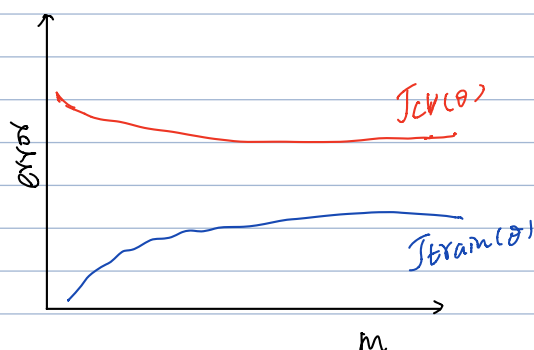
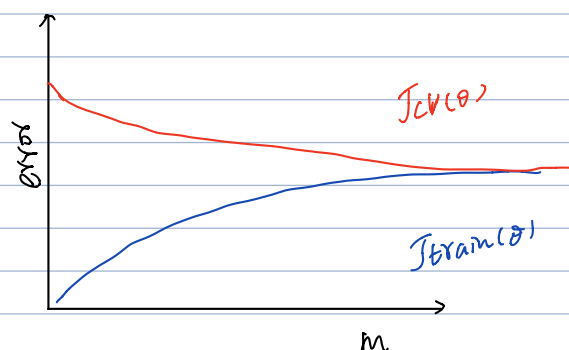


## Learning with large datasets.

In order to verify if the large datasets is likely to perform better than small datasets, we need to plot a learning curve for a range of value of  $m$  and verify that the algorithm has high variance problem when  $m$  is small



high variance  $\rightarrow$  increasing training examples may perform better



high bias  $\rightarrow$  increasing features or NN layers may work better.

## 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.  $\rightarrow$  加速收敛.

2. Repeat {

for  $i=1:m$  {

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

$\rightarrow$  every iteration only lets the algorithm fit  $i^{\text{th}}$  training example better.

for  $j=0:n$  {

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

}

The main problem of SGD is it doesn't always point to the global minimum, and it only wanders around the global minimum in some region, it doesn't get to the global minimum and stay there.

## Mini-batch Gradient Descent.

Batch GD: Use all  $m$  examples in each iteration

Stochastic GD: Use 1 example in each iteration.

Mini-batch GD: Use  $b$  example in each iteration.

$b = \text{batch-size}$ .

Repeat {

for  $i=1:m$  {

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

SGD convergence.

Check for convergence:

Batch GD:

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

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

SGD:

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

During learning, compute  $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$  before updating  $\theta$  using  $(x^{(i)}, y^{(i)})$ .  
Every 1000 iterations, compute  $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$  averaged over the last 1000 examples.

The learning rate can be slowly decreased over time:

$$\alpha = \frac{\text{const } 1}{\text{iterations} + \text{const } 2}$$

Map-reduce and Data Parallelism.

map reduce

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

