

Week 10 1
 Gradient Descent with Large Datasets

Learning with Large Datasets

問題設定: 間違いやすい語を classify する, 例 (to, too, two), {then, than}

"It's not who has the best algorithm that wins. It's who has the most data."
 「勝つのはアルゴリズムではなくデータ量だ」

$m = 100,000,000$

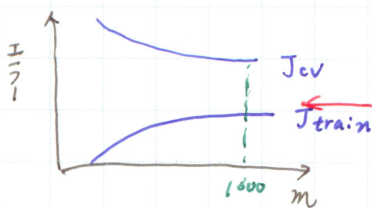
$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

最初から大量の m でトレーニングするのは良くない。たとえば 1000 個ほどのもの、と小さい数のデータをランダムに選んで試してみることは重要。

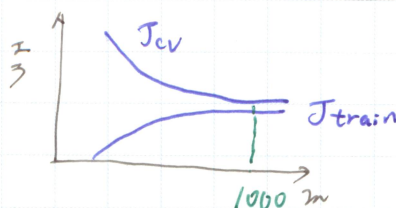
[講義中の Quiz]

supervised learning で m が大きい (1億)。 m が大きいでも全部でトレーニングした方が小数 (たとえば 4) でやるよりも良いと思うのでどうなのかな?

- X ① 試す必要はない。大きいことはいいことだ。
- X ② $J_{train}(\theta)$ を 1000 回ほどを横軸にしてプロットしてみる
- X ③ m を横軸にして learning curve ($J_{train}(\theta)$ と $J_{cv}(\theta)$) をプロットしてみよう。
 m が小さい時の high Bias であることを確認する。 m が小さいと underfit している
- O ④ m の範囲で learning curve をプロットして、 m が小さい時の high variance (overfit) であることを確認する。
 $m \rightarrow$ 大きい overfit を解消できる



overfit
 m を大きくすること
 役に立つ。



underfit
 m を増やしても解けな
 feature を増やせば
 neural network の hidden unit
 数を増やす方がよい。

Stochastic Gradient Descent

大量のデータセットに対しては gradient descent の計算量が問題となる。

Linear regression with gradient descent

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

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

Repeat {

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

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

}

$m = 300,000,000$ の場合 (3億)

Batch gradient descent

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. データセットをランダムにシャッフルする

2. Repeat {

これを10回繰り返す
ループさせる

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

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

(for $j=1, \dots, n$)

}

}

\sum が無いことに注意

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

[講義中の Quiz]

stochastic gradient descent (確率的 ~) にあてはまるものを2選べ

○ ① m がとても大きい時, 普通の gradient descent よりも速い

* ② コスト関数 $J_{\text{train}} = \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$ が batch ~ では

くりかえしの度に実行されるが stochastic ~ では必要ない

毎回くり返すのは $\theta_0 := 0; -\alpha \frac{\partial J}{\partial \theta_0}$ であってコストは毎回計算する必要はない

batch ~ では終了判定のために J の計算が必要。元々 $\frac{\partial J}{\partial \theta_0}$ の計算で $O(m)$ の計算を欠いているので J の計算を省略することはできない

* ③ stochastic ~ は linear regression だけに適用できる

単に偏微分を一部のデータについてだけやろうという話なので他のアルゴリズムでも適用できる

* ④ main loop の 前 に データセットをランダムにシャッフルしておくべき

~~main loop の中で、最初にランダムにシャッフルすべき~~

シャッフルは 1 回だけ行うべき

Week 10

4

Mini-Batch Gradient Descent

Batch gradient descent : 毎回のくり返しで全ての m データを使う

Stochastic ~ : 1個のデータ

Mini-Batch ~ : b 個のデータ
↑ mini-batch size, 例 10~100

Mini-Batch Gradient Descent

1. $b = 10, m = 1000$ とする

2. Repeat {

for $i = 1, \dots, (m-b+1)$ {

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

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

}

}

ベクトル化して並列化できる



[講義中の Quiz] m 個のデータに対して mini-batch ~ を適用する. $b = \text{mini-batch size}$.

これは batch ~ と同じなの？

$b = m$ のとき (答)

Stochastic Gradient Descent Convergence (確率的な最急降下法の収束)

収束している事を確認するためにやる事

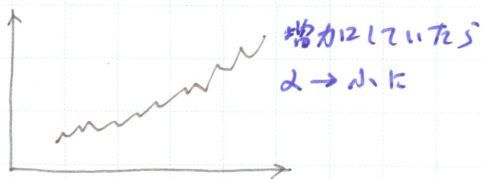
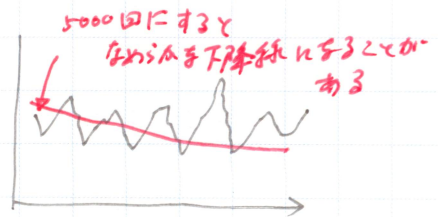
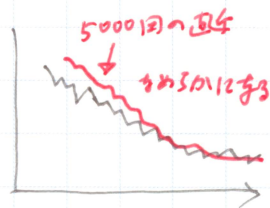
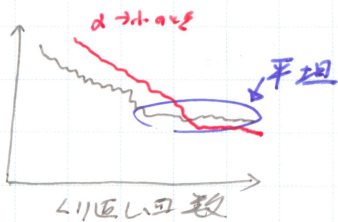
$$\text{cost}(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

学習中に $(x^{(i)}, y^{(i)})$ を用いて θ を更新する前に $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$ を計算する

約1000回のくり返し毎に、直近の1000回のexampleで平均した $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$ をプロットする

Checking for convergence (収束のチェック)

$\text{cost}(\theta, (x^{(i)}, y^{(i)}))$ の直近1000回の平均をプロットする



learning rate (学習率) α は典型的には定数だが、

θ を収束させるために少しずつ小さくする ことのできる

例
$$\alpha = \frac{\text{定数1}}{\text{くり返し回数} + \text{定数2}}$$

パラメータが増えすぎて面倒なので普通はやらない

[講義中の Quiz] stochastic ~ における α は?

- ① とても小さい α を選んでも不利な点はなく、学習が speed up する
× 逆. 基本的には slow down する
- ② α を減らして、十分手回し回数を走らせると、大至 α よりよいパラメータが
えられる ○
- ③ (local minimum) の回りを振動するからうまくのはずなく、収束させなければ、
 α を少しずつ小さくする.
⊖ × 「やりにくいけど、やると、とまってる」
- ④ 1000 回平均の cost (θ , $(x^{(i)}, y^{(i)})$) をプロットして、cost が減っていきながら
 α の値が悪いかも
○

Advanced Topics : Online Learning

連続な stream のデータが 入力され続けていて、そこから学習させる。

例1) shopping website : shipping service を提案 → ユーザは $y = \begin{cases} 1 \\ 0 \end{cases}$ を返す。
price

features x は ユーザの属性, 発送元, 発送先, 料金 を含む。
料金を最適化するために $p(y=1 | x; \theta)$ を学習する。

Repeat for ever {

Get (x, y) corresponding user

Update θ using (x, y)

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

}

}

例2) Product search (商品を探す)

ユーザは "Android phone 1080p camera" を探す。

店には 100 個の phone がある。10 個の検索結果を返す。

feature x は 名前か 検索ワードと何語マッチするか, $y =$
説明か

$$y = \begin{cases} 1 & \text{(ユーザがクリックした) click} \\ 0 & \text{otherwise through} \end{cases}$$

$p(y=1 | x; \theta)$ を学習する。

CTR 問題

Click Through Rate

例3) たくしんの news から、そのユーザに提供する news を 数個 選ぶ。

[講義中の Quiz] online learning algorithm の利点は?

- ① user の嗜好の変化に適合できる $P(y|x;\theta)$ がいつも変化する
- ✗ ② 学習率 α を選ぶ必要がない
- ③ データ・ストリークから学習できる。各 example を1度だけ使って、2度だけ使わない
- ✗ ④ 学習時に良い feature を選択する必要がない。

Advanced Topics: Map-reduce and data parallelism

Map-reduce

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

Machine 1: Use $(x^{(1)}, y^{(1)}), \dots, (x^{(100)}, y^{(100)})$ 最初の 1/4

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

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

$$tmp_j^{(2)} = \sum_{i=101}^{200} \dots$$

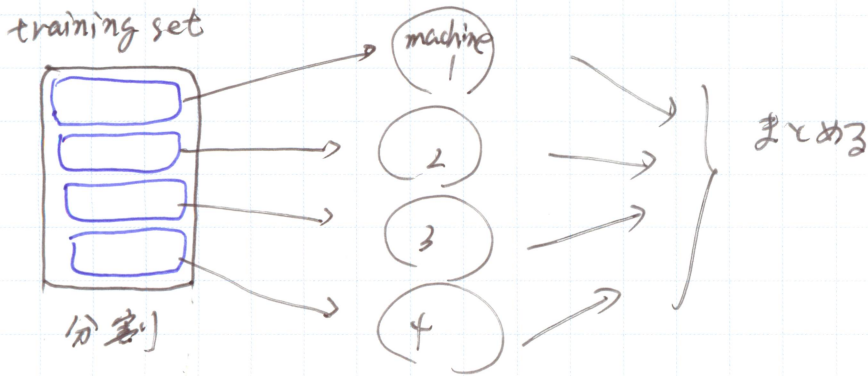
Machine 3: Use $(x^{(201)}, y^{(201)}), \dots, (x^{(300)}, y^{(300)})$

$$tmp_j^{(3)} =$$

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

$$tmp_j^{(4)} =$$

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



Map-reduce と training set に関する和 学習アルゴリズムの多くがこれを使っている
 computing sums of functions over training set

例 logistic regression

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

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

dataset を分割して 複数のマシンの実行する

1台のマシンの CPU が multi-core ならば分割する意味はある

ベクトル化したプログラムを記述しておくと、ライブラリが勝手に並列実行してくれることがある。

[講義中の Quiz] 10台で neural network の学習を行う (map-reduce を使う)。くり返し毎に各マシンで行うのは

- X ① forward と backward の propagation (それぞれ $\frac{1}{10}$ ずつ)
- ~~X~~ ② forward と backward の propagation (それぞれ $\frac{1}{10}$ の θ について) . $\frac{1}{10}$ の θ の微分 ?
- ~~X~~ ③ $\frac{1}{10}$ の θ の forward propagation . 一旦集めた、1台で back propagation を行う
- X ④ $\frac{1}{10}$ の θ の back propagation . 1台のマシンで forward propagation を行い、その後、