

Learning with large datasets

Stochastic Gradient Descent

- Linear regression with gradient descent

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

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

Repeat {

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

$$j = 0, 1, \dots, n$$

}

Batch gradient descent

Computationally expensive since need to load all m examples to do one small step scanning through the entire dataset periodically

- 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

2) Repeat { (usually 1 ~ 10 times repeating, depends on m) }

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

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

$$\text{for } i = 0, 1, \dots, n > \quad = \frac{\partial}{\partial \theta_i} \text{cost}(\theta, (x^{(i)}, y^{(i)}))$$

why stochastic gradient descent is much faster than gradient descent?

- for large dataset, gradient descent has to run through all training set to make one baby step of descent
- Stochastic gradient descent take one descent step only when run through one training example. Especially when dataset is larger, running through all training examples once is enough to get good training result.

o Mini-batch gradient descent

Batch Gradient Descent uses

Stochastic Gradient Descent uses

Mini Gradient Descent uses

$$\begin{bmatrix} \text{all } m \\ 1 \\ b \end{bmatrix}$$

Examples in each iteration

o Mini-batch gradient descent

1) choose b .

2) Repeat {

for $i = 1, 1+b, 1+2b, \dots$ {

$$\theta_j := \theta_j - \alpha \frac{1}{b} \sum_{k=1}^{b-1} (\hat{h}_\theta(x^{(k)}) - y^{(k)}) X_j^{(k)}$$

(for every $j = 0, 1, \dots, n$)

γ

γ

- Mini-batch gradient descent can be faster than stochastic gradient descent

* vectorization!

Mini-gradient descent use b examples for one step, can be vectorized, meaning partially parallelize !!!

Stochastic Gradient Descent Convergence

- Checking for convergence

$$\cdot \text{cost}(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{n} (\hat{y}_\theta(x^{(i)}) - y^{(i)})$$

• Putting the learning compute $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$ before updating using $(x^{(i)}, y^{(i)})$

• Every 100 iterations (say), print cost averaged over the last 100 examples processed by the algorithm.

* Stochastic gradient descent does not converge to the global minimum, but only oscillates around the global minimum

* Still need to choose the proper learning rate α .

* if want SGD to converge, can slowly decrease learning rate α over time

$$\text{e.g. } \alpha = \frac{\text{const 1}}{\text{iteration} + \text{const 2}}$$

Online learning

Eg1: Shipping service website. users sometimes choose your shipping service ($y=1$)

sometimes not ($y=0$)

Features x capture the properties of users, want to learn $P(y=1|x;\theta)$ to optimize the prize to get maximum users

Background: have large stream of user's data

Repeat {

Get (x, y) corresponding to user

update θ using (x, y)

$$\theta_j := \theta_j - \alpha (h_\theta(x) - y) \cdot x_j \quad (j=0, 1, \dots, n)$$

Discard (x, y) after using it once

}

*Advantage of online learning

- Can adapt to change of users' preferences

Eg2: Product Search (learning to search)

User searches for "Android phone 10mp camera"

Have 100 phones in store. Will return 10 results

x : features of phone that matches the query

$y = 1$ if user clicked the button

0 otherwise

learn $P(y=1 | x; \theta)$ [predicted CTR]

CTR: click through rate

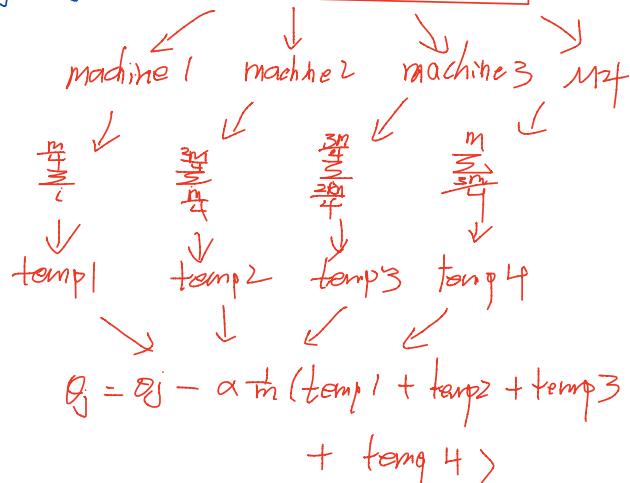
other examples: choosing special offers to users; customized selection of news articles; product recommendation

These examples can also be transferred to classical ML problem

Map-reduce and data parallelism

Map-reduce

$$\text{Batch-gradient descent } \theta_j = \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (\hat{h}_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$



training set



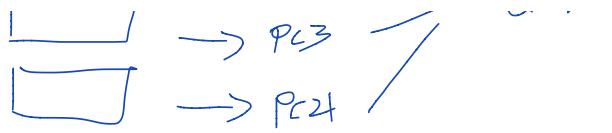
$\rightarrow PC_1$



$\rightarrow PC_2$



Combine results



Many learning algorithms can be expressed as computing sums of functions over the training set

- Multi-moves machines

PC_i → core: