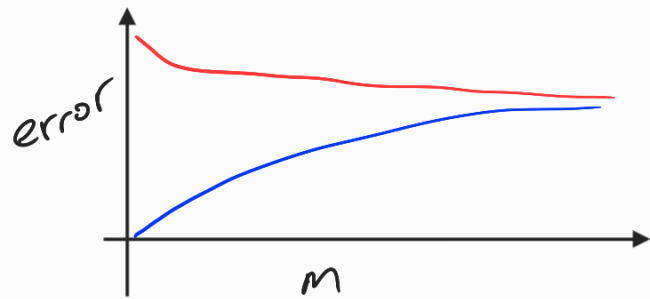
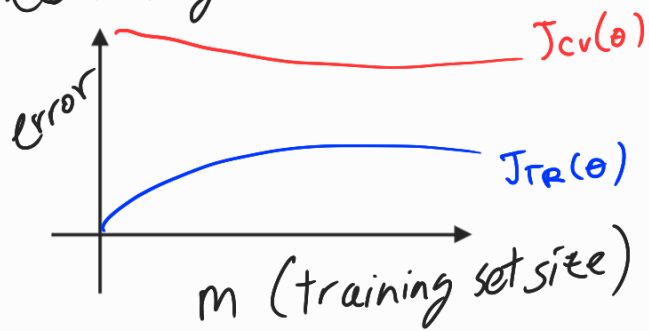


Before going to train with the large dataset, we might pick a small sample and plot the learning curves.



Stochastic gradient descent: 1 example \times iteration

1st: Randomly shuffle dataset

2nd: Repeat {

for ($i = 1 : m$) { for ($j = 0 : n$) {

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

} } }

Convergence:

Compute $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$ before updating θ using $(x^{(i)}, y^{(i)})$

Then, plot cost averaged

Increases \rightarrow smaller α

Noisy \rightarrow increase n° examples to average over

We can slowly decrease α on time $\left(\frac{\text{const } 1}{\text{iter number} + \text{const } 2} \right)$

Mini-batch gradient descent: b examples/iteration

for (i in range($1, n, b$)) {
 (...)

Convergence: plot $J(\theta)$ of the number of iterations.