

mini-batch gradient descent

$$X = \begin{bmatrix} x^{(1)} & \dots & x^{(k)} & | & x^{(i+1)} & \dots & x^{(2k)} & | & \dots & | & x^{(i)} & \dots & x^{(m)} \end{bmatrix} \quad (1)$$

$$Y = \begin{bmatrix} y^{(1)} & \dots & y^{(k)} & | & y^{(i+1)} & \dots & y^{(2k)} & | & \dots & | & y^{(i)} & \dots & y^{(m)} \end{bmatrix} \quad (2)$$

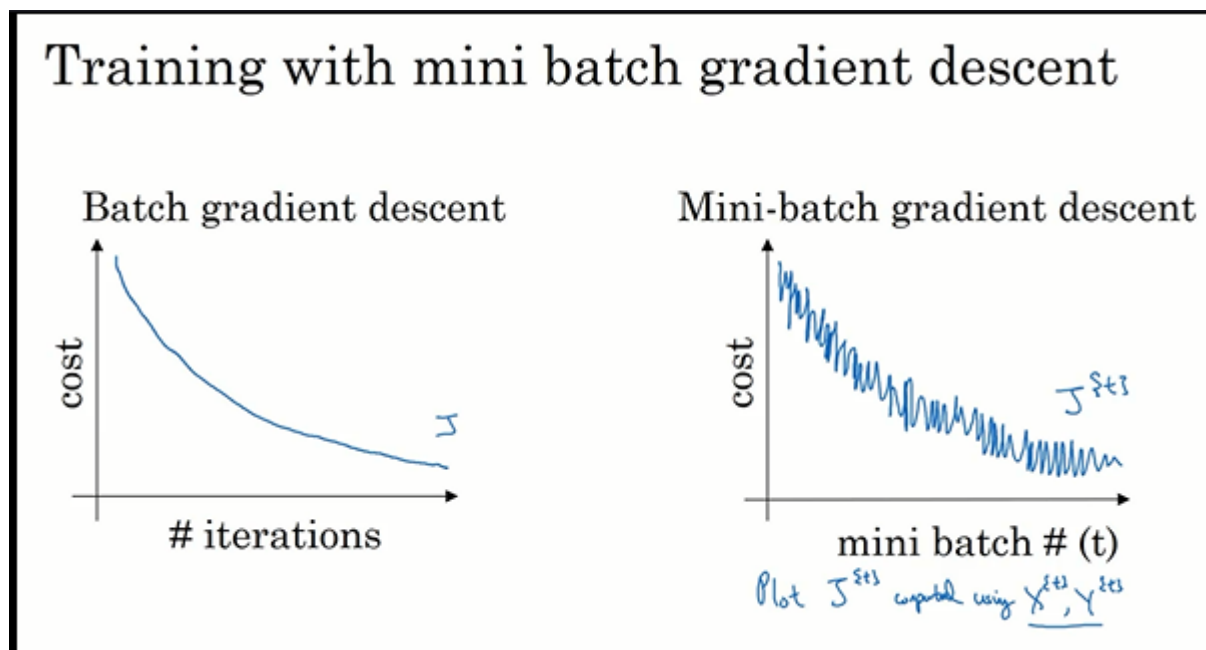
- mini-batches $t : X^{\{t\}}(n_x, \text{mini-batch size}), Y^{\{t\}}(1, \text{mini-batch size})$

$$\underbrace{x^{(1)} \dots x^{(k)}}_{X^{\{1\}}} \dots \underbrace{x^{(i)} \dots x^{(m)}}_{X^{\{j\}}}$$

$$\underbrace{y^{(1)} \dots y^{(k)}}_{Y^{\{1\}}} \dots \underbrace{y^{(i)} \dots y^{(m)}}_{Y^{\{j\}}}$$

- do not need to wait till the whole giant vector finish computing
- epoch**: one pass through the deep network

Training with mini-batch gradient descent



Choosing the mini-batch size

- If minibatch size = m : Batch gradient descent
 - take low-noise and large steps
 - too long per iteration
- If minibatch size = 1 : Stochastic gradient decent
 - oscillating around the optimal point
 - lose speedup from vectorization
- In practice: between 1 and m

- Small training set: use batch gradient descent
- typical mini-batch size: **64, 128, 256, 512, ...**
 - try several values to see if it optimizes the model
- make sure minibatch size fit in CPU\GPU memory