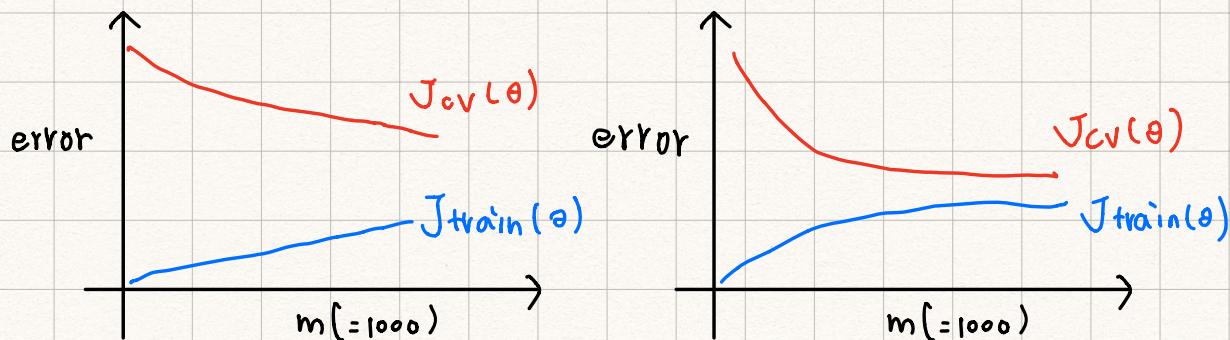


Gradient Descent with Large Datasets

Learning with Large Datasets

$m = 100000000$, 先畫出 $m=1000$ 的 learning curve 跟此圖形觀察加大 data 是否有用?



Machine learning and data

見 week 6 最後的幾頁敘述

Stochastic Gradient Descent (隨機梯度下降)

Stochastic GD

可以應用在任何以 GD 進行訓練的算法中

用來解決 GD 資料量太大所導致的計算量問題

\Rightarrow = sequential mode

Batch GD

vs

Stochastic GD

$J_{\text{train}}(\theta)$

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

Repeat {

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

(for every $j = 0 \cup n$)

}

$$\begin{array}{c} | \\ \text{cost}(\theta, (x^{(i)}, y^{(i)})) \end{array}$$

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

$J_{\text{train}}(\theta)$

1. Randomly shuffle

2. Repeat

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

以上兩種 GD 的收斂方式是不同的

但與向一直線往取最小值前進

靠近局部最小值處徘徊

Stochastic gradient descent

1. Randomly shuffle training examples

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

1~10 次就
收斂

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 通常介於 2 到 100 之間

Mini-batch gradient descent algorithm

Say $b=10$, $m=1000$

↑ 能借由向量化，某些庫
對向量計算有所優化
故計算速度↑

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 $j=0, \dots, n$)

}

Stochastic gradient descent convergence

確保此算法能收斂 and 選擇合適的學習率 α

Checking for convergence

Batch GD:

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

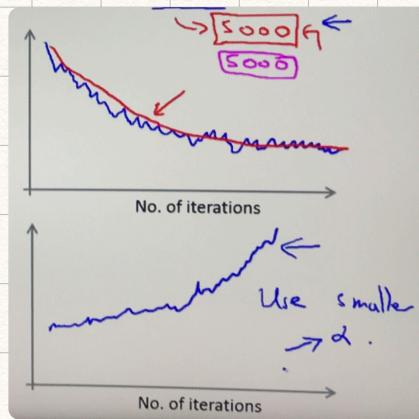
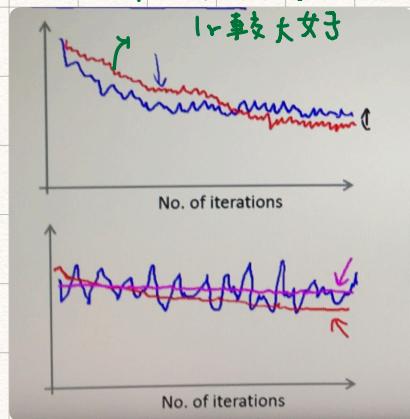
$$J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Stochastic GD:

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

在更新前計算 cost. every 1000 iterations
畫出 $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$ 的平均值

步驟小，結果可能比



Plot cost,
averaged over
the last
1000 examples
5000
⋮

Learning rate of Stochastic GD

Learning rate α is typically held constant.

將 α 設定成隨量化 = 次數修改 α 值

公式: $\alpha = \frac{\text{const 1}}{\text{iteration Num} + \text{const 2}}$

but 還要設置兩個常數，工作量大
較少見人用

Advanced Topics

Online learning

對大批進入的 data 進行學習 (在線學習)

- 特徵 x 是用戶的 properties
 - learn $p(y=1 | x; \theta)$ to optimize price
 T_1 表下單, T_0 表不下單
 - Repeat forever {
 - Get (x, y) corresponding to user.
 - Update θ using $(x, y) \leftarrow \text{data}$ 訓練完便丟棄
 - $\theta_j := \theta_j - \alpha (h_\theta(x) - y) \cdot x_j \quad (j=0, \dots, n)$}
- 此算法可以適應 user 的行為變化

Online learning example :

Product search (learning to search)

User searches for "Android phone 1080p camera"
Have 100 phones in store, will return
10 results.

$\left\{ \begin{array}{l} x = \text{features of phone, how many words in user query match} \\ y=1 \text{ if user clicks on link, } y=0 \text{ otherwise} \end{array} \right.$

Learn $p(y=1 | x; \theta)$. predict CTR
(click through rate)

Use to show user the 10 phones they're most likely to click in.

Other online learning examples:

$\left\{ \begin{array}{l} \text{customized selection of news articles} \\ \text{product recommendation} \\ \text{choosing special offers to show user} \end{array} \right.$

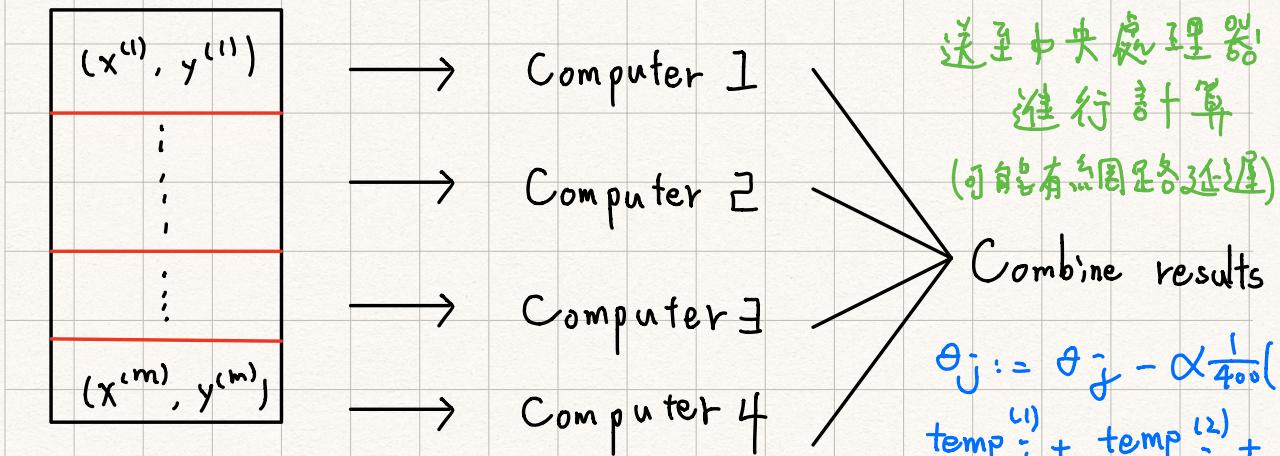
Map reduce and Data Parallelism

數據太大，而不能在一台計算機計算

多台 computer 同計算

Map - reduce

to batch-gradient-descent algorithm



Map-reduce and summation over the training set

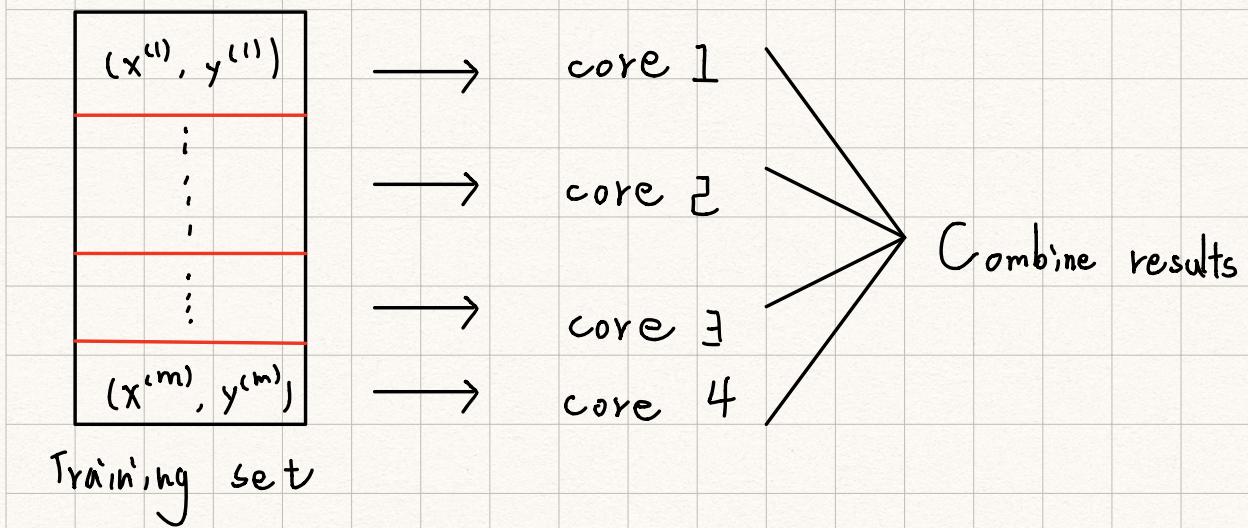
許多算法可以拆成多塊後再進行計算
(logistic regression)

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

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

Multi - core machines

只有一台pc也能做map-reduce (多核CPU)



如果 function 庫能自動實現多核計算
便不需要 map-reduce