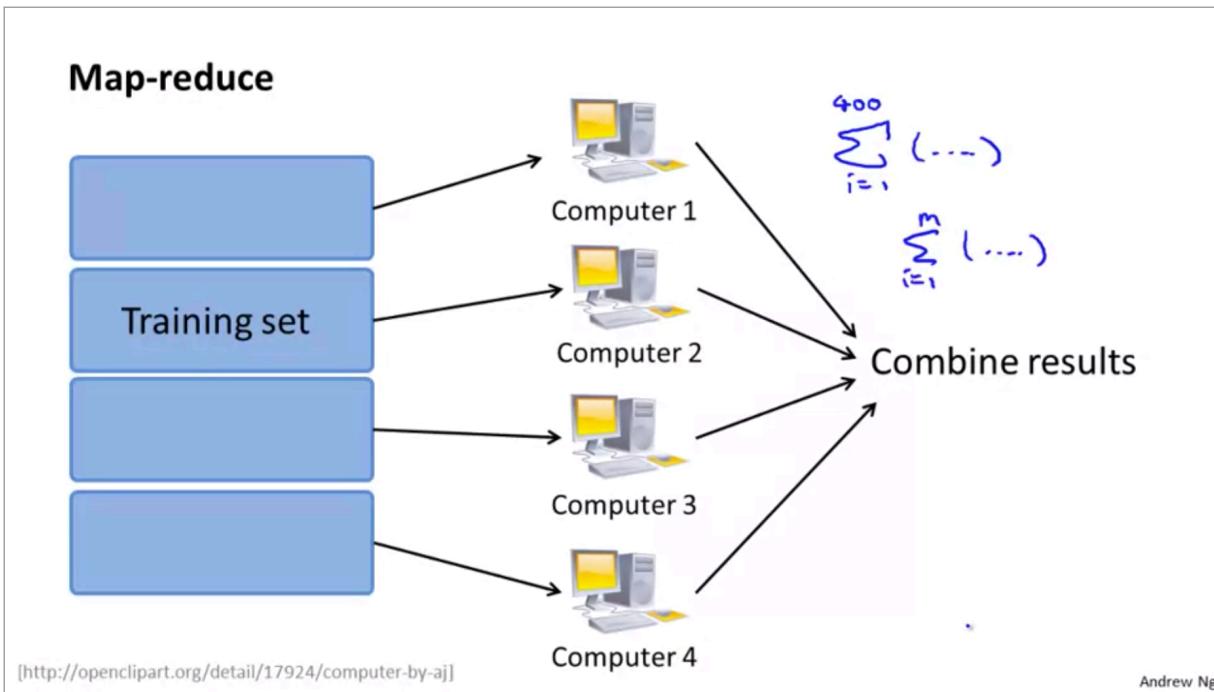
- Machine 4:

$$(x^{(301)}, y^{(301)}), \dots, (x^{(400)}, y^{(400)})$$

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

Combine temp:

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



Map-reduce and summation over the training set

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:

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

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

\nwarrow

$\text{temp}^{(i)}$ $\text{temp}_j^{(i)} \leftarrow$

Andrew Ng

Compute forward propagation and back propagation on 1/n of the data to compute the derivative with respect to that 1/n of the data.

More broadly, by taking other learning algorithms and expressing them in sort of summation form or by expressing them in terms of

computing sums of functions over the training set, you can use the MapReduce technique to parallelize other learning algorithms as well, and scale them to very large training sets.

Sometimes you can just implement your standard learning algorithm in a vectorized fashion and not worry about parallelization and numerical linear algebra libraries could take care of some of it for you.

댓글 0건 [hackerwins blog](#) [Disqus' Privacy Policy](#) 1 로그인 ▾

추천 Tweet 공유 최신순 ▾



토론 시작

[다음으로 로그인](#)

[또는 디스커스에 가입하세요.](#)

이름

1등으로 댓글 달기

구독 당신의 사이트에 Disqus 추가하기 | [Disqus 추가추가](#) [Do Not Sell My Data](#)



Top



Buy me a coffee

© 2020 hackerwins. Made with Jekyll using the [Tale](#) theme.