# ML with large datasets.

for checking if we should use 100 million examples
or not. We plot a learning curve



$J_{cv}(\theta)$

error

$J_{train}(\theta)$

m

(High variance)

Adding examples improve
efficiency

$J_{cv}(\theta)$ 

error

m

(High bias)

Adding examples wont
do anything much.

# Stochastic Gradient Descent

for linear regression with gradient descent,

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

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

Repeat {

$\downarrow$ steps computationally expensive,

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

for every $j = 0, \ldots n$

}

Stochastic gradient descent

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

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

1) Randomly shuffle dataset.

2) Repeat {

　for $i = 1, \ldots, m$

　{

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

　　for $j = (0, \ldots, m)$

　}

}

Mini batch gradient descent

we use $b$ examples in each iteration.

$b$ = mini-batch size

if $b = 10$

↪ $\theta_j := \theta_j - \alpha \frac{1}{10} \sum_{k=i}^{i+9} (h_\theta(x^{(k)}) - y^{(k)}) x_j^{(k)}$
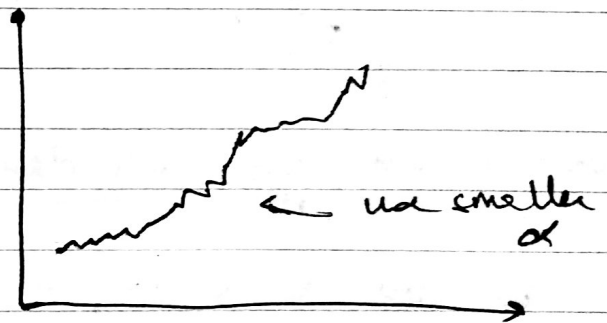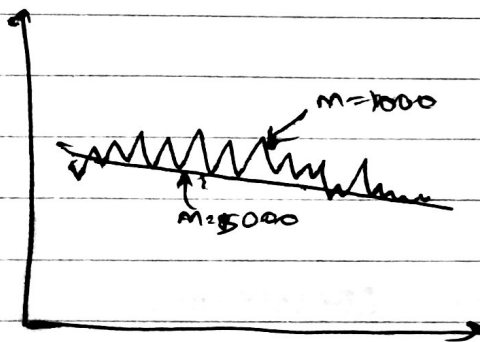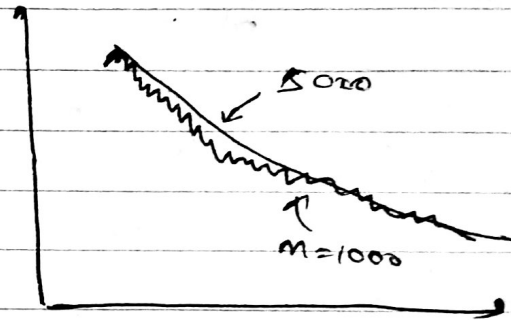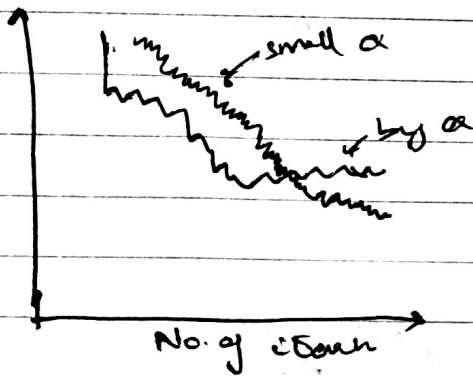
$i = 1 + 10$

Stochestic gradient descent convergence

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

During learning, compute $cost(\theta, (x^{(i)}, y^{(i)}))$ before updating $\theta$ using $(x^{(i)}, y^{(i)})$.

→ Every 1000 iterations plot cost $(\theta, (x^{(i)}, y^{(i)}))$ averaged over the last 1000 examples processed by algorithm.

You might see different trends



Stochastic gradient descent, if you want it to converge to a global minimum, then you could slow decrease $\alpha$ over time.

$$\text{if } \alpha = \frac{Cost1}{iteration\ No + const2}$$

## Online Learning

We may use logistic or neural.
for this, we use ~~this~~ regression.

we want to learn $p(y=1 \mid x; \theta)$

$\uparrow$ price

Repeat forever {
    get $(x, y)$ corresponding to user.
    update $\theta$ using $(x, y)$
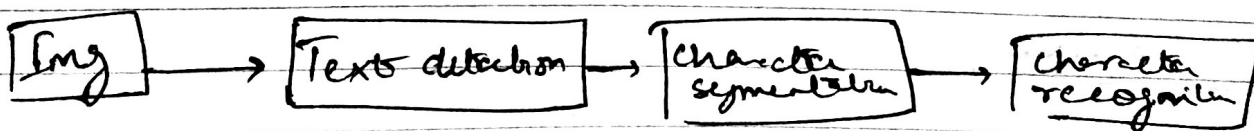    $\rightarrow \theta_j = \theta_j - \alpha(h_\theta(x) - y) \cdot x_j \quad \dots (j = 0 \dots n)$

~~Map Reduce~~

## PHOTO OCR

(Photo optical character recognition).

## Photo OCR pipeline

1) Text detection.
2) character segmentation.
3) character classification
4) check validity $\longrightarrow$ ( Clearing $\approx$ clearing)
        however, it is ignored in this video.

| Img | $\rightarrow$ | Text detection | $\rightarrow$ | Character segmentation | $\rightarrow$ | character recognition |

For pedestrian detection, we use a classifier ( $82 \times 36$ pixel)
    $y = 1$ contains imgs of pedestrian
    $y = 0$ contains random imgs
ok build a neural network or something like that.

# Getting lots of Data

→ You can take a low bias for getting a high performance algo, take a low bias algo and train it on a messier training set.

Artificial data synthesis — Creating new data from scratch.

2 ways

1) use new set.

2). Synthesizing by introducing distortions to existing.

→ Make sure that you have a low bias. — Plot learning curve for it.

→ Increase the no. of hidden units until you have a low bias

## To get More data

- Artificial data synthesis
- collect/ label it yourself.
- "crowd source" ( Eg Amazon Mechanical Turk)

## CEILING ANALYSIS

- Estimate errors due to each component.

# face Recognition