

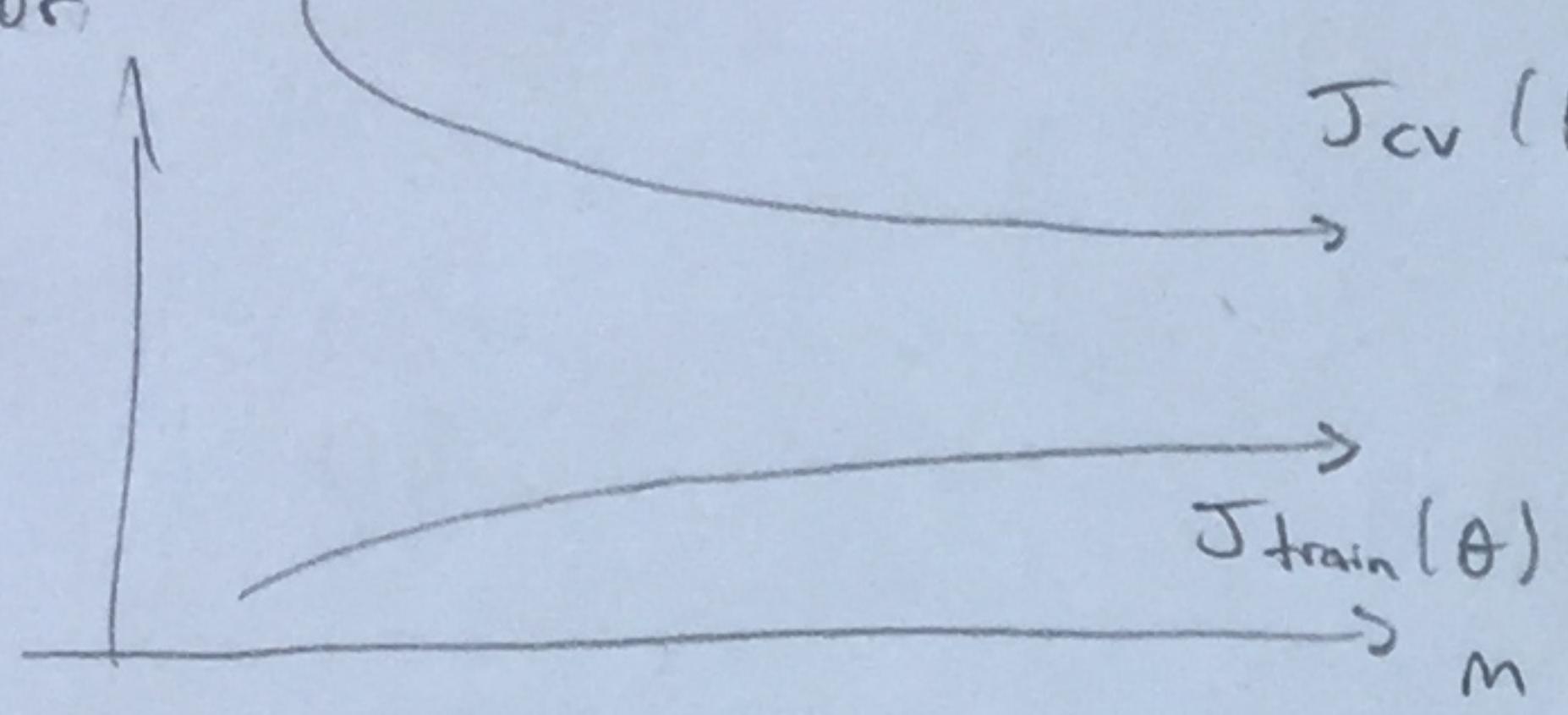
Large Scale Machine Learning:

Large performance improvements w/ more data!

IF  $m = 10^8 \rightarrow \theta_j$  updates w/ gradient descent run  $10^8$  times each step!

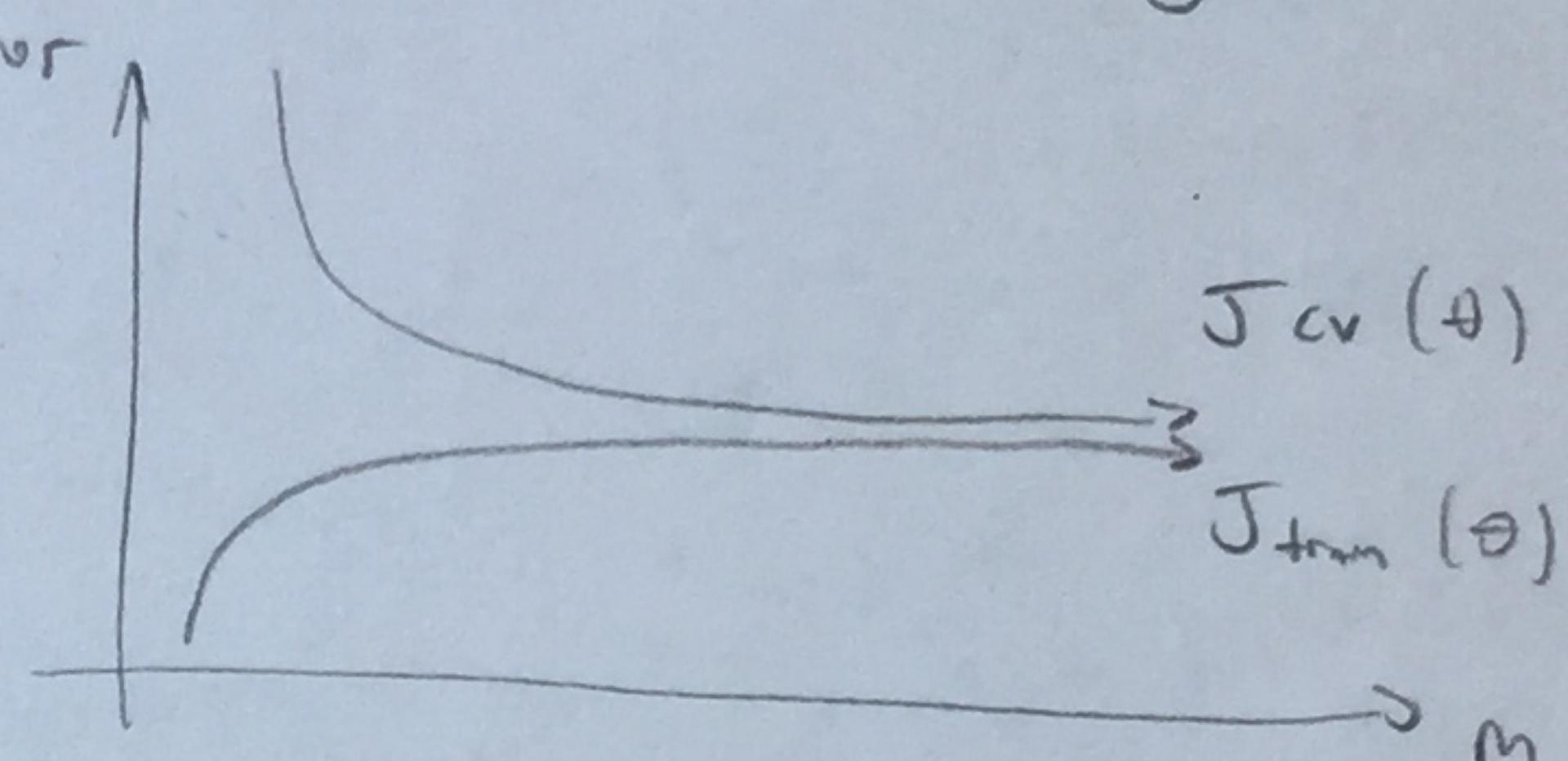
$$(\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)})$$

$\Rightarrow$  recall, check if  $m = 10^8$  better than  $m = 10^3$  w/ learning curve plot:



HIGH VARIANCE (overfit),

$\hookrightarrow$  increase m!

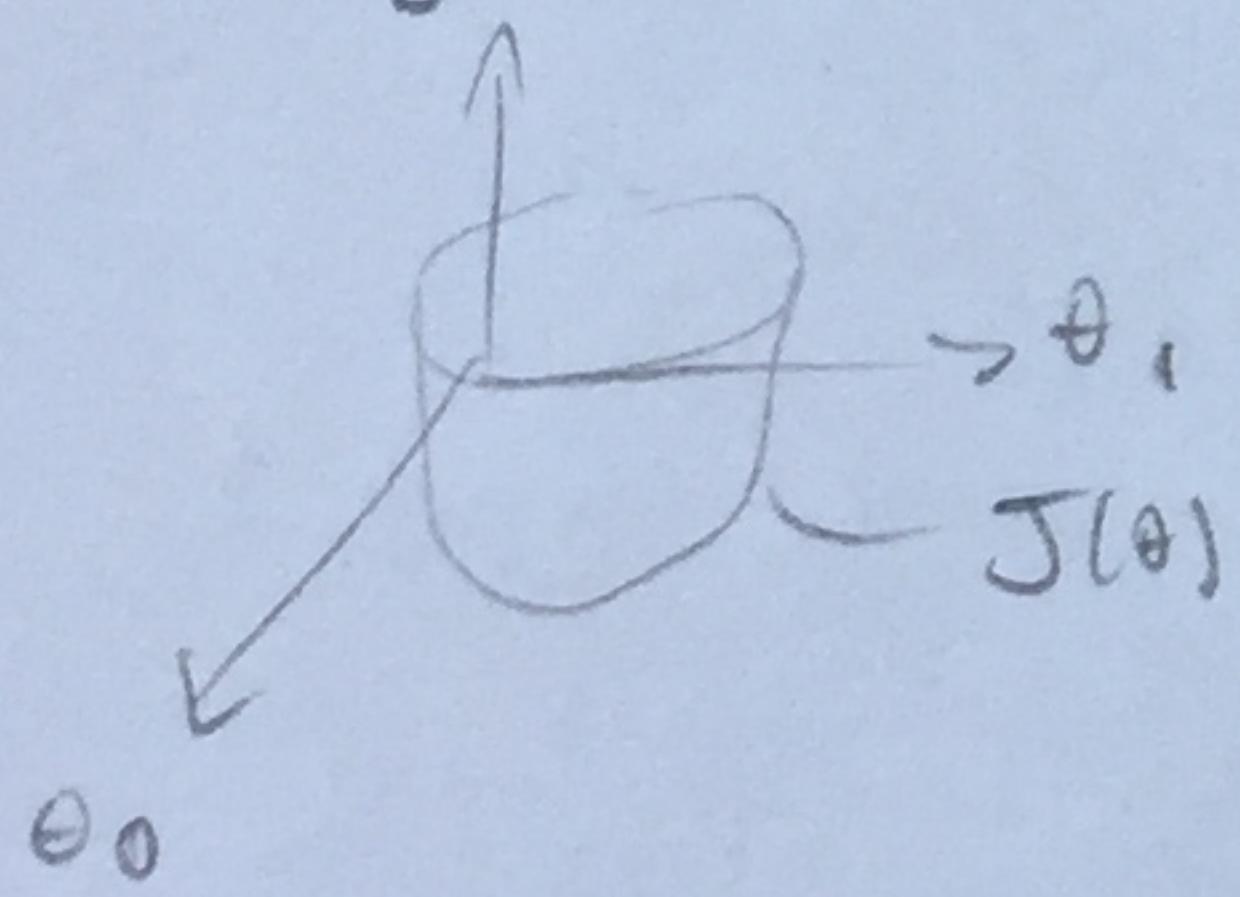


HIGH BIAS

$\hookrightarrow$  increase # features rather than m

Stochastic Gradient Descent:

e.g. Lin Reg Model:  $h_\theta(x) = \sum_{j=0}^n \theta_j x_j$ ,  $J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$



$\Rightarrow$  gradient descent to minimize: (Batch Gradient Descent)

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \text{ for } j = 0 - n$$

$\hookrightarrow$  m large  $\Rightarrow$  computationally

Stochastic Gradient Descent:

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

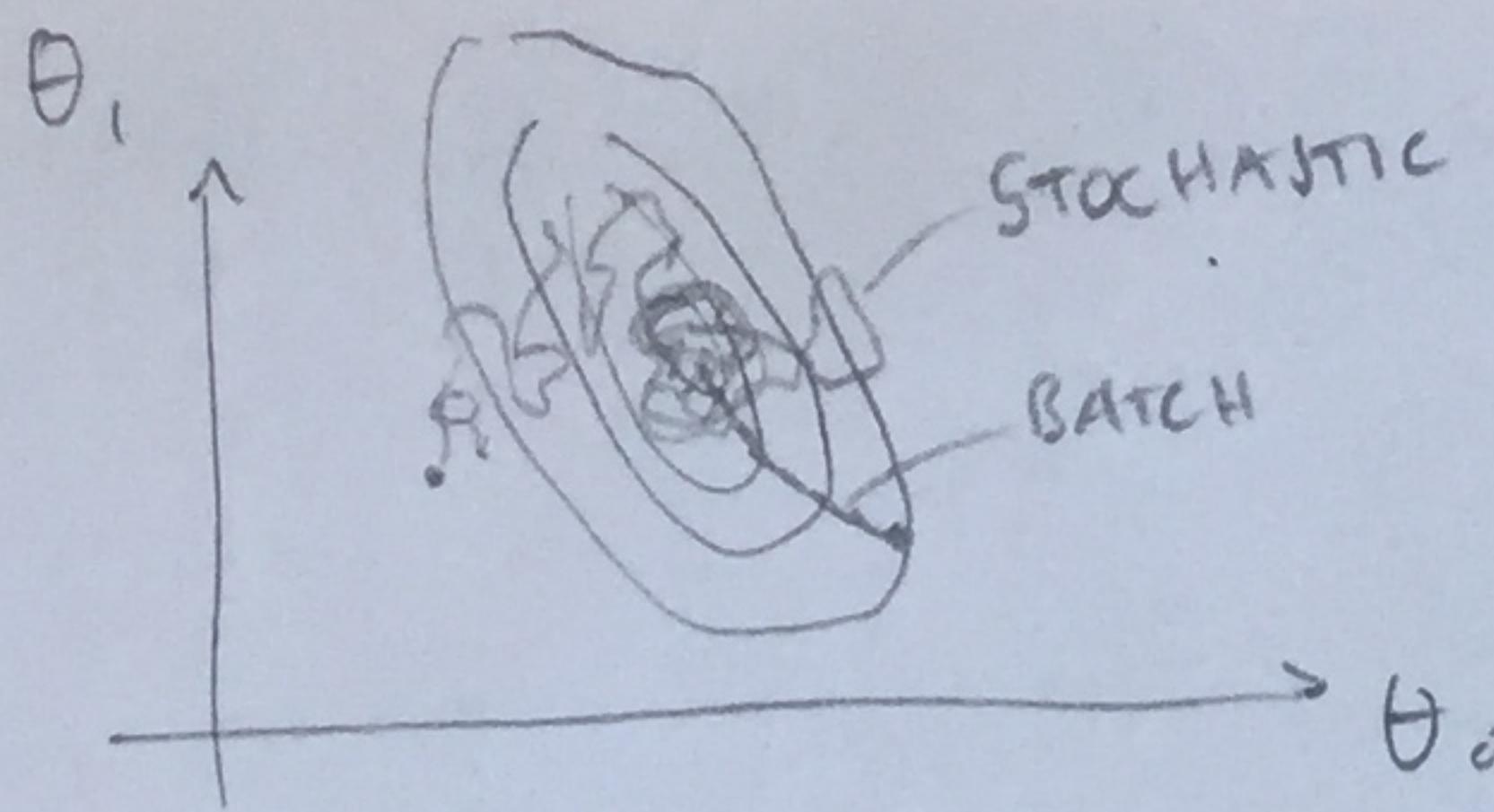
now:  $J_{\text{train}}(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(\theta, (x^{(i)}, y^{(i)}))$

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

① Randomly reorder training dataset

② For  $i = 1 \rightarrow m$ :  $\theta_j := \theta_j - \alpha \frac{(h_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}}{m}$   $\rightarrow$  update gradient based on single example!

$\Rightarrow$  wanders around global minimum  $\rightarrow$  good enough (often)



$\Rightarrow$  Stochastic faster, but more roundabout & doesn't converge

### Mini-Batch Gradient Descent:

Batch: use all  $m$  examples each iteration

Stochastic: 1 example each iteration

Mini-Batch: use  $b$  examples each iteration ( $b = \text{mini batch size}$ )  $\rightarrow$  typically  $b \in [2, 100]$

$$\hookrightarrow \theta_j := \theta_j - \alpha \frac{1}{b} \sum_{k=i}^{i+b-1} (h_\theta(x^{(k)}) - y^{(k)}) x_j^{(k)} \quad \text{for } i = 1, 11, 21 \text{ etc.} \\ (\text{if } b=10)$$

$\Rightarrow$  vectorization w/  $b$  allows faster performance

$\Rightarrow$  must tune  $b$  &  $\alpha$

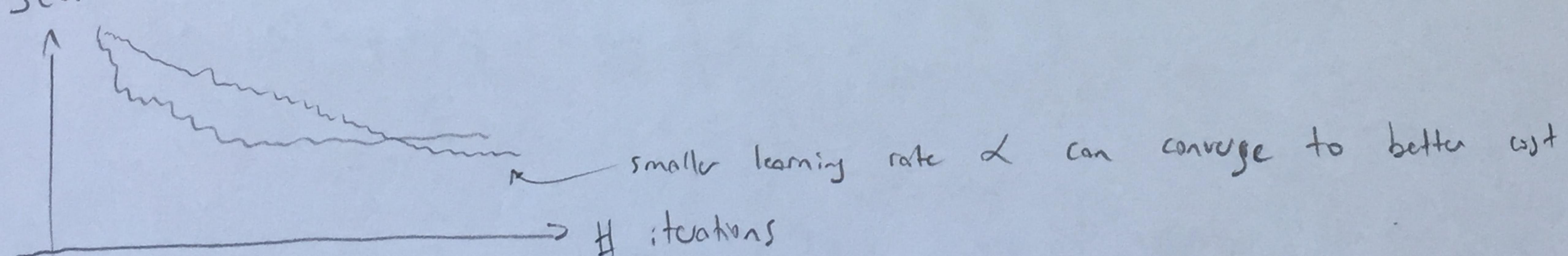
### Convergence:

Batch: plot  $J$  as  $f(\# \text{ iterations}) \rightarrow$  confirm  $J$  decreases each step

Stochastic:  $J$  doesn't decrease each step!

$\hookrightarrow$  compute cost  $(\theta, (x^{(i)}, y^{(i)})$  before updating  $\theta$  each step

then, every  $n$  iterations (eg 1000), plot cost  $(\theta, (x^{(i)}, y^{(i)})$  averaged over last  $n$  processed examples  $\rightarrow$  running estimate of cost



$\hookrightarrow$  increase  $b$  to get smoother lines

$\hookrightarrow$  if looks like diverging, decrease  $\alpha$  (diverging)

$\hookrightarrow$  can set  $\alpha$  to decrease over time:  $\alpha = \frac{\alpha_1}{\text{step \#} + \alpha_2}$

(have to play w/  $\alpha_1, \alpha_2$ )

Online Learning: best for continuous data streams

e.g. Run a shipping service where users visit, then use ( $y=1$ ) or don't use ( $y=0$ ) the service after specifying their package & destination & see a price.

$$\hookrightarrow p(y=1|x;\theta)$$

→ Repeat forever:

Get  $(x, y)$  from user

Update  $\theta$  using  $(x, y)$  ( $\theta_j := \theta_j - \alpha (h_\theta(x) - y) x_j, j = 0 \rightarrow n$ )

→ Can keep data or not (just update w/ new parameters)

↪ adapts to user preferences over time

e.g. 2. Product search → Search for phone, return top 10 results

$x$  = features of phone (# words matching, specs, etc.)

$y=1$  if user clicks link (else,  $y=0$ )

Learn  $p(y=1|x;\theta)$  → Predicted Click-Through Rate (CTR)

Return top 10  $p(y=1)$ . phones

Map Reduce & Data Parallelism:

→ Some algorithms too big to run on one machine

e.g. Batch G.D. w/  $m=400$

↪ machine 1 uses  $(x^{(1)}, y^{(1)}) \rightarrow (x^{(100)}, y^{(100)})$

$m_2$  uses  $101 - 200$

$m_3$  uses  $201 - 300$

$m_4$  uses  $301 - 400$

$$\hookrightarrow \text{temp}_j^{(1)} = \sum_{i=1}^{100} (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

EXACT SAME! ↴

↪ send all  $\text{temp}_j^{(1-4)}$  to master server:  $\theta_j := \theta_j - \frac{1}{400} (\text{temp}_j^{(1)} + \dots + \text{temp}_j^{(4)})$

⇒ anything expressed as sum over training set can use map reduce to speed up algorithm (look for  $\sum_{i=1}^m$ )

⇒ if machine has multiple cores, can use map reduce on one machine (parallelizing)  
↪ no network latency