

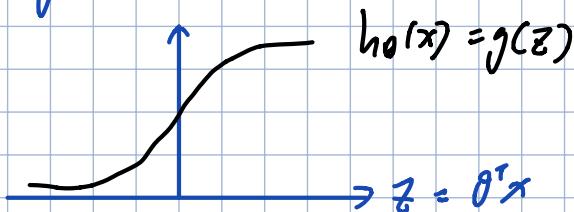
## Week 7

---

### Support Vector Machine (SVM): optimisation objective

an alternative view of logistic regression, recall

$$\rightarrow h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

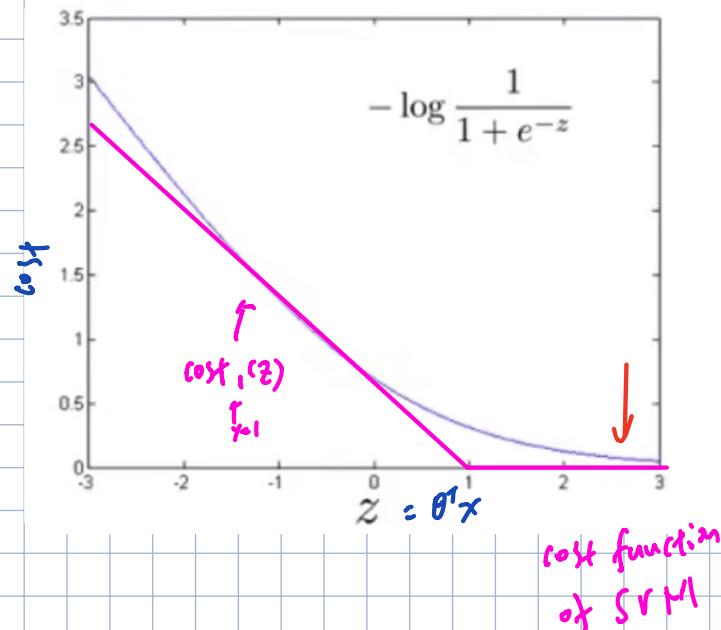


If  $y=1$ , we want  $h_{\theta}(x) \approx 1$ ,  $\theta^T x \gg 0$

If  $y=0$ , we want  $h_{\theta}(x) \approx 0$ ,  $\theta^T x \ll 0$

$$\begin{aligned} \rightarrow \text{Cost of example: } & -y \log(h_{\theta}(x)) + (1-y) \log(1-h_{\theta}(x)) \\ & = -y \log \frac{1}{1 + e^{-\theta^T x}} - (1-y) \log \left(1 - \frac{1}{1 + e^{-\theta^T x}}\right) \end{aligned}$$

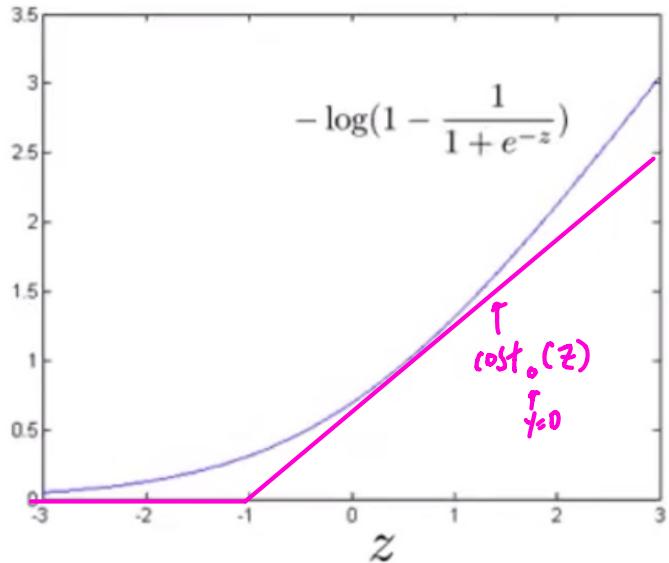
If  $y=1$  (want  $\theta^T x \gg 0$ ):



when logistic regression sees a positive example, ( $y=1$ ), it tries to set  $\theta^T x$  to be very large

→ give us some computational advantage

If  $y=0$  (want  $\theta^T \pi \ll 0$ ):



logistic regression cost function:

$$\min_{\theta} \frac{1}{m} \left[ \sum_{i=1}^m \underbrace{y^{(i)} (-\log h_{\theta}(x^{(i)}))}_{\text{for SVM: } \text{cost}_1(\theta^T \pi^{(i)})} + \underbrace{(1-y^{(i)}) (\log(1-h_{\theta}(x^{(i)})))}_{\text{for SVM: } \text{cost}_0(\theta^T \pi^{(i)})} \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

A                                    B

SVM cost function

$$\min_{\theta} \frac{1}{m} C \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T \pi^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T \pi^{(i)}) + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

convolution

$$\min_u (u-5)^2 + 1 \rightarrow u=5$$

$$\min_u 10 \times (u-5)^2 + 1 \rightarrow u=5$$

if  $\pi$  is big,  $B$  is high  
weights

logistic regression:  $A + \lambda B$   
SVM:  $C A + B$  if  $C$  is big,  $B$  is low  
weight  
 $\sum_i C = Y_n$

SVM cost function

$$\min_{\theta} C \sum_{i=1}^m [y^{(i)} \text{cost}_1(\theta^T \pi^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T \pi^{(i)})] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

unlike logistic regression, SVM doesn't output the probability.

SVM just makes a prediction of  $y$  being equal to 1 or 0 directly.

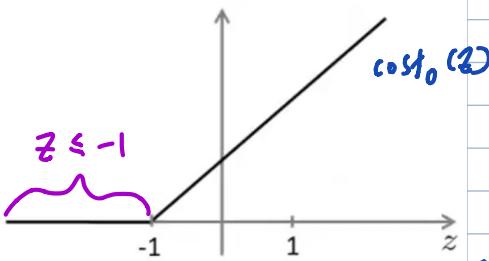
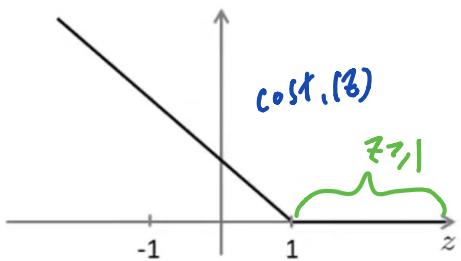
hypothesis:

$$h_{\theta}(x) = \begin{cases} 1 & \text{if } \theta^T x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

## SVM: Large Margin intuition

### Support Vector Machine

$$\min_{\theta} C \sum_{i=1}^m [y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)})] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$



→ If  $y = 1$ , we want  $\theta^T x \geq 1$  (not just  $\geq 0$ )

→ If  $y = 0$ , we want  $\theta^T x \leq -1$  (not just  $< 0$ )

previously  $\theta^T x \geq 0$  barely  
 $\theta^T x < 0$  right

extra safety factor in SVM

## SVM Decision Boundary

Consider  $C = 100,000$  (a very large value)

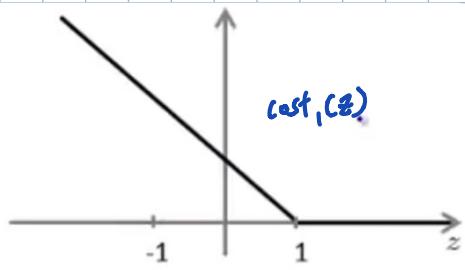
$$\min_{\theta} C \sum_{i=1}^m [y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)})] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

↳ highly motivated = 0  
(if  $C$  is very large)

whenever  $y^{(i)} = 1$ :

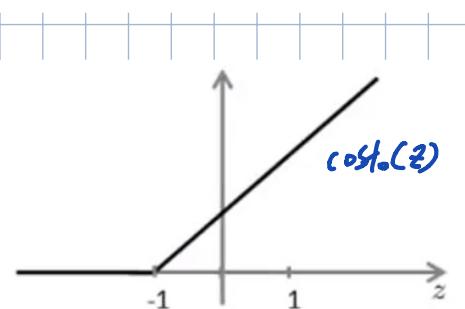
if we want 1st term above  $\approx 0$ , we need

$$\theta^T x^{(i)} \geq 1$$



whenever  $y^{(i)} = 0$ :

$$\theta^T x^{(i)} \leq -1$$



therefore, assuming we try to make first term  $\approx 0$ , what we left is the optimisation problem:

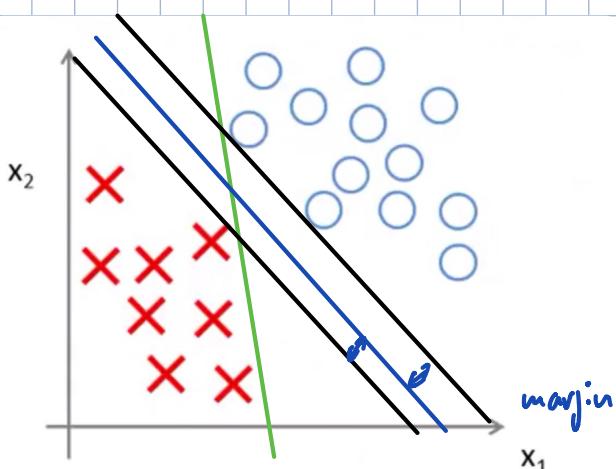
$$\min C \times 0 + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

this will be subject to the constraint that

$$\theta^T x^{(i)} \geq 1 \quad \text{if } y^{(i)} = 1$$

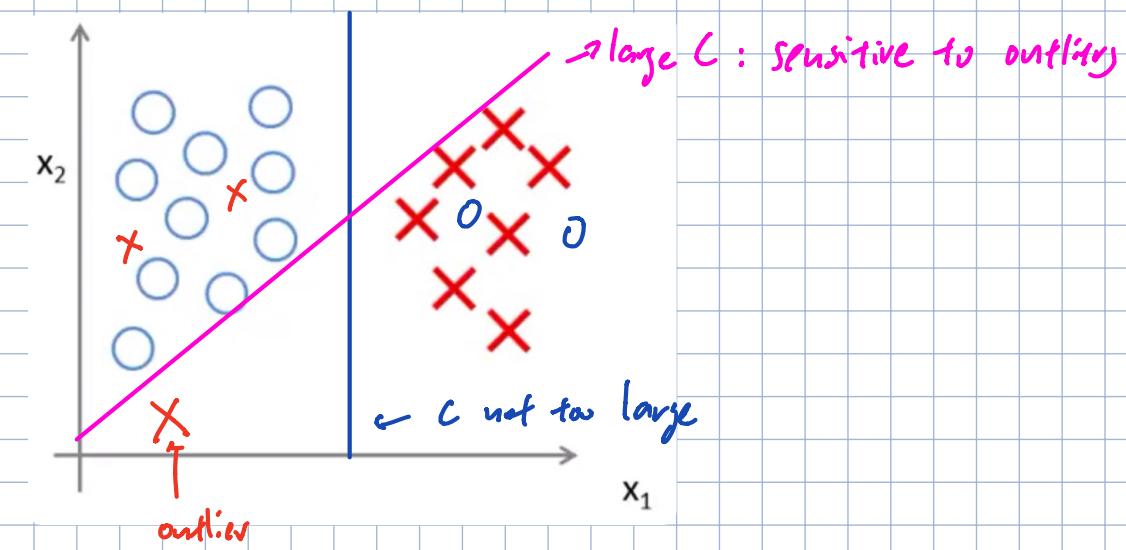
$$\theta^T x^{(i)} \leq -1 \quad \text{if } y^{(i)} = 0$$

### SVM. Linearly separable case



$\Rightarrow$  large margin classifier

# Large margin classifier in presence of outliers

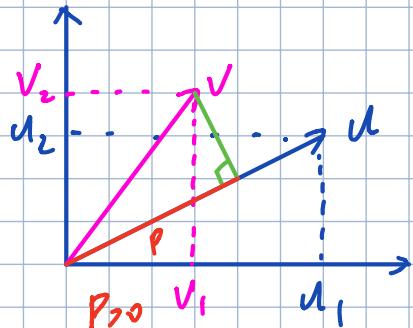


## SVM: the maths

### Vector Inner Product

$$u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \quad v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

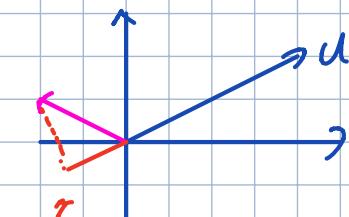
$$\|u\| = \text{length/normal of vector } u \\ = \sqrt{u_1^2 + u_2^2} \in \mathbb{R}$$



$\rightarrow p = \text{length of projection of } v \text{ onto } u = \|\mathbf{p}\|$

$$u^T v = p \cdot \|u\| = v^T u$$

$$= u_1 v_1 + u_2 v_2$$



$$u^T v = p \cdot \|u\|$$

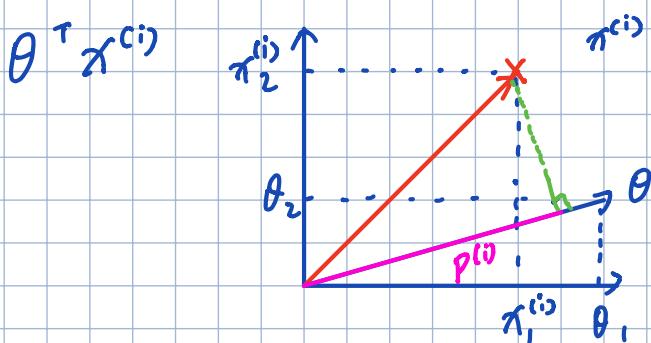
$p < 0 \quad \text{angle between } u \text{ and } v \geq 90^\circ$

## SVM Decision Boundary

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} (\theta_1^2 + \theta_2^2) = \frac{1}{2} \left( \sqrt{\theta_1^2 + \theta_2^2} \right)^2 = \frac{1}{2} \|\theta\|^2$$

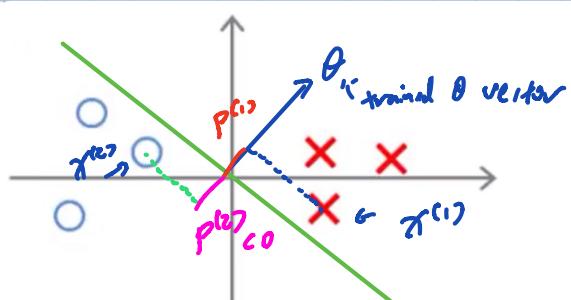
subject to  $\theta^T x^{(i)} \geq 1 \quad \text{if } y^{(i)} = 1$   
 $\theta^T x^{(i)} \leq -1 \quad \text{if } y^{(i)} = 0$

$\Rightarrow$  simplification here:  $\theta_0 = 0$ ,  $n=2$  (for now)



$$\begin{aligned} \theta^T x^{(i)} &= p^{(i)} \cdot \|\theta\| \\ &= \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} \end{aligned}$$

$\theta_0 = 0$  means that decision boundary pass through  $(0,0)$



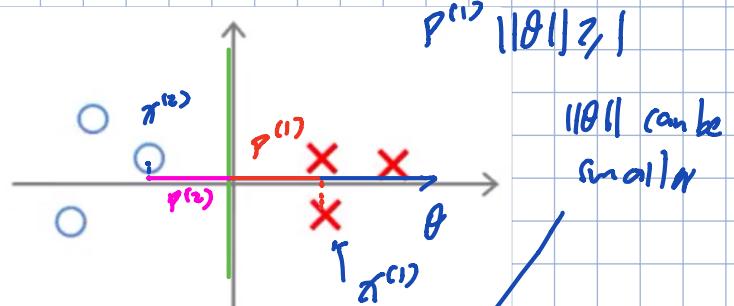
$p^{(i)}$  is small number

but we want  $p^{(i)} \cdot \|\theta\| \geq 1$  for  $y^{(i)} = 1$

if  $p^{(i)}$  is small, then  $\|\theta\|$  has to be large

same for

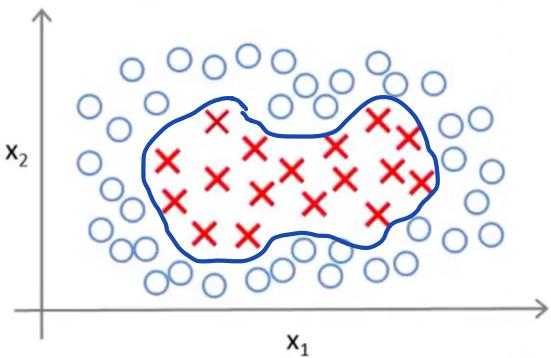
$p^{(i)} \cdot \|\theta\| \leq -1$  for  $y^{(i)} = 0$



objective f  
 $\min_{\theta}$

## SVM: Kernels

### Non-linear Decision Boundary



Predict  $y = 1$  if

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 + \theta_5 x_2^2 + \dots \geq 0$$

$$h_{\theta}(x) = \begin{cases} 1 & \text{if } \theta_0 + \theta_1 x_1 + \dots \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

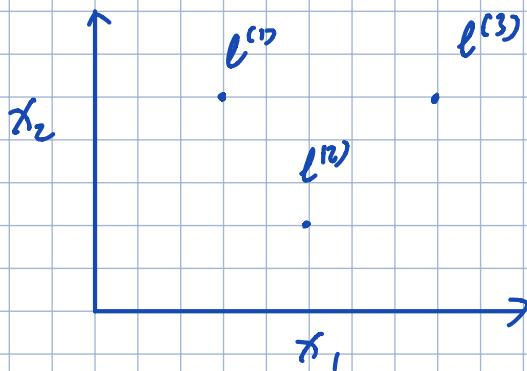
we define  $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \dots$

where  $f_1 = x_1, f_2 = x_2, f_3 = x_1 x_2, f_4 = x_1^2, f_5 = x_2^2$

⇒ is there a different/better choice of the features  $f_1, f_2, \dots$ ?

↙  
computationally expensive

### Kernel



Given  $x$ , compute new feature depending on proximity to landmarks  $l^{(1)}, l^{(2)}, l^{(3)}$

$$\text{Given } x: \text{ Define } f_1 = \text{similarity}(x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$$

$$f_2 = \text{similarity}(x, l^{(2)}) = \exp\left(-\frac{\|x - l^{(2)}\|^2}{2\sigma^2}\right)$$

$$f_3 = \text{similarity}(x, l^{(3)}) = \exp\left(-\frac{\|x - l^{(3)}\|^2}{2\sigma^2}\right)$$

Kernel (Gaussian)

$$\Downarrow K(x, l^{(i)})$$

## Kernels and Similarity

$$f_1 = \text{similarity}(\pi, \ell^{(1)}) = \exp\left(-\frac{\|\pi - \ell^{(1)}\|^2}{2\sigma^2}\right)$$

→ If  $\pi \approx \ell^{(1)}$ :  $f_1 \approx \exp\left(-\frac{0^2}{2\sigma^2}\right) \approx 1$

→ If  $\pi$  is far from  $\ell^{(1)}$ :

$$f_1 = \exp\left(-\frac{(\text{large number})^2}{2\sigma^2}\right) \approx 0$$

Each landmark defines a new feature (Given a training example  $\ell^{(i)}$ )

$$\ell^{(1)} \rightarrow f_1$$

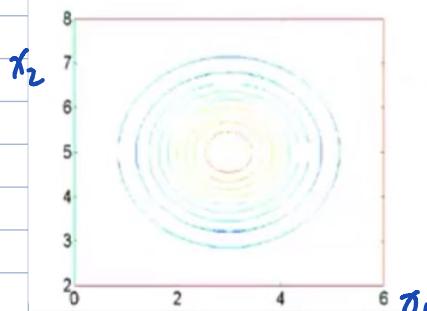
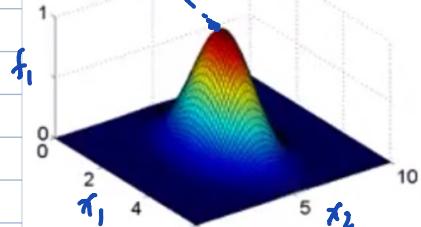
$$\ell^{(2)} \rightarrow f_2$$

$$\ell^{(3)} \rightarrow f_3$$

e.g.  $\ell^{(1)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$ ,  $f_1 = \exp\left(-\frac{\|\pi - \ell^{(1)}\|^2}{2\sigma^2}\right)$

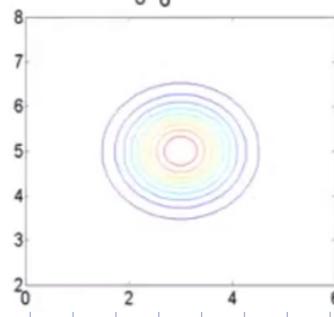
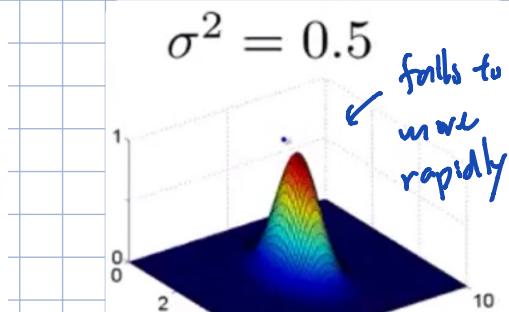
when  $\pi = \begin{bmatrix} 3 \\ 5 \end{bmatrix} \Rightarrow f_1 = 1$

$$\sigma^2 = 1$$

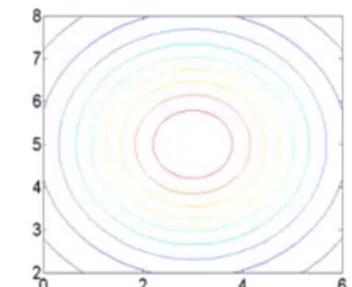
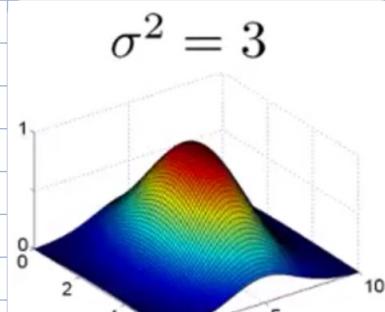


$$\sigma^2 = 0.5$$

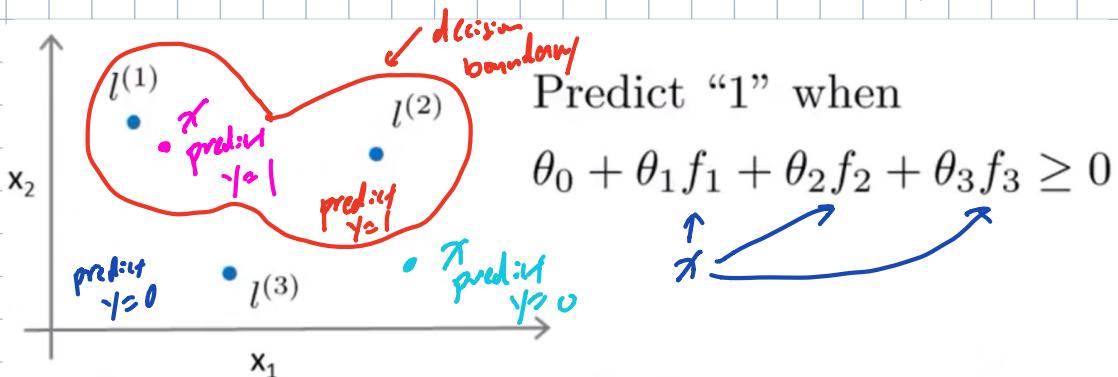
falls to 0  
more rapidly



$$\sigma^2 = 3$$



Given the training example  $X$ , we are going to compute  $f_1, f_2, f_3$



Predict "1" when

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$$

Say after training we get  $\theta_0 = -0.5, \theta_1 = 1, \theta_2 = 1, \theta_3 = 0$

$$\begin{aligned} & \theta_0 + \theta_1 \times 1 + \theta_2 \times 0 + \theta_3 \times 0 \\ &= -0.5 + 1 = 0.5 > 0 \end{aligned}$$

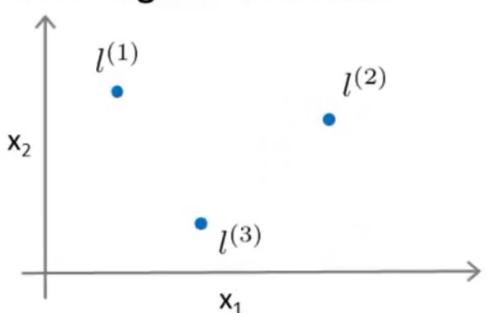
$$f_1 \approx 1, f_2 \approx 0, f_3 \approx 0$$

$$f_1, f_2, f_3 \approx 0 \Rightarrow \theta_0 + \theta_1 f_1 + \dots \approx \theta_0 = -0.5 < 0$$

Here, if  $x$  is near  $l^{(1)}$  and  $l^{(2)}$   $\Rightarrow$  predict  $y=1$  ( $\because \theta_3=0$ )  
is far from ..  $\Rightarrow$  predict  $y=0$

## Choosing landmarks

### Choosing the landmarks



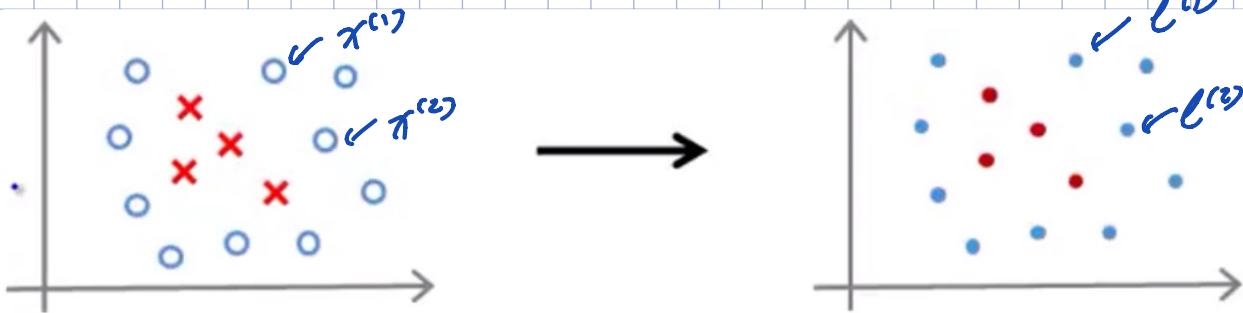
Given  $x$ :

$$\begin{aligned} f_i &= \text{similarity}(x, l^{(i)}) \\ &= \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right) \end{aligned}$$

Predict  $y = 1$  if  $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$

Where to get  $l^{(1)}, l^{(2)}, l^{(3)}, \dots$ ?

In practice, for every training example we have, we're putting landmarks at the same locations as the training examples.



↳  $\ell^{(i)} \Rightarrow \text{SVM with kernels}$

: → Given  $(x^{(1)}, y^{(1)})$ ,  $(x^{(2)}, y^{(2)})$ , ...,  $(x^{(m)}, y^{(m)})$ ,  
 $\ell^{(m)}$  choose  $\ell^{(1)} = x^{(1)}$ ,  $\ell^{(2)} = x^{(2)}$ , ...,  $\ell^{(m)} = x^{(m)}$

→ Given example  $x$ , ( $x$  can be in training/test/CV set)

$$f_1 = \text{similarity}(x, \ell^{(1)}) \quad x^{(1)}$$

$$f_2 = \text{similarity}(x, \ell^{(2)})$$

$$\dots$$

$$\left. \begin{array}{l} f = \begin{bmatrix} f_0 \\ f_1 \\ \vdots \\ f_m \end{bmatrix} \\ f_0 = 1 \end{array} \right\} \text{feature vector}$$

e.g. for training example  $(x^{(i)}, y^{(i)})$ :

$$x^{(i)} \xrightarrow{\text{map}} f_i^{(i)} = \text{sim}(x^{(i)}, \ell^{(1)})$$

$$\in \mathbb{R}^{n+1} \quad f_2^{(i)} = \text{sim}(x^{(i)}, \ell^{(2)}) \quad x^{(i)}$$

$$\vdots \quad f_i^{(i)} = \text{sim}(x^{(i)}, \ell^{(i)}) = 1 \quad (\text{for Gaussian Kernel})$$

$$f_m^{(i)} = \text{sim}(x^{(i)}, \ell^{(m)})$$

$$\downarrow \text{map}$$

$$f^{(i)} = \begin{bmatrix} f_0^{(i)} \\ f_1^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix} \in \mathbb{R}^{m+1}$$

## SVM with Kernels

Hypothesis: Given  $\pi$ , compute features  $f \in \mathbb{R}^{m+1}$

predict " $y=1$ " if  $\theta^T f \geq 0$  (i.e.  $\theta_0 + \theta_1 f_1 + \dots + \theta_m f_m \geq 0$ )  
(assume we already have a learning set of  $\theta$ )

$$\theta \in \mathbb{R}^{m+1}$$

training for  $\theta$ :

$$\min_{\theta} C \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T f^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T f^{(i)}) + \frac{1}{2} \sum_{j=1}^m \theta_j^2$$

in practice

$$\sum_{j=1}^m \theta_j^2 = \theta^T \theta \leftarrow \theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_m \end{bmatrix} \quad \text{ignore } \theta_0$$

$$\hookrightarrow \theta^T M \theta^T \leftarrow \text{rescaled}$$

$\hookrightarrow$  able to scale to much bigger training set

why not kernels idea to other algorithms (logistic regression)?

we can, but computational tricks that apply for SVM don't generalise well to other algorithms. (run very slow!)

## SVM parameters

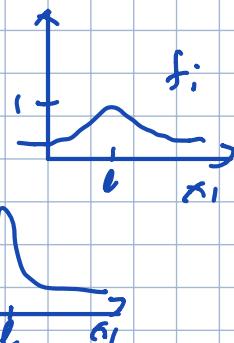
$C$  ( $= 1/\lambda$ ) : large  $C$ : low bias, high variance (small  $\sigma^2$ ), overfit

small  $C$ : high bias, low variance (large  $\sigma^2$ ), underfit

$\sigma^2$ : large  $\sigma^2$ : feature  $f_i$  very more smoothly.  
high bias, low variance

(Gaussian kernel)

small  $\sigma^2$ : feature  $f_i$  very less smoothly  
low bias, high variance



## SVM in practice: using an SVM

→ use SVM library to solve for  $\theta$

→ hyperparameters:

↳ C

↳ Kernel (similarity function)

e.g. No Kernel ("linear kernel")

↳ predict " $y=1$ " if  $\theta^T \pi > 0$

↳ use when n is large, m is small

e.g. Gaussian Kernel:

$$f_i = \exp\left(-\frac{\|\pi - \ell^{(i)}\|^2}{2\sigma^2}\right), \text{ where } \ell^{(i)} = \pi^{(i)}$$

Need to choose  $\sigma^2$

↳ if n is small, m is large

↳ implement Kernel (similarity) functions:

function  $f = \text{Kernel}(\pi_1, \pi_2)$

$$f = \exp\left(-\frac{\|\pi_1 - \pi_2\|^2}{2\sigma^2}\right)$$

return

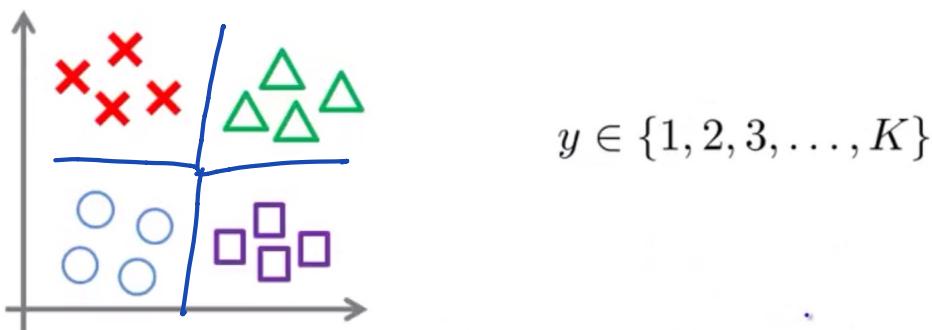
\* Note: Do perform feature scaling before using Gaussian Kernel.

### Other choice of Kernel

Not all similarity functions  $\text{similarity}(\pi, f)$  make valid kernels.  
(Need to satisfy technical condition called "Mercer's Theorem" to make sure SVM package's optimisation run correctly, and do not diverge.)

- ↙ dot product
- polynomial Kernel :  $K(x, l) = (x^T l)^2$  or  $(x^T l)^3$
- $(x^T l + 5)^2 \Rightarrow (x^T l + \text{const})^{\text{degree}}$
- ↳ perform worse than Gaussian most of the time
- ↳  $x, l$  strictly not negative
- string Kernel, chi-square Kernel, histogram intersection Kernel, ...

### Multi-class classification



Many SVM packages already have built-in multi-class classification functionality.

Otherwise, use one-vs.-all method. (Train  $K$  SVMs, one to distinguish  $y = i$  from the rest, for  $i = 1, 2, \dots, K$ ), get  $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(K)}$   
Pick class  $i$  with largest  $(\theta^{(i)})^T x$

$$\begin{array}{c} \uparrow \\ y=1 \\ \uparrow \\ y=2 \end{array}$$

### Logistic regression vs SVMs

Say  $n = \text{number of features } (x \in \mathbb{R}^{n+1})$ ,  $m = \text{number of training examples}$ .

if  $n$  is large (relative to  $m$ )!  $n \gg m$

→ use logistic regression, or SVM w/o Kernel ("linear Kernel")

if  $n$  is small,  $m$  is intermediate:

→ use SVM with Gaussian Kernel

if  $n$  is small,  $m$  is large

→ create/add more feature, then use logistic regression / SVM w/o Kernel.

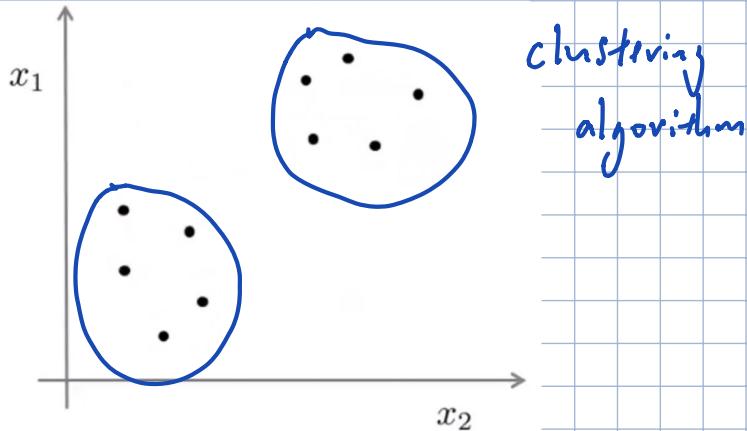
→ Neural Network likely to work well for most of these settings but may be slower to train.

SVM optimisation problem is a convex optimisation problem

↳ can always find global minimum.

## Week 8

### Unsupervised Learning



Training set:  $\{x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)}\}$   $y^{(i)}$  in training set

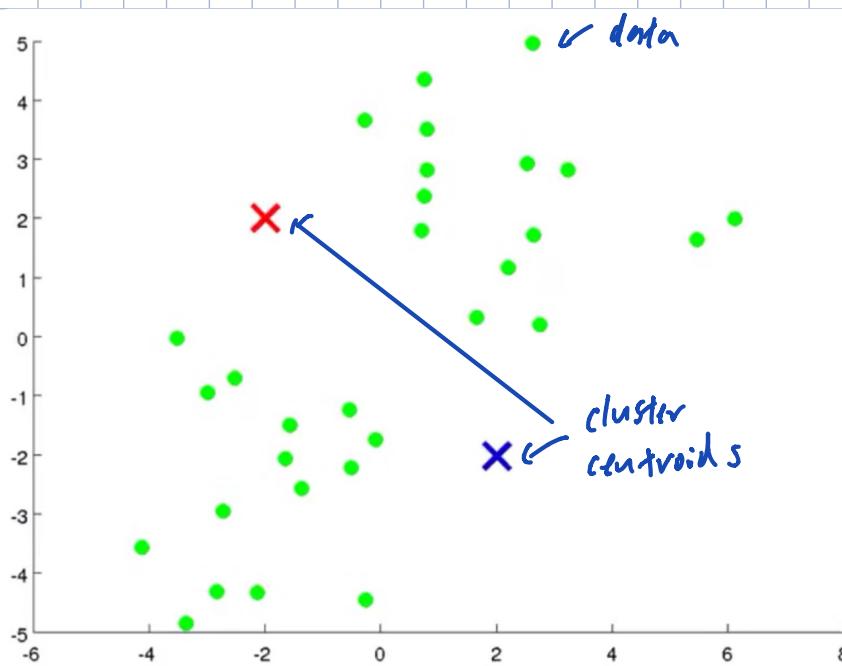
No labels

⇒ algorithm suppose to find some structure for us

### Unsupervised Learning Problem: Clustering

#### Clustering: K-means algorithm

e.g.



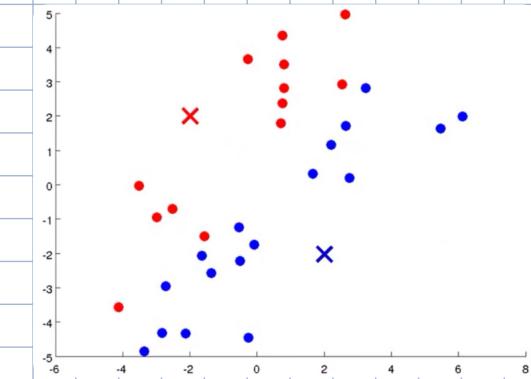
→ randomly initialise 2 points  
called **cluster centroid**

K-means is an iterative algorithm  
and it does 2 things:

- cluster assignment
- move centroid

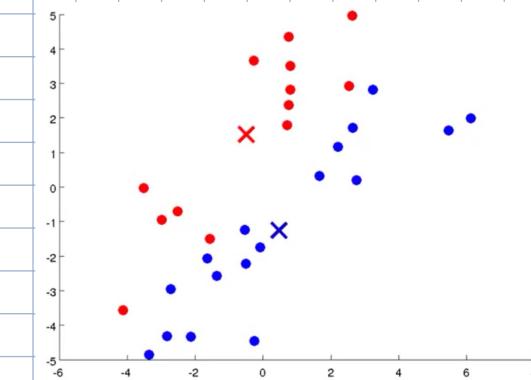
## cluster assignment

The algorithm will go through each of examples, depending on whether it is closer to the red or the blue cluster centroid, it is going to assign each of the data points to one of the cluster centroids.

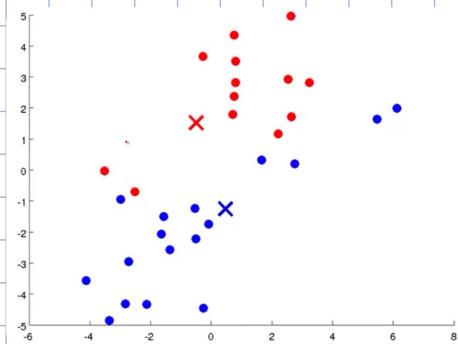


## Move centroids

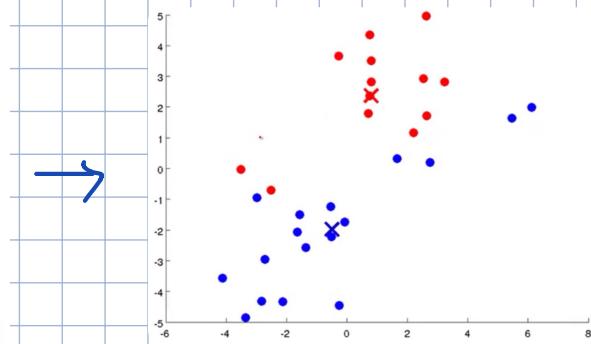
The 2 cluster centroid will be moved to the average of the points colored the same color.



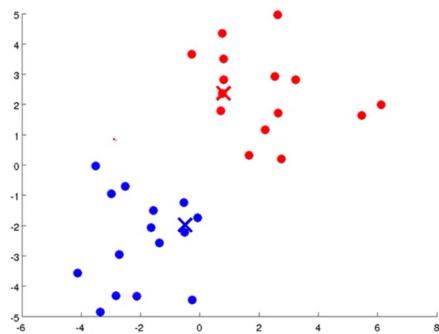
These 2 step repeats:



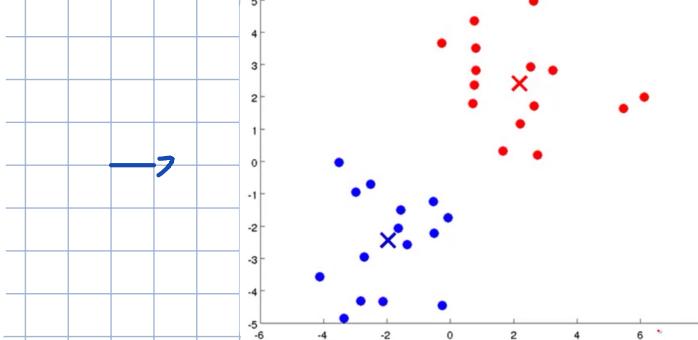
cluster assignment



move centroids



cluster assignment



move centroids

## K-means algorithm

→ takes 2 inputs:

①  $K$  (number of clusters)

② training set  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

$x^{(i)} \in \mathbb{R}^n$  (drop  $x_0 = 1$  convention)

→ Randomly initialise  $K$  cluster centroids  $M_1, M_2, \dots, M_K \in \mathbb{R}^n$

Repeat {

cluster assignment [ for  $i=1$  to  $m$   $c^{(i)}$  is a number from 1 to  $K$   
 $c^{(i)} :=$  index (from 1 to  $K$ ) of cluster centroid closest  
 to  $x^{(i)}$  i.e.  $\min_k \|x^{(i)} - M_k\|^2$  (find  $k$ )  $\Rightarrow C^{(i)}$  ]

for  $k=1$  to  $K$

$M_k :=$  average (mean) of points assigned to  
 cluster  $k$

}

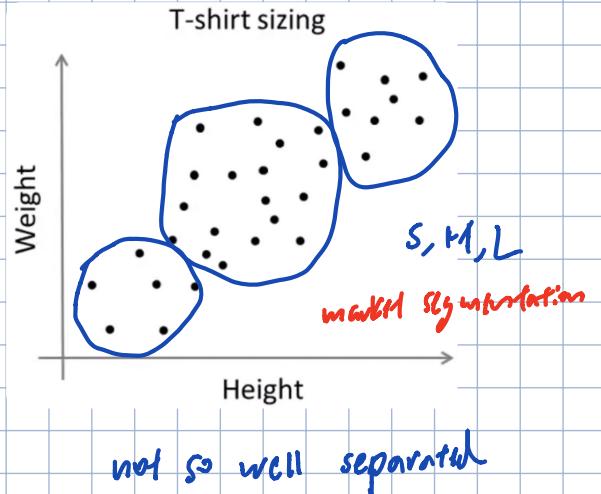
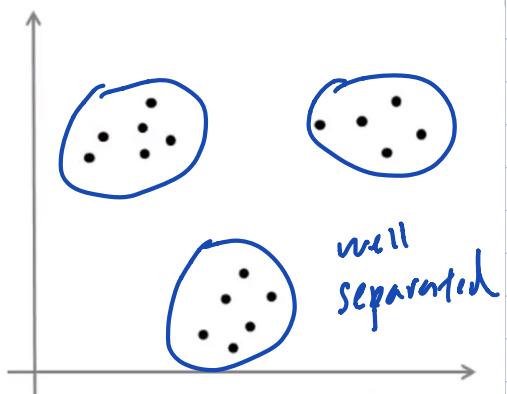
e.g. cluster 2 has training example  $x^{(1)}, x^{(5)}, x^{(6)}, x^{(10)}$

$\hookrightarrow C^{(2)} = 2, C^{(5)} = 2, C^{(6)} = 2, C^{(10)} = 2$

$\hookrightarrow M_2 = \left[ \frac{x^{(1)} + x^{(5)} + x^{(6)} + x^{(10)}}{4} \right] \in \mathbb{R}^n$

If there is no point  
 in cluster, eliminate  
 the cluster centroid

## K-means for non-separated clusters



## K-means: optimisation objective

For supervised learning, they all have an optimisation objective, or some cost function that algorithm was trying to minimise.

Same for K-means.

While K-means is running, it keeps track of 2 sets of variables.

→  $C^{(i)}$  = index of cluster ( $1, 2, \dots, K$ ) to which example  $x^{(i)}$  is currently assigned

→  $M_k$  = cluster centroid  $k$  ( $M_k \in \mathbb{R}^n$ )

One more notation:

$M_{c(i)}$  = cluster centroid of cluster to which example  $x^{(i)}$  has been assigned. ↳ if  $x^{(i)} \rightarrow 5$  (i.e. assigned to cluster 5)  
 $c^{(i)} = 5$

$$\hookrightarrow M_{c(i)} = M_5$$

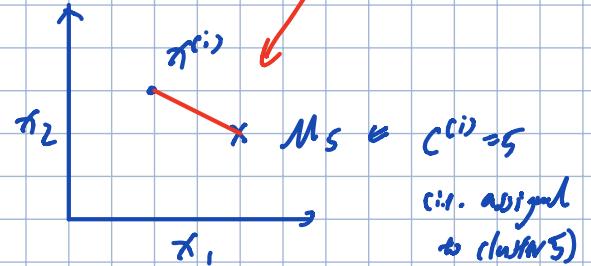
## Optimisation objective:

$$J(C^{(1)}, \dots, C^{(m)}, M_1, \dots, M_k) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - M_{c^{(i)}}\|^2$$

$$\min_{C^{(1)}, \dots, C^{(m)}} J(C^{(1)}, \dots, C^{(m)}, M_1, \dots, M_k)$$

$$\min_{M_1, \dots, M_k} J(C^{(1)}, \dots, C^{(m)}, M_1, \dots, M_k)$$

try to find  $C^{(i)}, M_k$  that minimize  $J$ .



This cost function  $J$  is sometimes called Distortion.

## K-means algorithm

Randomly initialize  $K$  cluster centroids  $\mu_1, \mu_2, \dots, \mu_K \in \mathbb{R}^n$

Repeat {  
    cluster assignment steps:  
        minimize  $J(\dots)$  wrt  $C^{(1)}, C^{(2)}, \dots, C^{(m)}$   
        (holding  $M_1, \dots, M_k$  fixed)}

for  $i = 1$  to  $m$

$c^{(i)} :=$  index (from 1 to  $K$ ) of cluster centroid  
closest to  $x^{(i)}$

for  $k = 1$  to  $K$

$\mu_k :=$  average (mean) of points assigned to cluster  $k$

}

move centroid steps:

minimize  $J(\dots)$  wrt  $M_1, \dots, M_k$

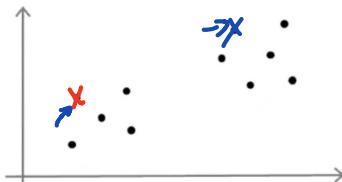
## K-means: Random Initialization

### Random initialization

Should have  $K < m$

Randomly pick  $K$  training examples.

Set  $\mu_1, \dots, \mu_K$  equal to these  $K$  examples.

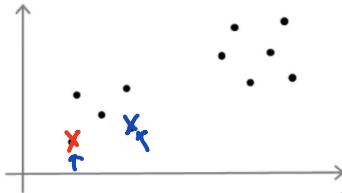


e.g. say  $K=2$ , pick 2 examples

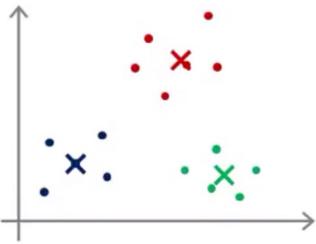
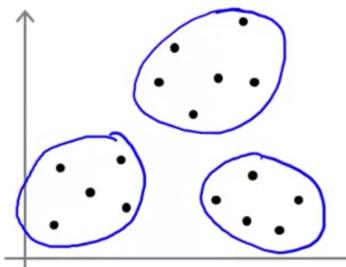
$$\Rightarrow M_1 = x^{(i)}$$

$$M_2 = x^{(j)}$$

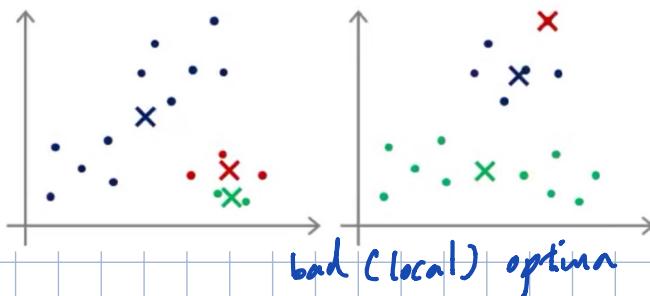
$\Rightarrow$  K-means can end up in different solutions depending on how initialisation gets.



Local optima



good optima (global)



Local minima:  
not minimising  
 $J$  very well.

$\Rightarrow$  we can try multiple random initialization

### Random initialisation

For  $i = 1$  to  $100$  {

Randomly initialise K-means.

Run K-means. Get  $C^{(1)}, \dots, C^{(m)}, M_1, \dots, M_K$ .

Compute cost function (distortion)

$J(C^{(1)}, \dots, C^{(m)}, M_1, \dots, M_K)$

}

Pick clustering that gave the lowest cost  $J(C^{(1)}, \dots, C^{(m)}, M_1, \dots, M_K)$

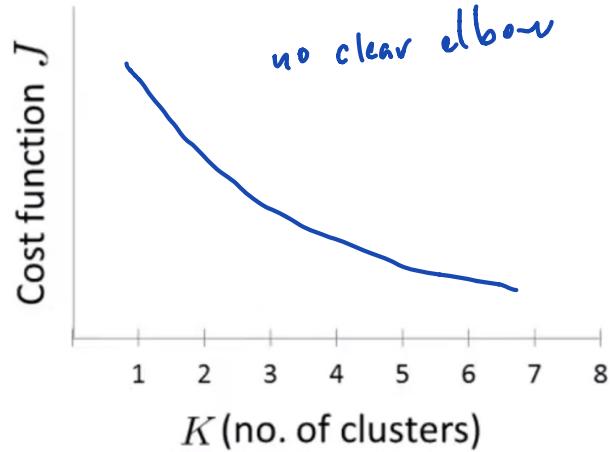
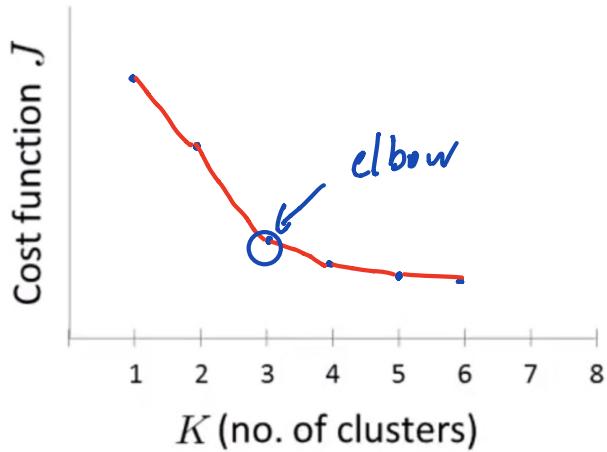
$\Rightarrow$  usually for small  $K$  (e.g. 2-10), random initialisation a few times might help. If large  $K$ , then not so much.

## K-means: choosing the number of clusters

commonly choose  $K$  by hand (through visualisation of data etc)

It is often ambiguous how many clusters there are in the data

### Elbow Method

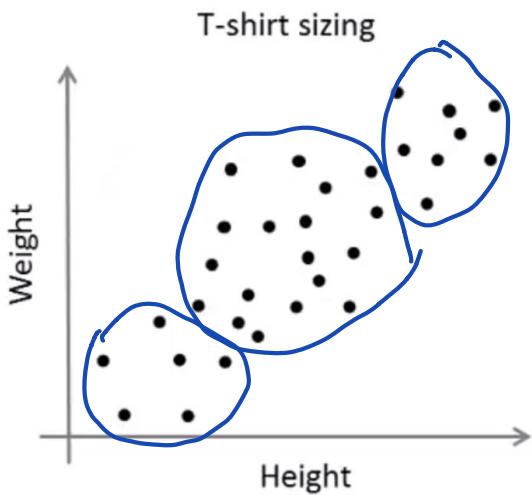


### Choosing the value of $K$

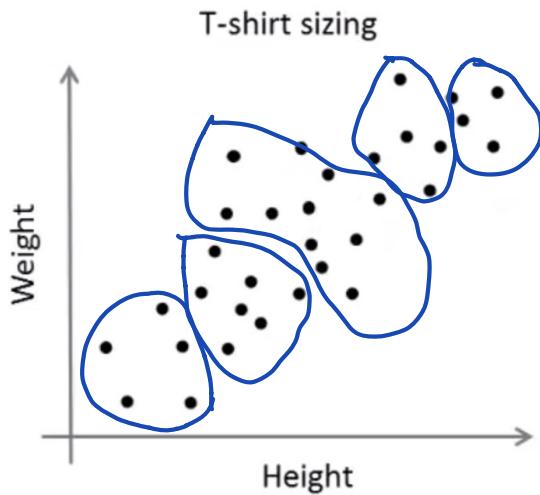
Sometimes, you're running K-means to get clusters to use for some later/downstream purpose. Evaluate K-means based on a metric for how well it performs for that later purpose.

$K=3 S, M, L$

E.g.



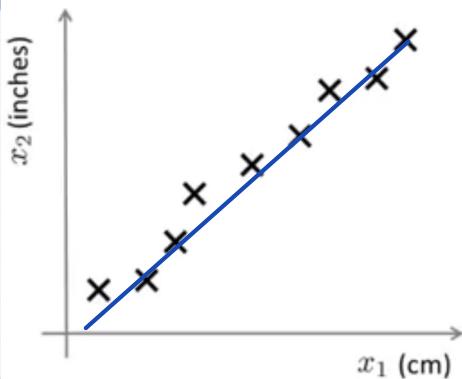
$K=5 XS, S, M, L, XL$



## Unsupervised Learning Problem: Dimensionality Reduction

### Motivation: Data Compression

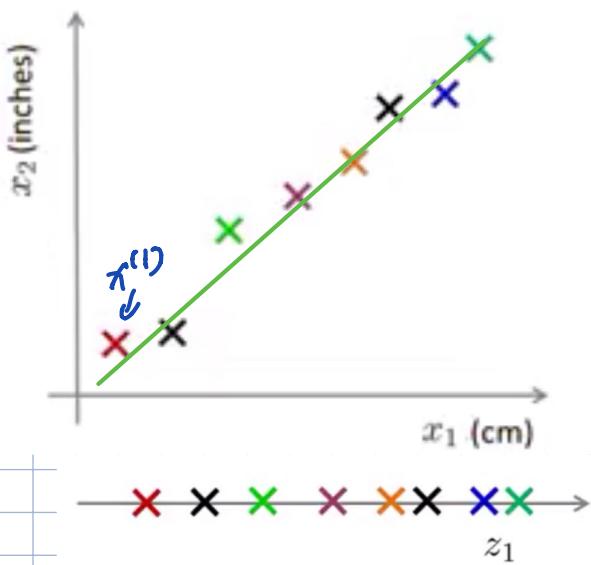
$2D \rightarrow 1D$



Reduce data from  
2D to 1D

Let's say a dataset with many features, unknown to us, 2 of the features are length in different unit. This gives a highly redundant representation  $\Rightarrow$  reduce the data (plotted above) to 1 D.

(Reasons that these 2 features are not on a perfectly line is due to round-off during data collection/processing)



Reduce data from  
2D to 1D

$$x^{(1)} \in \mathbb{R}^2 \rightarrow z^{(1)} \in \mathbb{R}^1$$

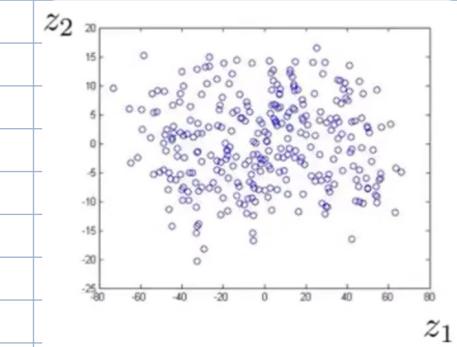
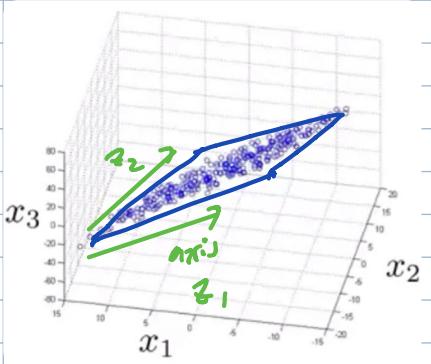
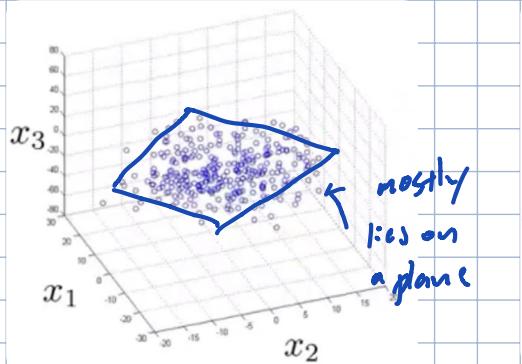
⋮

$$x^{(m)} \rightarrow z^{(m)}$$

↳ new feature  $z_1$ , - on the green line

more specifically here, by dimensionality reduction, we would find the direction where the data seems to lie (green line), and project all the data on that line.

$3D \rightarrow 2D$



$$x^{(i)} \in \mathbb{R}^3$$

project all the data on a plane

$$z^{(i)} \in \mathbb{R}^2$$

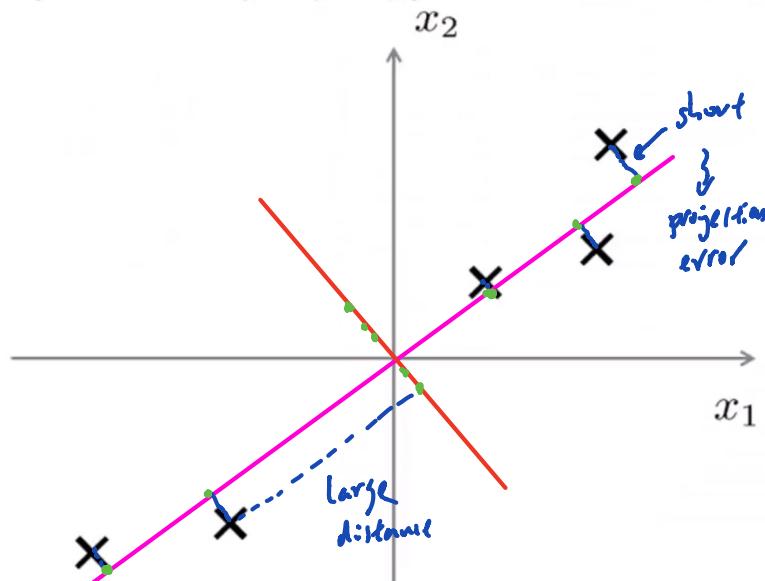
$$z = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \Rightarrow z^{(i)} = \begin{bmatrix} z_1^{(i)} \\ z_2^{(i)} \end{bmatrix}$$

### Motivation: Data Visualisation

Instead of  $x^{(i)} \in \mathbb{R}^{50}$ , we can have a different feature representation we can plot them out in  $\mathbb{R}^2$ . (new feature:  $z \in \mathbb{R}^2$ )

### Dimensionality Reduction: Principal Component Analysis (PCA)

#### Principal Component Analysis (PCA) problem formulation



say we have data set

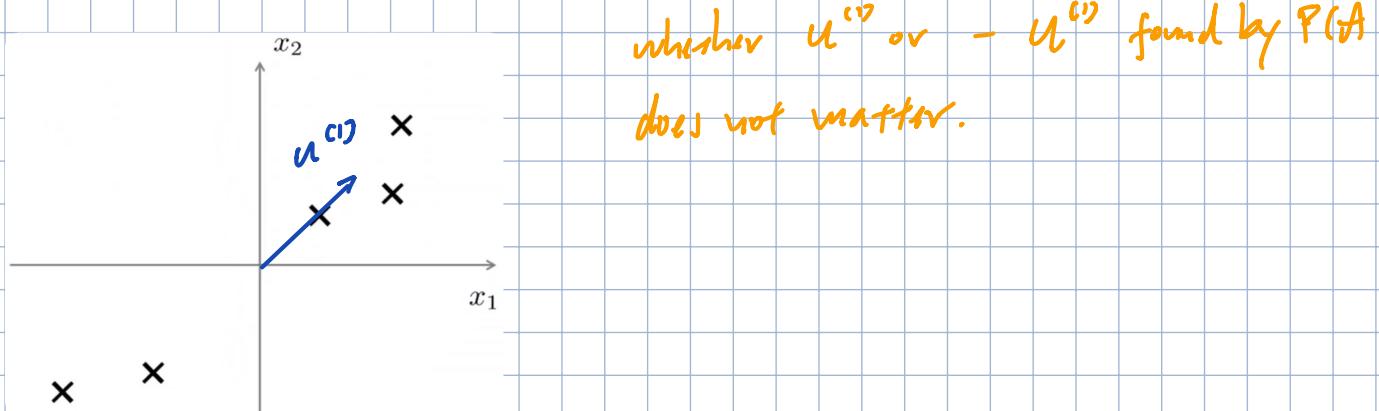
$$X \in \mathbb{R}^2$$

and would like  $2D \rightarrow 1D$

PCA will try to find the projection plane that minimizes projection error

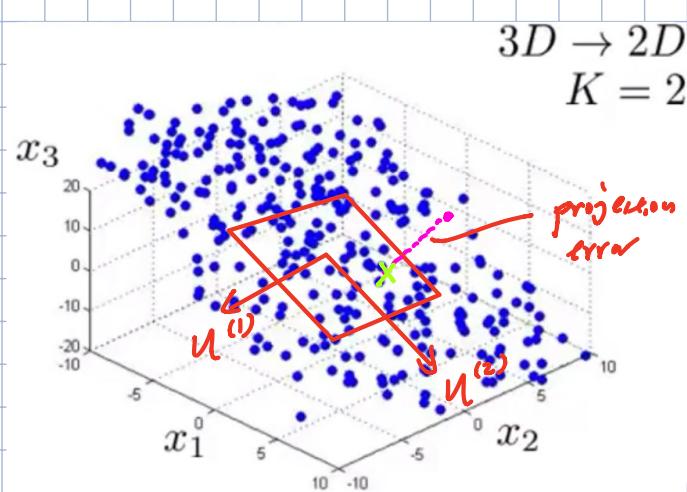
Before Applying PCA, it is standard practice to first perform mean normalisation or feature scaling so that the feature  $x_1$  and  $x_2$  should have 0 mean, and have comparable range of values.

Goal of PCA: Reduce from 2-dimension to 1-dimension: Find a direction (a vector  $u^{(1)} \in \mathbb{R}^n$ ) onto which to project the data so as to minimise the projection error.



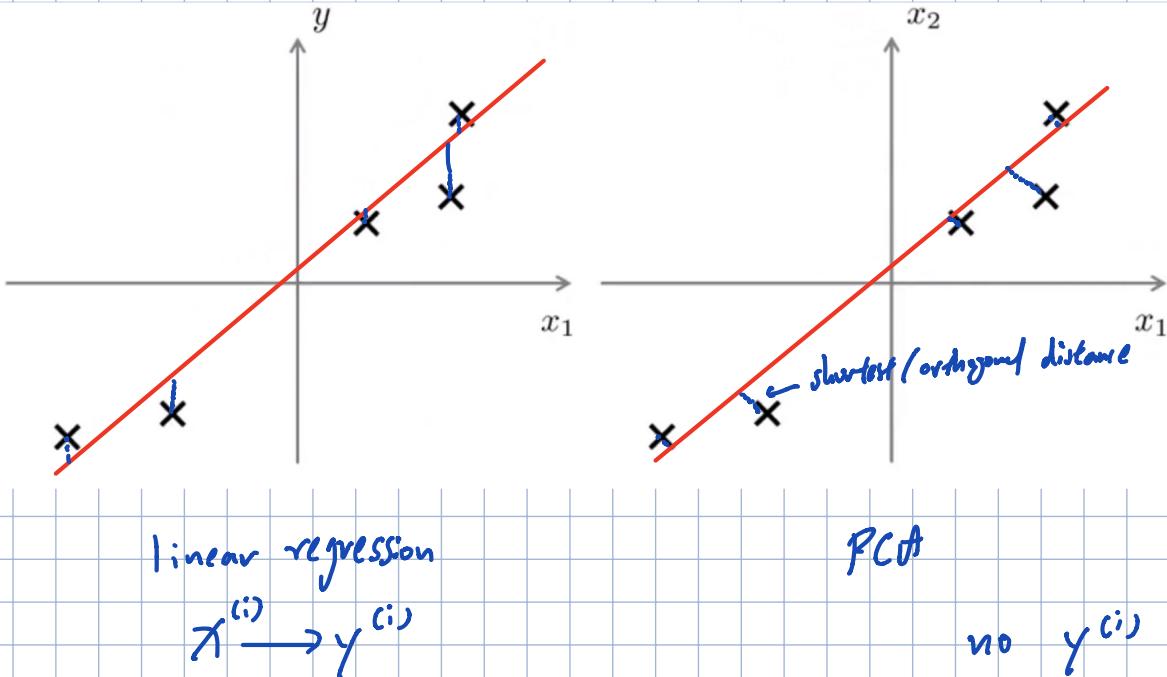
More generally:

Reduce from  $n$ -dimension to  $K$ -dimension: Find  $K$  vectors  $u^{(1)}, u^{(2)}, \dots, u^{(K)}$  onto which to project the data, so as to minimise the projection error.



i.e. project data onto the linear subspace spanned by this set of  $u^{(K)}$  vectors.

PCA is NOT linear regression:



### PCA: algorithm

Before applying PCA, there is a data pre-processing step we should always do.

Training set:  $\mathcal{X}^{(1)}, \mathcal{X}^{(2)}, \dots, \mathcal{X}^{(m)}$

→ pre processing (feature scaling / mean normalisation):

$$\hookrightarrow M_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)} \quad (\text{get mean of feature } x_j \text{, across all } m \text{ samples})$$

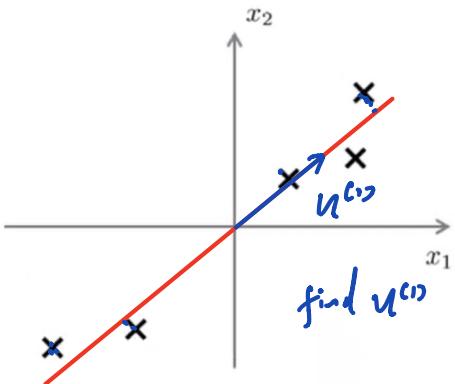
$\hookrightarrow$  replace each  $x_j^{(i)}$  with  $x_j^{(i)} - M_j$  (all samples's  $x_j$  feature - mean)  
(each feature is zero mean)

$\hookrightarrow$  If different features on different scales (e.g.  $x_1$  = size of house,  $x_2$  = number of bedrooms), scale feature to have comparable range of values.

$$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - M_j}{S_j}$$

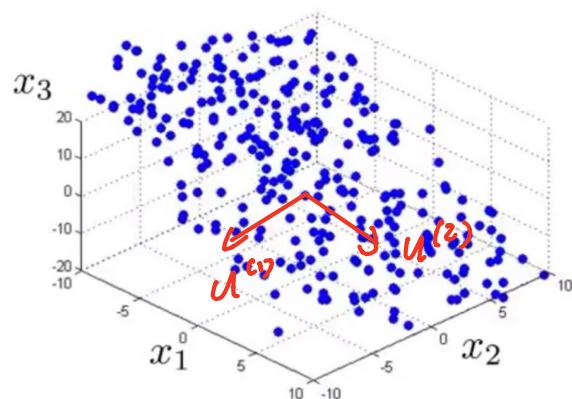
$S_j$  std dev of  $x_j$

Recall PCA does:



Reduce data from 2D to 1D

$$x^{(i)} \in \mathbb{R}^2 \rightarrow z^{(i)} \in \mathbb{R}$$



Reduce data from 3D to 2D

$$x^{(i)} \in \mathbb{R}^3 \rightarrow z^{(i)} \in \mathbb{R}^2$$

$$z = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$$

PCA needs to compute:

- (1)  $u^{(1)}, \dots, u^{(k)}$
- (2)  $z^{(1)}, \dots, z^{(m)}$

### PCA algorithm

Reduce data from n-dimensions to k-dimensions

↳ Compute "covariance matrix":

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)}) (x^{(i)})^T$$

$$\begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix} \quad [x_1 \dots x_m]$$

↳ Compute "eigenvectors" of matrix  $\Sigma$ :

$$[U, S, V] = svd(\Sigma) \quad \begin{array}{l} \text{singular value decomposition} \\ \text{T} \quad m \times m \text{ matrix} \end{array}$$

$$U = \begin{bmatrix} | & | & | \\ u^{(1)} & u^{(2)} & \dots & u^{(n)} \\ | & | & | \end{bmatrix} \quad U \in \mathbb{R}^{n \times n}$$

$n$   
 $k$

↳ select the first k vectors  $\Rightarrow u^{(1)}, \dots, u^{(k)}$

Next we need find the  $Z$ :  $x \in \mathbb{R}^n \rightarrow z \in \mathbb{R}^k$

→ construct a matrix use the 1st K columns of U:

$$U_{\text{reduced}} = \begin{bmatrix} | & | & | \\ u^{(1)} & u^{(2)} & \dots & u^{(K)} \\ | & | & | \end{bmatrix} \text{ an } n \times K \text{ matrix}$$

$$Z^{(i)} = U_{\text{reduced}}^T X^{(i)} \quad X^{(i)} = \begin{bmatrix} -(U^{(1)})^T \\ -(U^{(2)})^T \\ \vdots \\ -(U^{(n)})^T \end{bmatrix} \quad \begin{bmatrix} x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix}$$

$n \times 1$

$K \times n$

$K \times 1$

for 1 example

$$z_j = (u^{(j)})^T x$$

### Summary

→ After mean normalisation (ensure every feature has zero mean) and optionally feature scaling!

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)}) (x^{(i)})^T$$

$$[U, S, V] = \text{svd}(\Sigma);$$

$$U_{\text{reduce}} = U(:, 1:K);$$

$$Z = U_{\text{reduce}}' \times X; \quad X \in \mathbb{R}^n \quad X \neq 1$$

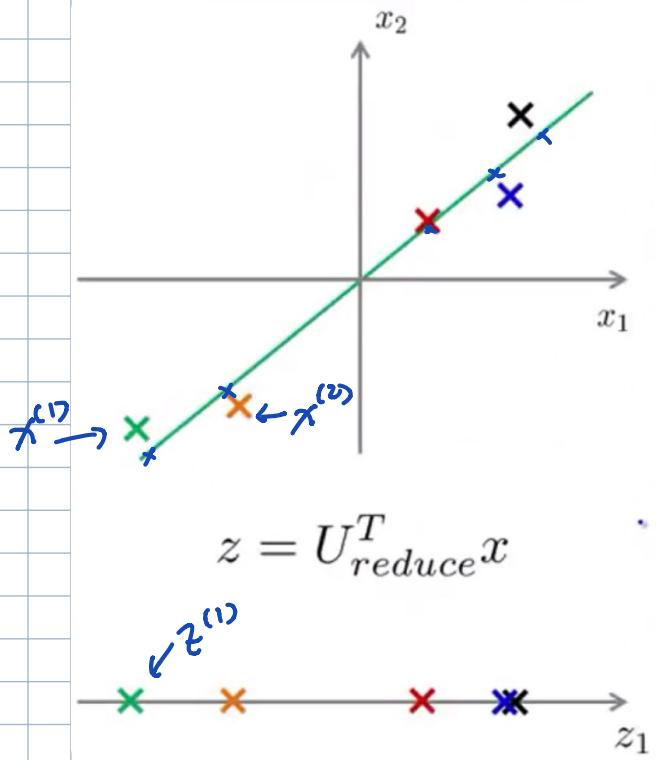
if  $X = \begin{bmatrix} -x^{(1)} \\ -x^{(2)} \\ \vdots \\ -x^{(n)} \end{bmatrix}$   $X' = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(n)} \end{bmatrix}$

verticalised

$$\Sigma = (\frac{1}{m}) \times X' \times X;$$

\* This procedure gives the projection of data onto K-dimensional subspace that actually minimise the squared projection error.

## PCA: Reconstruction from Compressed Representation

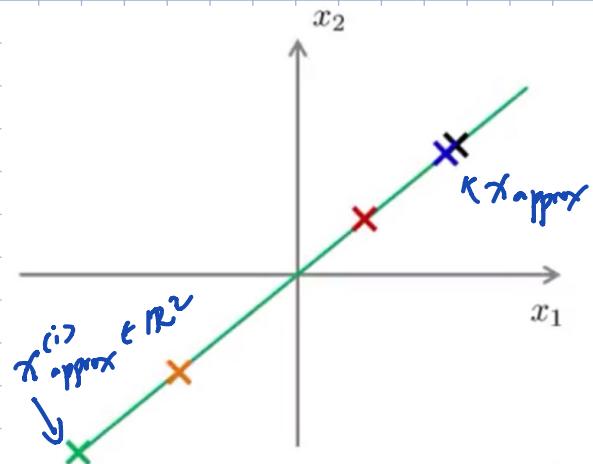


Given points on  $\mathbb{Z}_1$ , we would like to recover points in  $\mathbb{R}^2$ .

$$\text{i.e. } \mathbb{Z} \in \mathbb{R} \rightarrow x \in \mathbb{R}^2$$

$$x_{\text{approx}} = U_{\text{reduce}} \cdot \mathbb{Z}$$

$$\mathbb{R}^n \quad \underbrace{\begin{matrix} n \times k \\ n \times 1 \end{matrix}}_{k \times 1} \quad \approx x$$



## PCA: choosing the number of principle components

In PCA, we take  $n$ -dimensional features and reduce them to some  $k$ -dimensional feature representation.

$k$ : a hyperparameter that specify no. of principle components that we've retained.

what PCA is trying to do is that it tries to minimise the average squared projection error.

→ Average squared projection error:  $\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2$

→ Total variation in the data:  $\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2$   
 In average, how far are my training examples from the vectors that are all zeros?

+ avg length of the training triangles

Typically, choose K to be the smallest value so that

$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01 \quad (1\%)$$

↳ 99% of variance is retained

### Implementation

Algorithm:

try PCA with K=1

compute  $U_{\text{reduced}}, z^{(1)}, z^{(2)}, \dots, z^{(m)}$ ,

$x_{\text{approx}}^{(1)}, \dots, x_{\text{approx}}^{(m)}$

check if

$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01 ?$$

if no good, try K=2

:

till  $\leq 0.01$

not efficient!

$$[U, S, V] = \text{svd}(\text{Sigma})$$

$$S = \begin{bmatrix} s_{11} & s_{21} & \dots & 0 \\ 0 & s_{22} & \dots & s_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & s_{nn} \end{bmatrix}$$

n × n  
diag matrix

For given K:

$$1 - \frac{\sum_{i=1}^K s_{ii}}{\sum_{i=1}^n s_{ii}} \leq 0.01$$

Set K from 1 → n, to ensure this 99.99%  
 call svd only once.

## Advice for applying PCA

use PCA to speed up a learning algorithm

→ Supervised learning speed up

$$(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})$$

Say  $\mathbf{x}^{(i)} \in \mathbb{R}^{10,000}$  i.e. high dimensional feature

e.g. in computer vision  → 10,000 features, 1 for each pixel

we first take our labelled training set and extract just the inputs:

$$\hookrightarrow \text{unlabelled dataset: } \mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)} \in \mathbb{R}^{10,000}$$

↓ PCA

$$\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(m)} \in \mathbb{R}^{1000}$$

↪ New training set:

$$(\mathbf{z}^{(1)}, y^{(1)}), (\mathbf{z}^{(2)}, y^{(2)}), \dots, (\mathbf{z}^{(m)}, y^{(m)})$$

proceed with supervised learning.

e.g. logistic regression:  $h_{\theta}(\mathbf{z}) = \frac{1}{1 + e^{-\theta^T \mathbf{z}}}$

↪ If new test example  $\mathbf{x}^{(j)}$ , then use the same mapping to get the corresponding  $\mathbf{z}^{(j)}$ , then use the trained hypothesis  $h_{\theta}(\mathbf{z})$  to predict.

\* Mapping  $\mathbf{x}^{(i)} \rightarrow \mathbf{z}^{(i)}$  should be defined by running PCA only on the training set. This mapping can be applied as well to  $\mathbf{x}_{cv}^{(i)}$  and  $\mathbf{x}_{test}^{(i)}$  in the cross validation and test sets.

i.e.  $U$  reduced is a parameter learned by PCA through training step.

### Application of PCA

→ compression

↳ reduce memory/disk needed to store data

↳ speed up learning algorithm

→ choose  $K$  by % of variance retained

→ visualisation

→ choose  $K=2$  or  $K=3$

### Bad use of PCA: To prevent overfitting

Reasoning: use  $Z^{(i)}$  instead of  $X^{(i)}$  to reduce the number of features to  $K$  ( $c_n$ ). Thus - fewer features, less likely to overfit.

This might work ok, but is NOT a good way to address overfitting.

use regularisation instead.

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

PCA throw away info on  $\hat{y}$ , w/o knowing  $y$  (unsupervised!)

Regularisation knows  $y$  value during minimisation, it is better.

Other times PCA is used where it shouldn't be

Design of ML systems! (the plan)

- Get training set  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$
- Run PCA to reduce  $x^{(i)}$  in dimension to get  $z^{(i)}$
- Train logistic regression on  $\{(z^{(1)}, y^{(1)}), \dots, (z^{(m)}, y^{(m)})\}$
- Test on test set:
  - Map  $x_{\text{test}}^{(i)}$  to  $z_{\text{test}}^{(i)}$ .
  - Run  $h_\theta(z)$  on  $\{(z_{\text{test}}^{(1)}, y_{\text{test}}^{(1)}), \dots, (z_{\text{test}}^{(m)}, y_{\text{test}}^{(m)})\}$

How about doing the whole thing w/o PCA?

Before implementing PCA, first try the original data  $x^{(i)}$ .

Only if that doesn't do what we want, then use PCA and consider using  $z^{(i)}$ .

## Week 9

### Anomaly Detection: motivation

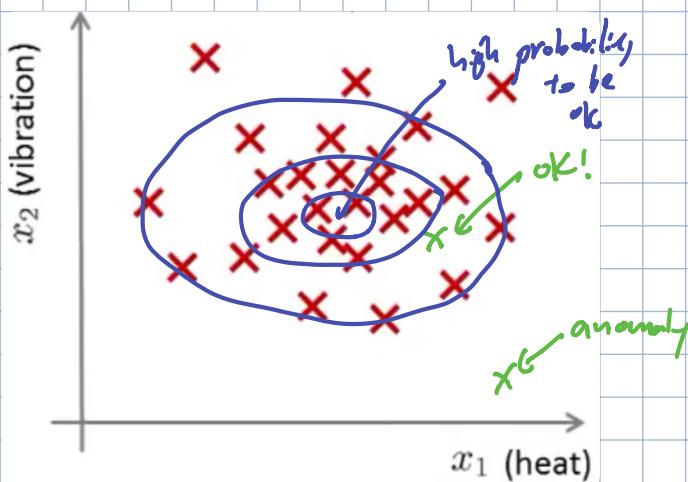
e.g. Aircraft engine features:

$x_1$  = heat generated

$x_2$  = vibration intensity

...

Data set:  $\{x^{(1)}, x^{(2)}, \dots, x^{(n)}\}$



Say now we have a new aircraft engine, and have some set of features  $x_{\text{test}}$ .

→ We want to know if this engine is anomalous, i.e. if engine should go further checking

More formally, in the anomaly detection problem, we are given some data set:  $\{x^{(1)}, x^{(2)}, \dots, x^{(n)}\}$  (non-anomalous examples)

Question: is  $x_{\text{test}}$  anomalous?

Approach: given the unlabelled training set, build a model:  $p(x)$ , i.e. probability of  $x$  (where  $x$  are the features of aircraft engines)

↳ then if  $p(x_{\text{test}}) < \epsilon$  → flag as anomaly

$p(x_{\text{test}}) \geq \epsilon \rightarrow \text{OK}$

\*  $p(\mathbf{x})$  is actually the normalized probability density as parameterised by the feature vector,  $\mathbf{x}$ .

$\epsilon$  is a threshold condition on the probability density.

↳ Determination of the actual probability need integration of this density over the appropriate extent of phase space.

Anomaly Detection example:

→ Fraud detection:

$\mathbf{x}^{(i)}$  = features of user i's activity.

model  $p(\mathbf{x})$  from data. (what is the probability that different users behave in different ways?)

Identify unusual users by checking which have  $p(\mathbf{x}) < \epsilon$

e.g.  $x_1$ : how often this user logs in

$x_2$ : pages visited

$x_3$ : no. of transactions

:

→ Manufacturing

→ Monitoring computers in a data center:

$\mathbf{x}^{(i)}$  = features of machine i

$x_1$  = memory use,  $x_2$  = no. of disk access

$x_3$  = CPU load,  $x_4$  = CPU load/network traffic

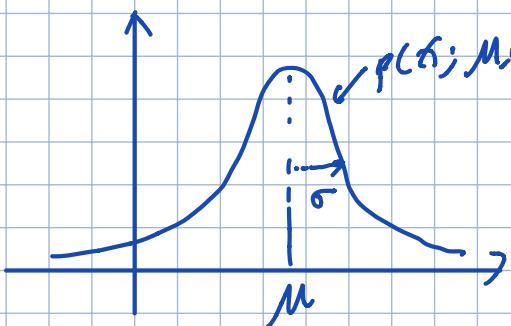
## Anomaly Detection: Gaussian distribution

Also called Normal Distribution.

Say  $\pi \in \mathbb{R}$  (i.e. a row number). If  $\pi$  is a distributed Gaussian with mean  $\mu$ , variance  $\sigma^2$ . ( $\sigma$ : standard deviation)

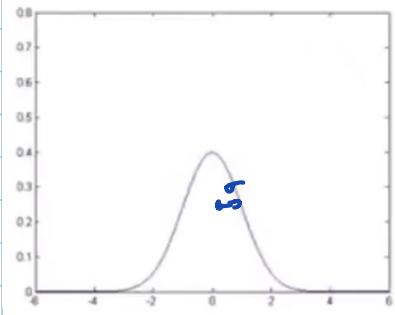
$$\pi \sim \mathcal{N}(\mu, \sigma^2)$$

T "distributed as"



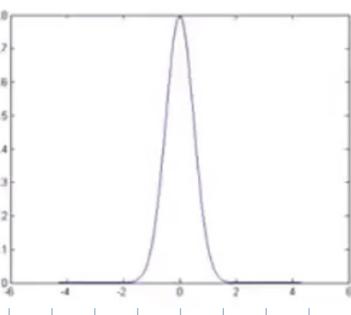
$$p(\pi; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(\pi-\mu)^2}{2\sigma^2}\right)$$

$$\mu = 0, \sigma = 1$$

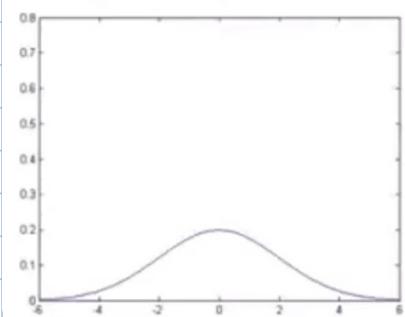


$$\mu = 0, \sigma = 0.5$$

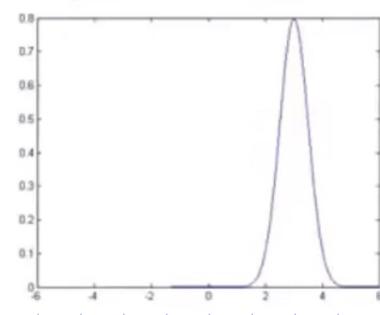
$$\sigma^2 = 0.25$$



$$\mu = 0, \sigma = 2$$



$$\mu = 3, \sigma = 0.5$$



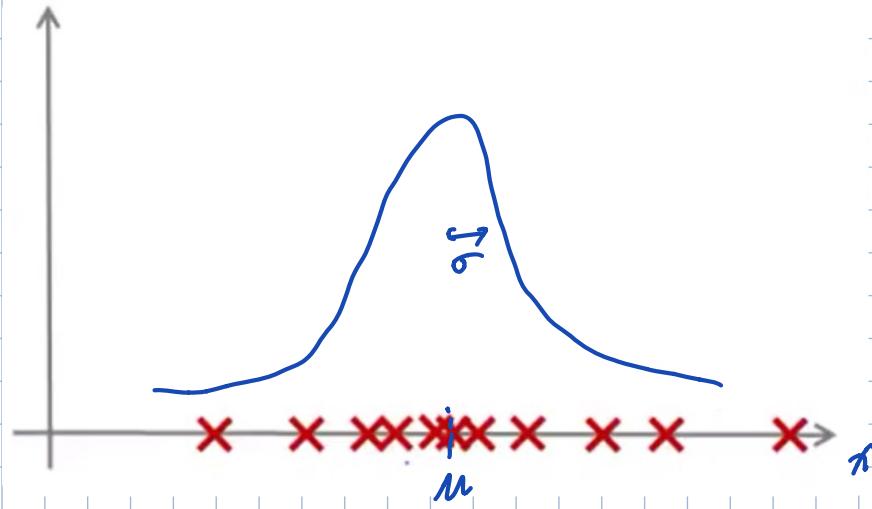
## Parameter estimation problem

Dataset:  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$   $x^{(i)} \in \mathbb{R}$

Let's say we suspect these examples form a Gaussian distribution.

i.e.  $x^{(i)} \sim \mathcal{N}(\mu, \sigma^2)$

↳ Given the dataset, try to estimate what are the values of  $\mu$  and  $\sigma^2$ .



$$\begin{cases} \mu = \frac{1}{m} \sum_{i=1}^m x^{(i)} \\ \sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^2 \end{cases}$$

they are also the maximum likelihood estimates of the prior of  $\mu$  and  $\sigma^2$ .

↓ sometimes use  
 $m-1$ , not much difference  
in practice.

## Anomaly Detection: Algorithm

Let's say we have an unlabeled training set of  $m$  examples!

training set:  $\{x^{(1)}, \dots, x^{(m)}\}$ ,

each example is  $x \in \mathbb{R}^n$

The way we are going to address anomaly detection, is going to model  $p(x)$  from the dataset.

$$\hookrightarrow p(x) = p(x_1; M_1, \sigma_1^2) p(x_2; M_2, \sigma_2^2) \dots p(x_n; M_n, \sigma_n^2)$$

$$= \prod_{j=1}^n p(x_j; M_j, \sigma_j^2)$$

↙ we have to assume the feature  $x_i$  is distributed according to a Gaussian distribution, with some mean and variance:

$$x_i \sim N(M_i, \sigma_i^2)$$

$$x_2 \sim N(M_2, \sigma_2^2)$$

\* features by right should be independent when using this equation, but in practice it works just fine.

### The Anomaly Detection Algorithm

①. choose features  $x_i$  that we think might be indicative of anomalous examples

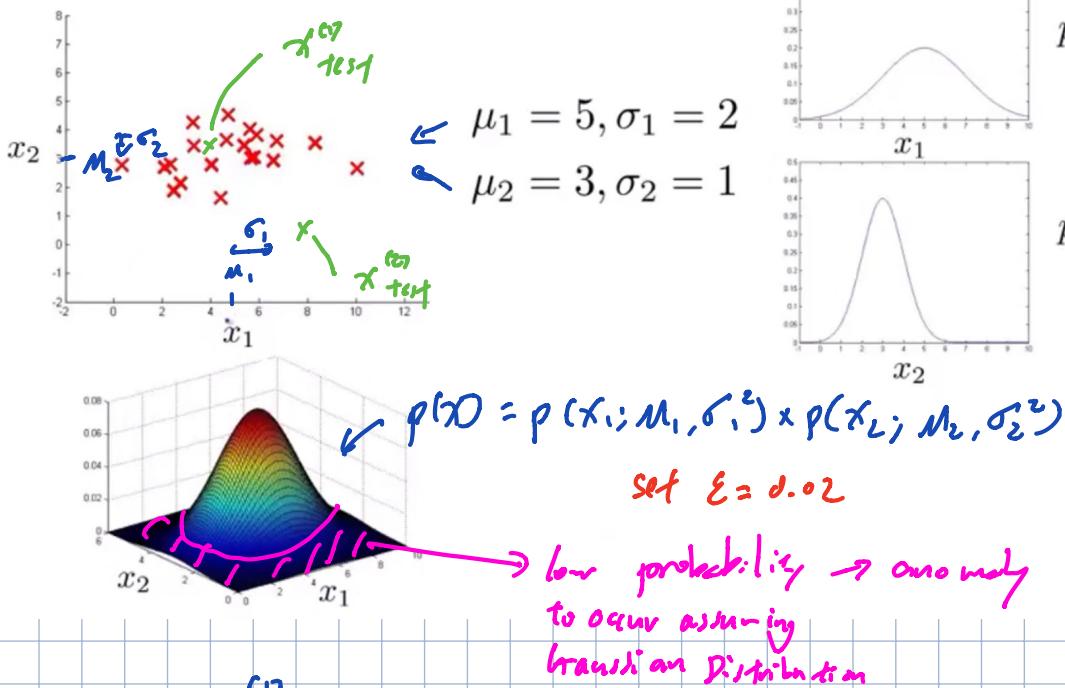
②. fit parameters  $M_1, \dots, M_n, \sigma_1^2, \dots, \sigma_n^2$  { $x^{(1)}, \dots, x^{(m)}$ }  
 $p(x_j; M_j, \sigma_j^2)$   $M_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$  i.e. for each feature across all m training samples  
 $\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - M_j)^2$  if  $M = \begin{bmatrix} M_1 \\ \vdots \\ M_n \end{bmatrix} = \frac{1}{m} \sum_{i=1}^m x^{(i)}$  (vectorize)

③. Given new example  $x$ , compute  $p(x)$ :

$$p(x) = \prod_{j=1}^n p(x_j; M_j, \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi} \sigma_j} \exp\left(-\frac{(x_j - M_j)^2}{2\sigma_j^2}\right)$$

Anomaly if  $p(x) < \epsilon$

## Anomaly detection example



$$\Rightarrow p(x_{\text{test}}^{(1)}) = 0.0426 > \epsilon$$

$$p(x_{\text{test}}^{(2)}) = 0.0021 < \epsilon \quad (\text{anomaly})$$

### Anomaly Detection: developing / evaluating an anomaly detection system

when developing a learning algorithm (choosing features, etc.), making decisions is much easier if we have a way of evaluating our algorithm.

Assume we have some labeled data, of anomalies and non-anomalies examples ( $y=0$  if normal,  $y=1$  if anomalous)

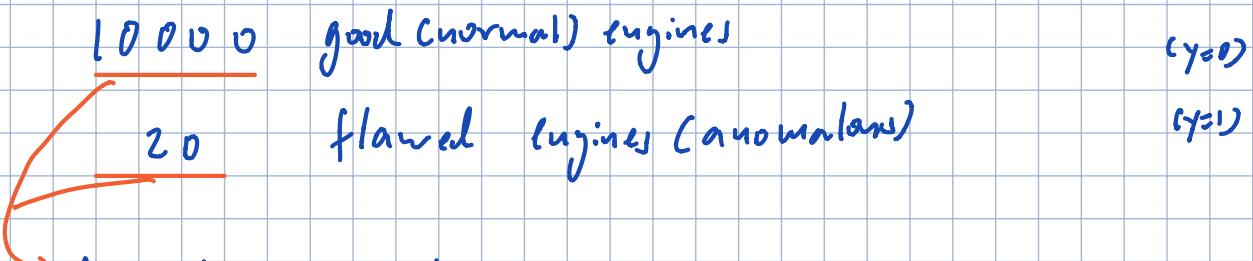
(previously we treat anomaly detection as an unsupervised problem)

$\Rightarrow$  training set:  $x^{(1)}, x^{(2)}, \dots, x^{(m)}$  (assume normal examples / not anomalies)

$\Rightarrow$  cross-validation set:  $(x_{cv}^{(1)}, y_{cv}^{(1)}), \dots, (x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$

$\Rightarrow$  test set:  $(x_{\text{test}}^{(1)}, y_{\text{test}}^{(1)}), \dots, (x_{\text{test}}^{(m_{\text{test}})}, y_{\text{test}}^{(m_{\text{test}})})$

## c-5. aircraft engine



A typical way of splitting:

training set: 6000  $\xrightarrow{M_1, \sigma_1^2, \dots, M_n, \sigma_n^2}$  good engines  $\rightarrow p(x) = p(x_1; M_1, \sigma_1^2) \dots p(x_n; M_n, \sigma_n^2)$

(CV: 2000 good engines ( $y=0$ ), 10 anomalies ( $y=1$ ))

test: 2000 good engines ( $y=0$ ), 10 anomalies ( $y=1$ )

Alternatively, some people do this:

train set: 6000 good engines

(CV: 4000 good engines ( $y=0$ ), 10 anomalies ( $y=1$ ))

test: 4000 good engines ( $y=0$ ), 10 anomalies ( $y=1$ )

## Algorithm Evaluation

→ Fit model  $p(x)$  on training set  $\{x^{(1)}, \dots, x^{(m)}\}$

→ On a cross validation/test example  $x_i$ , predict

$$y = \begin{cases} 1 & \text{if } p(x) < \varepsilon \text{ (anomaly)} \\ 0 & \text{if } p(x) \geq \varepsilon \text{ (normal)} \end{cases}$$

↳ Possible evaluation metrics

↳ confusion matrix

↳ Precision / recall

↳  $F_1$ -Score

Classification accuracy is not a good

metric

↳ very skewed data

→ can also use CV set to choose  $\varepsilon$ .

## Anomaly Detection vs Supervised Learning

### Anomaly Detection

- Very small no. of positive examples ( $y=1$ ). ( $0-20$  is common)
- Large no. of negative examples ( $y=0$ )
- Many different "types" of anomalies.  
Hard for any algorithm to learn from positive examples what anomaly looks like; future anomaly may look very different.

e.g. Fraud detection

manufacturing

monitoring machines in data center

### Supervised learning

- Large no. of positive and negative examples
- Enough positive examples for algorithm to get a sense of what positive examples are like; future positive example likely to be similar to the ones in training set.

e.g. Spam email

weather predictions

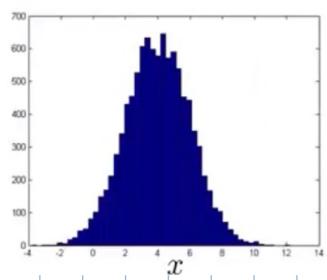
Cancer classification

## Anomaly Detection: Choosing features to use

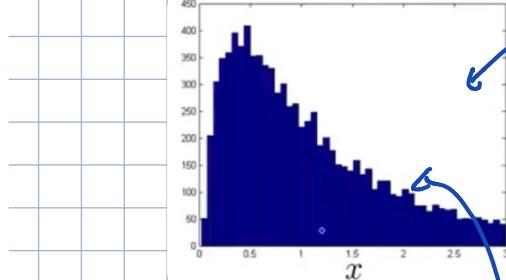
Previously we mention that the features are modeled using Gaussian distribution:  $p(x_i; \mu_i, \sigma_i^2)$

We can plot a histogram of the data to make sure they are (sort of) Gaussian. (Even if the data is non-Gaussian, model usually works fine)

good

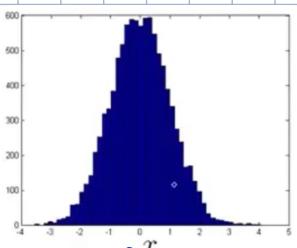


bad



need some transformation

$\log(x)$



or  $x_3 \leftarrow \sqrt{x_3}$  i.e.  $x_1 \leftarrow \log(x_1)$

$x_4 \leftarrow \sqrt[3]{x_4}$

How we come up with features for an anomaly detection?  
→ via error analysis procedure

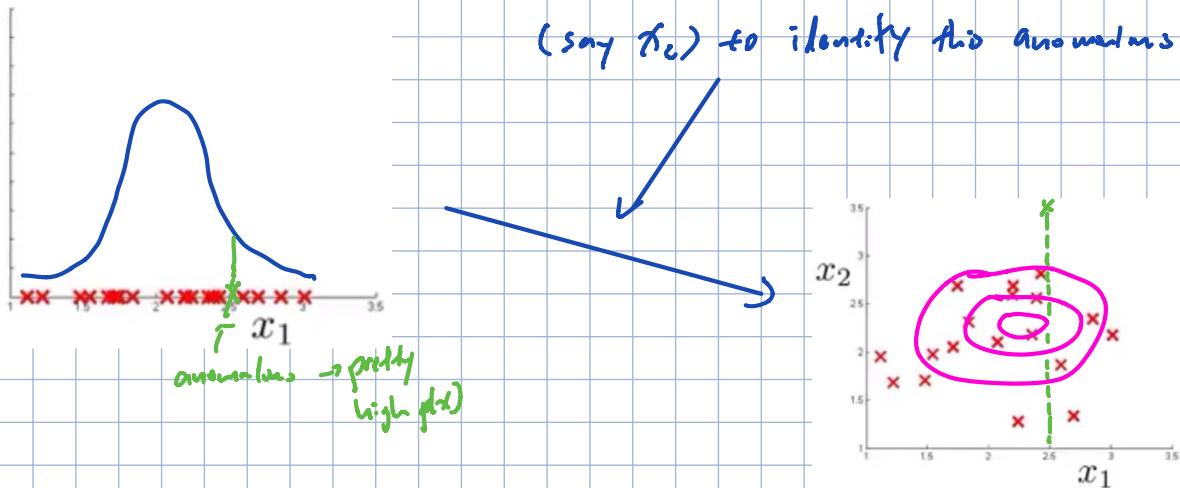
Goal: want  $p(x)$  large for normal examples  $\pi$ .  
 $p(x)$  small for anomalous examples  $\pi$ .

Most common problem:

$p(x)$  is comparable (say, both large) for normal and anomalous examples.

e.g.

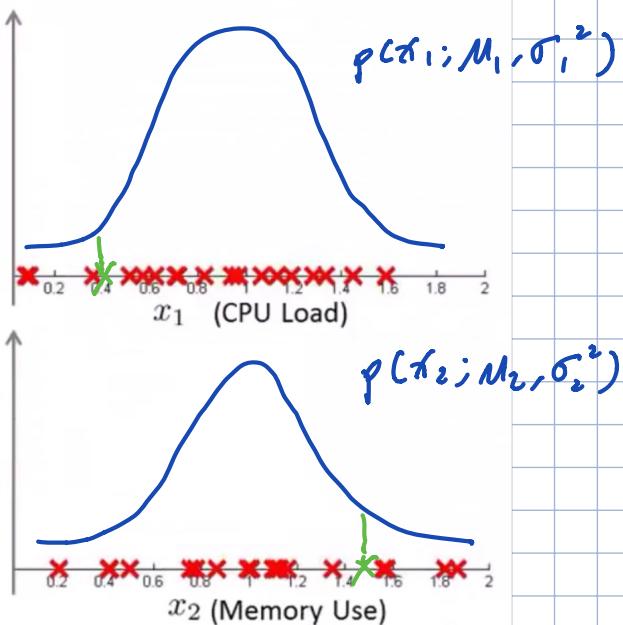
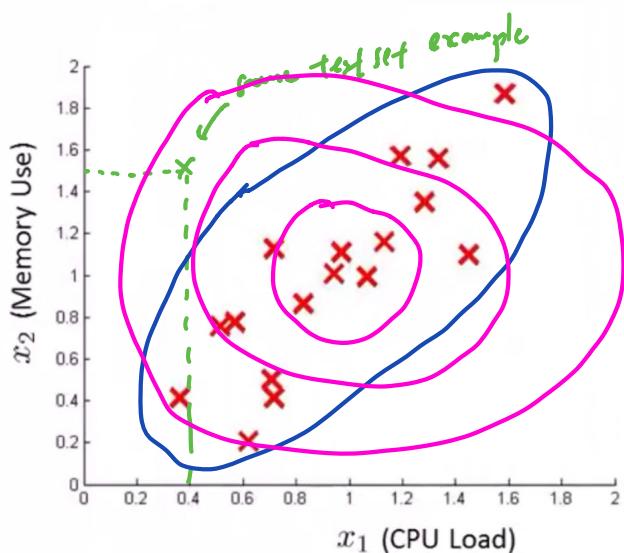
⇒ look at this anomalous example, construct a new feature (say  $x_0$ ) to identify this anomalous example



\* Choose features that might take on unusually large or small in the event of anomaly.

## Multivariate Gaussian Distribution

### Motivating example: Monitoring machines in a data center



this test set example does not look that bad.

### Multivariate Gaussian (Normal) distribution

$$X \in \mathbb{R}^n$$

Don't model  $p(x_1), p(x_2), \dots$ , etc. separately.

Model  $p(X)$  all in one go.

Parameters:  $\mu \in \mathbb{R}^n$ ,  $\Sigma \in \mathbb{R}^{n \times n}$  (covariance matrix)

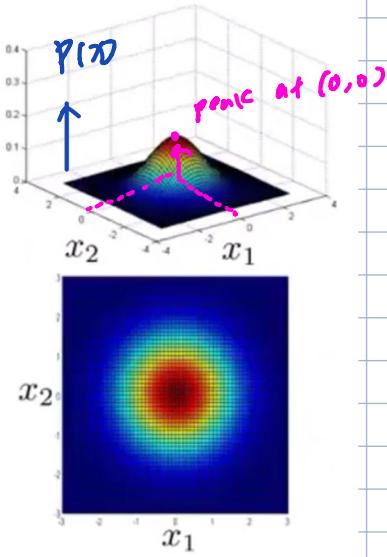
$$p(X; \mu, \Sigma) =$$

$$\frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (X - \mu)^T \Sigma^{-1} (X - \mu)\right)$$

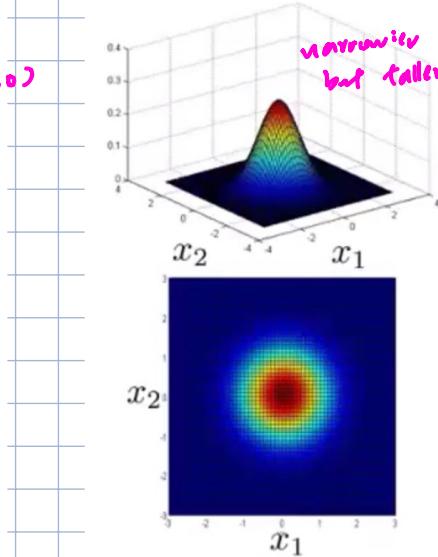
$\downarrow$   
 determinant  
 of  $\Sigma$

## Multivariate Gaussian (Normal) examples

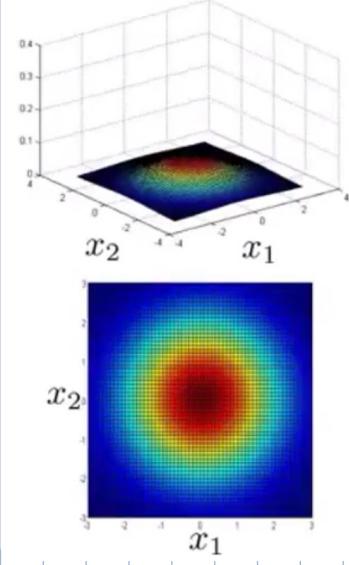
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



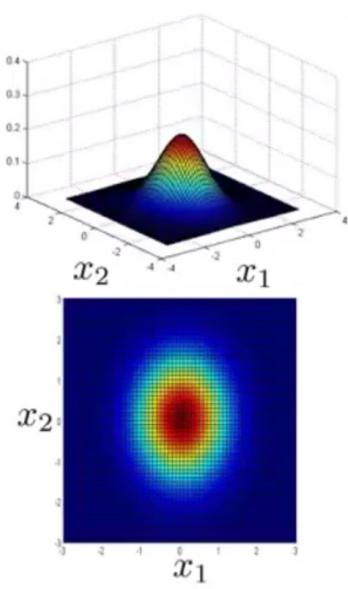
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 0.6 \end{bmatrix}$$



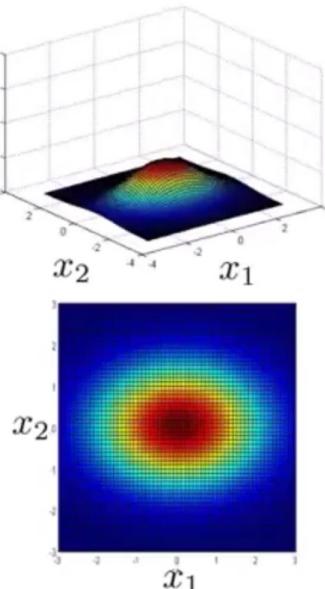
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$



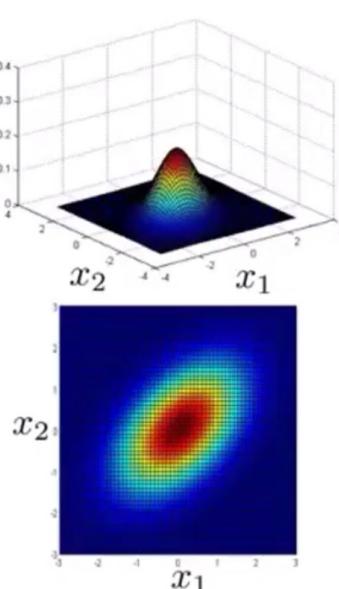
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 1 \end{bmatrix}$$



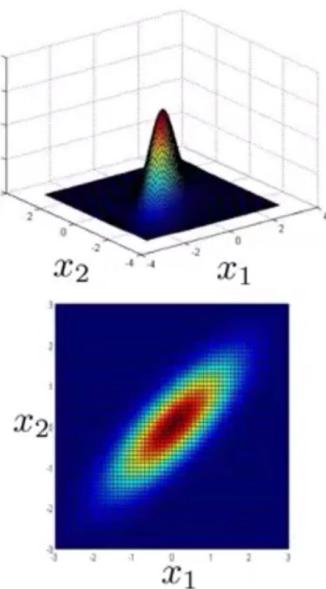
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$$



$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$

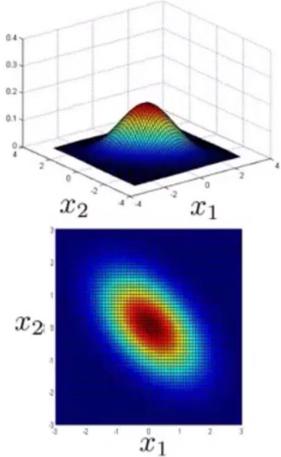


$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$

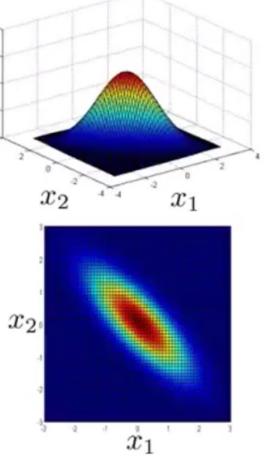


can be modeled  
for correlation  
between  $x_1$  and  $x_2$

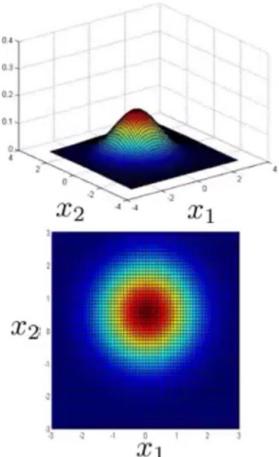
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & -0.5 \\ 0.5 & 1 \end{bmatrix}$$



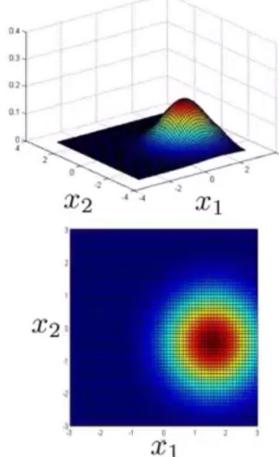
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & -0.8 \\ -0.8 & 1 \end{bmatrix}$$



$$\mu = \begin{bmatrix} 0 \\ 0.5 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



$$\mu = \begin{bmatrix} 1.5 \\ -0.5 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



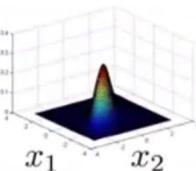
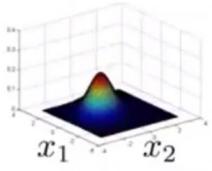
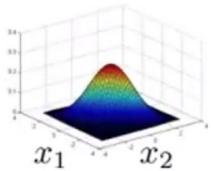
### Anomaly Detection: multivariate Gaussian Distribution

#### Multivariate Gaussian (Normal) distribution

Parameters  $\mu, \Sigma$

$$\mu \in \mathbb{R}^n, \Sigma \in \mathbb{R}^{n \times n}$$

$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left( -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right)$$



Parameter fitting:

Given training set:  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\} \quad x \in \mathbb{R}^n$

If I think my examples comes from a multivariate Gaussian Distribution, how do we estimate  $\mu$  and  $\Sigma$

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)(x^{(i)} - \mu)^T$$

## Anomaly detection with the multivariate Gaussian distribution algorithm:

① Fit model  $p(x)$  by setting

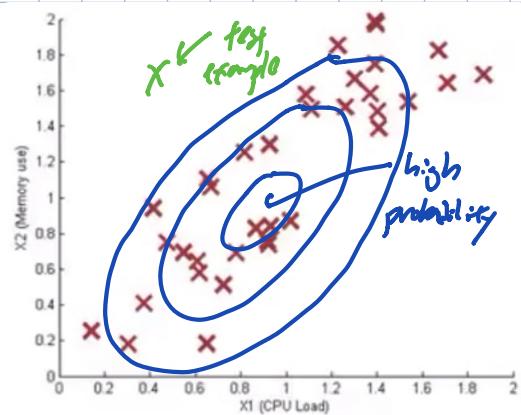
$$M = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - M)(x^{(i)} - M)^T$$

② Given a new example  $x$ , compute

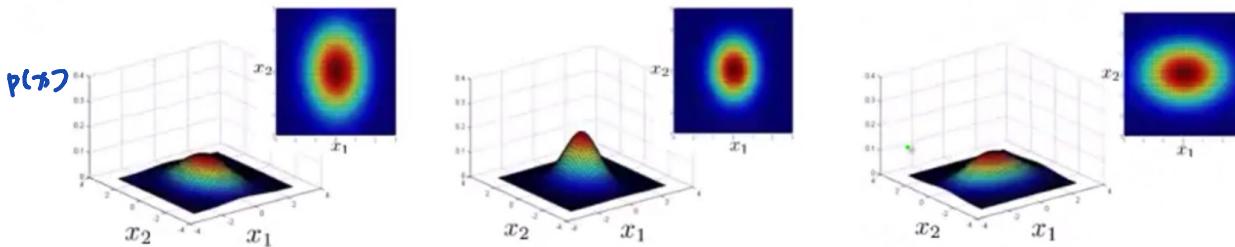
$$p(x) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - M)^T \Sigma^{-1} (x - M)\right)$$

Flag anomaly if  $p(x) < \epsilon$



## Relationship to original model

Original model:  $p(x) = p(x_1; \mu_1, \sigma_1^2) \times p(x_2; \mu_2, \sigma_2^2) \times \cdots \times p(x_n; \mu_n, \sigma_n^2)$



corresponds to a special case of multivariate Gaussian distribution,

by constraining  $p(x)$  so that the contour of  $p(x)$  on the axes is aligned.



i.e. no

$$p(x; M, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - M)^T \Sigma^{-1} (x - M)\right)$$

$$\text{where } \Sigma = \begin{bmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_n^2 \end{bmatrix}$$

(i.e. no correlation between features)

## original (Gaussian) model

$$p(x_1; \mu_1, \sigma_1^2) \times \dots \times p(x_n; \mu_n, \sigma_n^2)$$

- Manually create features to capture anomalies where  $x_1, x_2$  take unusual combinations of values.
- Computationally cheaper (scales better to large  $n$ ).
- OK even if  $m$  (training set size) is small.

## Multivariate Gaussian

$$p(x; \mu, \Sigma)$$

- Automatically captures correlation between features.
- Computationally more expensive (e.g.  $\Sigma^{-1}$  is expensive for large  $n$ )
- must have  $m \geq n$ , or else  $\Sigma$  is non-invertible.  
 ↳ or redundant features.  
 (e.g.  $x_1 = x_2$  or  $x_3 = x_4 + x_5$ )

## Recommender System: Problem formulation

### Example: Predicting movie ratings

User rates movies using one to five stars <sup>zero</sup>



Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)
Love at last	5	5	0	0
Romance forever	5	?	?	0
Cute puppies of love	?	4	0	?
Nonstop car chases	0	0	5	4
Swords vs. karate	0	0	5	?

$n_u = \text{no. of users}$

$n_m = \text{no. of movies}$

so.  $n_u = 4$

$n_m = 5$  here

$r(i, j) = 1$  if user  $j$  has rated movie  $i$

$y(i, j)$  = rating given by user  $j$  to movie  $i$   
 (defined only if  $r(i, j) = 1$ )

Given  $r(i,j)$  and  $\theta^{(i,j)}$ , can we predict the "?" rating?

## Recommender systems: Content Based Recommendation

Movie	Alice (1) $\theta^{(1)}$	Bob (2) $\theta^{(2)}$	Carol (3) $\theta^{(3)}$	Dave (4) $\theta^{(4)}$	$x_1$ (romance)	$x_2$ (action)
1 Love at last	5	5	0	0	0.9	0
2 Romance forever	5	?	?	0	1.0	0.01
3 Cute puppies of love	?	4	0	?	0.99	0
4 Nonstop car chases	0	0	5	4	0.1	1.0
5 Swords vs. karate	0	0	5	?	0	0.9

$x_1$ : a feature that measures the degree to which a movie is a romantic movie.

$x_2$ : a feature that measures the degree to which a movie is an action movie.

$$\theta_0 = 1$$

then we can have a movie vector

i.e.  $\boldsymbol{x}^{(1)} = \begin{bmatrix} 1 \\ 0.9 \\ 0 \end{bmatrix}$  for the first movie

In order to make predictions, we can train for each user  $j$ , learn a parameter  $\theta^{(j)} \in \mathbb{R}^3$ . (i.e. separate linear regression problem)  $\theta^{(5)} \in \mathbb{R}^{n+1}$

Predict user  $j$  as rating movie  $i$  with  $(\theta^{(j)})^T \boldsymbol{x}^{(i)}$  stars.

e.g.  $\boldsymbol{x}^{(3)} = \begin{bmatrix} 1 \\ 0.99 \\ 0 \end{bmatrix}$  for "Cute puppies of love"

we want to predict what Alice would rate this movie

say we have  $\theta^{(1)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}$  after training, then our prediction will be

$$(\theta^{(1)})^T \boldsymbol{x}^{(3)} = 4.95$$

## Problem formulation

More formally:

$r(i,j) = 1$  if user  $j$  has rated movie  $i$  (0 otherwise)

$y^{(i,j)}$  = rating by user  $j$  on movie  $i$  (if defined)

$\theta^{(j)}$  = parameter vector for user  $j$

$x^{(i)}$  = feature vector for movie  $i$

For user  $j$ , movie  $i$ , predicted rating:  $(\theta^{(j)})^T(x^{(i)})$

$m^{(j)}$  = no. of movie rated by user  $j$

To learn  $\theta^{(j)}$ :

$$\min_{\theta^{(j)}} \frac{1}{2m^{(j)}} \sum_{i: r(i,j)=1} ((\theta^{(j)})^T(x^{(i)}) - y^{(i,j)})^2 + \frac{\lambda}{2m^{(j)}} \sum_{k=1}^n (\theta_k^{(j)})^2$$

⚡  
 sum all  $i$   
 that  $r(i,j)=1$   
 i.e.  $j$  user  
 has rated

$\theta^{(j)} \in \mathbb{R}^{n+1}$

\* For easier implementation on Recommender

To learn  $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(m)}$ :

$$\min_{\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(m)}} \frac{1}{2} \sum_{j=1}^m \sum_{i: r(i,j)=1} ((\theta^{(j)})^T(x^{(i)}) - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^m \sum_{k=1}^n (\theta_k^{(j)})^2$$

⚡ obtain separate  $\theta$  for each user

## Optimization objective:

To learn  $\theta^{(j)}$  (parameter for user  $j$ ):

$$\min_{\theta^{(j)}} \frac{1}{2} \sum_{i:r(i,j)=1} \left( (\theta^{(j)})^T x^{(i)} - y^{(i,j)} \right)^2 + \frac{\lambda}{2} \sum_{k=1}^n (\theta_k^{(j)})^2$$

To learn  $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(n_u)}$ :

$$\min_{\theta^{(1)}, \dots, \theta^{(n_u)}} \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i:r(i,j)=1} \left( (\theta^{(j)})^T x^{(i)} - y^{(i,j)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2$$

## Optimization algorithm:

$$\min_{\theta^{(1)}, \dots, \theta^{(n_u)}} \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i:r(i,j)=1} \left( (\theta^{(j)})^T x^{(i)} - y^{(i,j)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2$$

$J(\theta^{(1)}, \dots, \theta^{(n_u)})$

Gradient descent update:

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} \quad (\text{for } k = 0)$$

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \left( \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \theta_k^{(j)} \right) \quad (\text{for } k \neq 0)$$

$$\frac{\partial}{\partial \theta_k^{(j)}} J(\theta^{(1)}, \dots, \theta^{(n_u)})$$

This content base approach assumes that we have available to us features for the different movie (i.e. is this action movie?)

## Recommender system: Collaborative filtering

Problem motivation:

sometimes we do not have feature about movie content

we change the problem a bit and suppose we do not know of these value about movie content.

Movie	Alice (1) $\theta^{(1)}$	Bob (2) $\theta^{(2)}$	Carol (3) $\theta^{(3)}$	Dave (4) $\theta^{(4)}$	$x_1$ (romance)	$x_2$ (action)
$\pi^{(1)}$ Love at last	5	5	0	0	?	?
Romance forever	5	?	?	0	?	?
Cute puppies of love	?	4	0	?	?	?
Nonstop car chases	0	0	5	4	?	?
Swords vs. karate	0	0	5	?	?	?

lets say we've gone to each of our users, and each of our user has told us how much they like romantic / action movie ( $\pi_0=1$ )

$$\theta^{(1)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}, \theta^{(2)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}, \theta^{(3)} = \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix}, \theta^{(4)} = \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix}$$

i.e. each user tells us  $\theta^{(i)}$

$\hookrightarrow$  it is possible to infer  $x_1$  and  $x_2$  from each movie

Mathematically :

what feature vector should  $\pi^{(1)}$  be, so that  $(\theta^{(1)})^T \pi^{(1)} \approx 5$

$$\pi^{(1)} = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}$$

$(\theta^{(1)})^T \pi^{(1)} \approx 5$   
 $(\theta^{(2)})^T \pi^{(1)} \approx 0$   
 $(\theta^{(3)})^T \pi^{(1)} \approx 0$   
 $(\theta^{(4)})^T \pi^{(1)} \approx 0$

## Optimisation algorithm

→ Given  $\theta^{(1)}, \dots, \theta^{(n_m)}$ , to learn  $\pi^{(i)}$ :

$$\min_{\pi^{(i)}} \frac{1}{2} \sum_{j:r(i,j)=1} ((\theta^{(j)})^T \pi^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{k=1}^n (\pi_k^{(i)})^2$$

→ Given  $\theta^{(1)}, \dots, \theta^{(n_m)}$ , to learn  $\pi^{(1)}, \dots, \pi^{(n_m)}$ :

$$\min_{\pi^{(1)}, \dots, \pi^{(n_m)}} \frac{1}{2} \sum_{i=1}^{n_m} \sum_{j:r(i,j)=1} ((\theta^{(j)})^T \pi^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (\pi_k^{(i)})^2$$

## Collaborative Filtering

Given  $\pi^{(1)}, \dots, \pi^{(n_m)}$  (and movie rating)

can estimate  $\theta^{(1)}, \dots, \theta^{(n_m)}$

Given  $\theta^{(1)}, \dots, \theta^{(n_m)}$ ,

can estimate  $\pi^{(1)}, \dots, \pi^{(n_m)}$

We can randomly guess some  $\theta$ , then get  $\pi$ , then estimate better  $\theta$ .

# The collaborative filtering algorithm

## Collaborative filtering optimization objective

→ Given  $x^{(1)}, \dots, x^{(n_m)}$ , estimate  $\theta^{(1)}, \dots, \theta^{(n_u)}$ :

$$\min_{\theta^{(1)}, \dots, \theta^{(n_u)}} \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2$$

→ Given  $\theta^{(1)}, \dots, \theta^{(n_u)}$ , estimate  $x^{(1)}, \dots, x^{(n_m)}$ :

$$\min_{x^{(1)}, \dots, x^{(n_m)}} \frac{1}{2} \sum_{i=1}^{n_m} \sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2$$

Minimizing  $x^{(1)}, \dots, x^{(n_m)}$  and  $\theta^{(1)}, \dots, \theta^{(n_u)}$  simultaneously:

$$J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}) = \frac{1}{2} \sum_{(i,j):r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2$$

instead of iteratively go from  $\theta \rightarrow x \rightarrow \theta \rightarrow x \dots$   
more efficient algorithm

Andrew Ng

## Collaborative filtering algorithm

→ 1. Initialize  $x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}$  to small random values.

→ 2. Minimize  $J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)})$  using gradient descent (or an advanced optimization algorithm). E.g. for every  $j = 1, \dots, n_u, i = 1, \dots, n_m$ :

$$x_k^{(i)} := x_k^{(i)} - \alpha \left( \sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) \theta_k^{(j)} + \lambda x_k^{(i)} \right)$$

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \left( \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \theta_k^{(j)} \right)$$

3. For a user with parameters  $\theta$  and a movie with (learned) features  $x$ , predict a star rating of  $\theta^T x$ .

$$(\theta^{(j)})^T x^{(i)}$$

user j on movie i

# Recommender System Vectorisation: Low Rank Matrix Factorisation

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	
Love at last	5	5	0	0	
Romance forever	5	?	?	0	
Cute puppies of love	?	4	0	?	
Nonstop car chases	0	0	5	4	$y_{ci,j}$
Swords vs. karate	0	0	5	?	

$Y = \begin{bmatrix} 5 & 5 & 0 & 0 \\ 5 & ? & ? & 0 \\ ? & 4 & 0 & ? \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & 0 \end{bmatrix}$

$n_m = 5$        $\uparrow 5 \times 4$

$n_u = 4$

Predicted ratings:

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 \\ 5 & ? & ? & 0 \\ ? & 4 & 0 & ? \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & 0 \end{bmatrix} \quad \begin{bmatrix} (\theta^{(1)})^T(x^{(1)}) & (\theta^{(2)})^T(x^{(1)}) & \dots & (\theta^{(n_u)})^T(x^{(1)}) \\ (\theta^{(1)})^T(x^{(2)}) & (\theta^{(2)})^T(x^{(2)}) & \dots & (\theta^{(n_u)})^T(x^{(2)}) \\ \vdots & \vdots & \vdots & \vdots \\ (\theta^{(1)})^T(x^{(n_m)}) & (\theta^{(2)})^T(x^{(n_m)}) & \dots & (\theta^{(n_u)})^T(x^{(n_m)}) \end{bmatrix}$$

$$\underset{\uparrow}{(\theta^{(j)})^T(x^{(i)})} \quad \underset{\uparrow}{X(H)^T}$$

Vectorisation:

$$X = \left[ \begin{array}{c} \cdots (x^{(1)})^T \cdots \\ \cdots (x^{(2)})^T \cdots \\ \vdots \\ \cdots (x^{(n_m)})^T \cdots \end{array} \right] \quad (H) = \left[ \begin{array}{c} \cdots (\theta^{(1)})^T \cdots \\ \cdots (\theta^{(2)})^T \cdots \\ \vdots \\ \cdots (\theta^{(n_u)})^T \cdots \end{array} \right]$$

Also called Low Rank Matrix factorization algorithm.

Additionally, we can use the learned features in order to find related movies.

For each product  $i$ , we learn a feature vector  $x^{(i)} \in \mathbb{R}^n$

e.g.  $x_1 = \text{romance}$ ,  $x_2 = \text{action}$ , ...

How to find movie  $j$  related to movie  $i$ ?

Small  $\|x^{(i)} - x^{(j)}\| \rightarrow$  movie  $i$  and  $j$  are similar

↳ 5 most similar movies to movie  $i$ :

↳ 5 movie  $j$  with the smallest  $\|x^{(i)} - x^{(j)}\|$

### Implementation Details: Mean Normalisation

#### Users who have not rated any movies

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	Eve (5)
Love at last	5	5	0	0	?
Romance forever	5	?	?	0	?
Cute puppies of love	?	4	0	?	?
Nonstop car chases	0	0	5	4	?
Swords vs. karate	0	0	5	?	?

↖ unwatched any movie

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 & ? \\ 5 & ? & ? & 0 & ? \\ ? & 4 & 0 & ? & ? \\ 0 & 0 & 5 & 4 & ? \\ 0 & 0 & 5 & 0 & ? \end{bmatrix}$$

$$\min_{\substack{x^{(1)}, \dots, x^{(n_m)} \\ \theta^{(1)}, \dots, \theta^{(n_u)}}} \frac{1}{2} \sum_{(i,j): r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2$$

say  $n=2$ , i.e. 2 features: romance, action

$\theta^{(5)} \in \mathbb{R}^2$ , since Eve has not rated any movie, so no movies that  $r(i,j)=1$ .

↳ only last term has effect in determining  $\theta^{(5)}$ , i.e.

$$\frac{\lambda}{2} [(\theta_1^{(5)})^2 + (\theta_2^{(5)})^2] \rightarrow \theta^{(5)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

thus  $(\theta^{(5)})^T x^{(i)} = 0$  Not useful

## Mean Normalization:

$$Y = \begin{bmatrix} \rightarrow [5 & 5 & 0 & 0 & ?] \\ \rightarrow [5 & ? & ? & 0 & ?] \\ ? & 4 & 0 & ? & ? \\ 0 & 0 & 5 & 4 & ? \\ \rightarrow [0 & 0 & 5 & 0 & ?] \end{bmatrix}$$

$$\mu = \begin{bmatrix} 2.5 \\ 2.5 \\ 2 \\ 2.25 \\ 1.25 \end{bmatrix} \rightarrow Y = \begin{bmatrix} 2.5 & 2.5 & -2.5 & -2.5 & ? \\ 2.5 & ? & ? & -2.5 & ? \\ ? & 2 & -2 & ? & ? \\ -2.25 & -2.25 & 2.75 & 1.75 & ? \\ -1.25 & -1.25 & 3.75 & -1.25 & ? \end{bmatrix}$$

$\uparrow$   
average rating  
for each movie

$\uparrow$   
Subtract the mean  
 $\downarrow$   
Learn  $\theta^{(0)}, x^{(0)}$

For user  $j$ , on movie  $i$ , predict:

$$(\theta^{(j)})^T (\underbrace{x^{(i)}}_{\text{compute back}}) + \mu_j$$

User 5 (env):

$$\theta^{(5)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$(\theta^{(5)})^T (\underbrace{x^{(i)}}_0) + \mu_5$$

i.e. if we never rates any movies - just use the mean.

## week 10

### Learning with large datasets

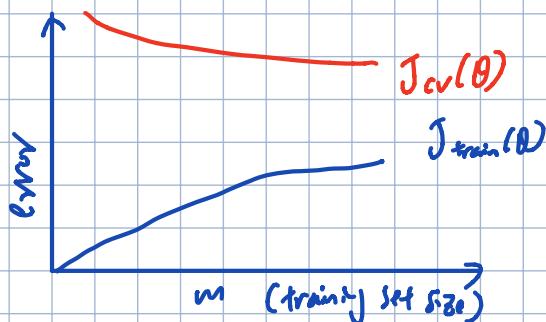
Says,  $m = 100,000,000$  (training samples)

Gradient descent for linear / logistic regression:

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

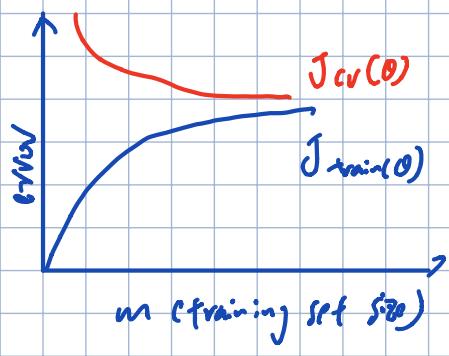
$\alpha \approx 100 \cdot 10^{-1}$

For sanity check, if  $m$  is small is also good, plot learning curve



high variance learning problem.

increase in weight help



high bias problem

increase in doesn't help much

should add extra feature, or extra hidden layer.

## Stochastic Gradient Descent

Gradient descent becomes computationally expensive when data set is big

### Linear regression with gradient descent

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

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

Gradient Descent:

Repeat {

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

↑ expensive  
if m is large

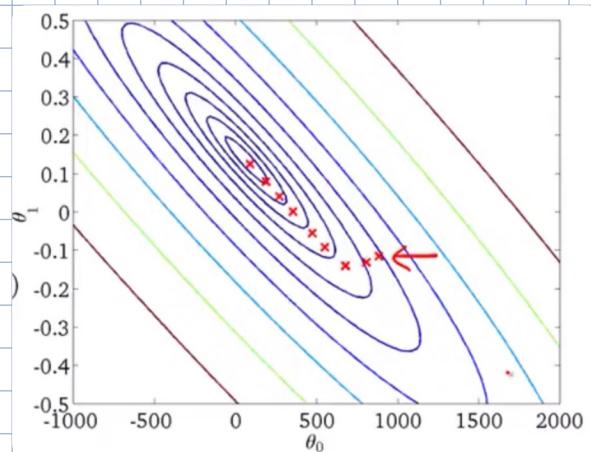
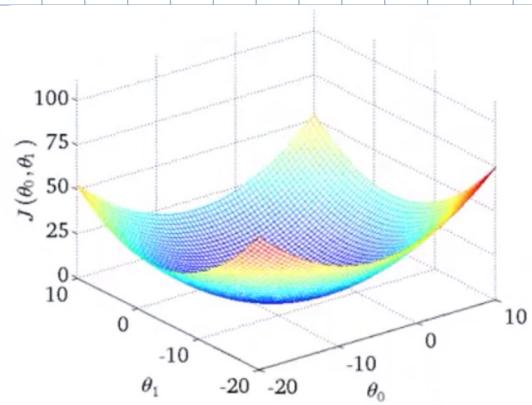
(for every  $j = 0, \dots, n$ )

$$\frac{\partial}{\partial \theta_j} J_{\text{train}}(\theta)$$

}

↳ this is called batch gradient descent.

↳ looking at all of the training examples.



## 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 the dataset

②. Repeat {

for  $i = 1, \dots, m$  }

maybe

1-to  
times

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

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

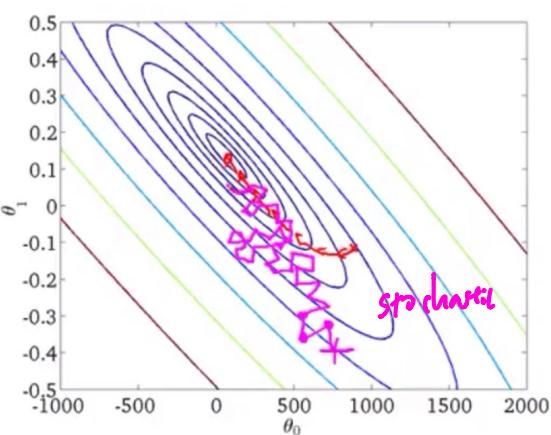
}

}

i.e. the algorithm will only look at  $(x^{(i)}, y^{(i)})$ , and the cost of just this training example.

then fit the 2nd training example  $\rightarrow$  then continue on for the entire training set.

this approach modifies the parameter list by list, rather than in one go.



Stochastic approach might not reach the global minimum in the same sense as batch.

It will also wandering around the minimum instead of converge to minimum.

on a single example

## Mini-batch Gradient Descent

so far:

Batch gradient descent: Use all m examples in each iteration.

Stochastic gradient descent: use 1 example in each iteration.

mini-batch gradient descent: use b examples in each iteration.

b = mini-batch size (e.g. b=10)

e.g.

Get  $b=10$  examples  $(x^{(i)}, y^{(i)}), \dots, (x^{(i+9)}, y^{(i+9)})$

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

$$i := i + 10$$

## Mini-batch gradient descent

Say  $b = 10, m = 1000$ .

Repeat {

for  $i = 1, 11, 21, 31, \dots, 991$  {

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

(for every  $j = 0, \dots, n$ )

}

}

mini-batch is likely to overshoot  
stochastic only if we have  
poor regularisation implementation

## Stochastic Gradient Descent Convergence

check for convergence!

→ Batch gradient descent:

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

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

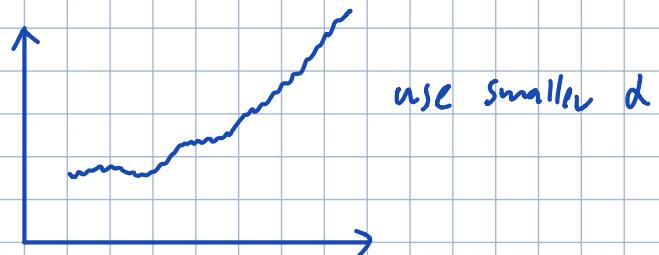
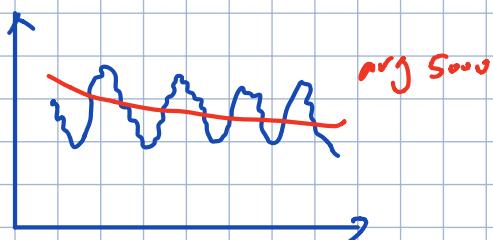
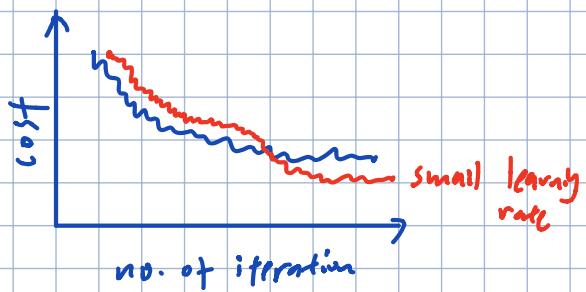
if  $m$  is small,  $J_{\text{train}}(\theta)$  can be computed rather efficiently.

→ Stochastic gradient descent:

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

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

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



the learning rate  $\alpha$  is typically held constant. Can slowly decrease  $\alpha$  over time if we want  $\theta$  to converge. (E.g.  $\alpha = \frac{\text{const1}}{\text{iterationNumber} + \text{const2}}$ )

∴ for stochastic,  $\theta$  may converge, just wander around global minimum

## Online Learning

Suppose we run a shipping service website where user comes, specifies origin and destination, you offer to ship their package for some asking price, and users sometimes choose to use your shipping service ( $y=1$ ), sometimes not ( $y=0$ )

Features to capture properties of user, of origin/destination and asking price. We want to learn  $p(y=1 | x; \theta)$  to optimise price.

{  
capture price.

For online learning and logistic regression:

Repeat forever     $\checkmark$  origin/destination/price rec.  
 $\checkmark$  do it

Get  $(x, y)$  corresponding to user.

Update  $\theta$  using  $(x, y) \leftarrow (x^{(t)}, y^{(t)})$ , i.e. just this example

$$\theta_j := \theta_j - 2(\text{log}(x) - y) \cdot x_j \quad (j=0, \dots, n)$$

Look at 1 example at a time, and discards after using.

Can adapt to changing user preference.

other online learning examples:

Product Search (learning to search)

User searches for "Android phone 1080p camera"

Have 100 phones in store. Will return 10 results.

For each phone and given a specific user query, we can construct a feature vector  $x_i$ .

$\chi$  = features of phone, how many words in user query match name of phone, how many words in query match description of phone, etc

$y=1$  if user click on link,  $y=0$  otherwise

learn  $p(y=1 | \chi; \theta)$

use to show user 10 phones they're most likely to click on.

### Map-Reduce and Data Parallelism

Say we want to fit a linear regression model or logistic regression model.

$$\text{Batch gradient descent: } \theta_j = \theta_j - \alpha \frac{1}{400} \sum_{i=1}^{400} (\text{log}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Machine 1: Use  $(x^{(1)}, y^{(1)}), \dots, (x^{(100)}, y^{(100)})$

$x^{(1)}, y^{(1)}$
:
$x^{(m)}, y^{(m)}$

$$\text{temp}_j^{(1)} = \sum_{i=1}^{100} (\text{log}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

Same

Machine 2: Use  $(x^{(101)}, y^{(101)}), \dots, (x^{(200)}, y^{(200)})$

$$\text{temp}_j^{(2)} = \sum_{i=101}^{200} (\text{log}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

Machine 3

Machine 4

$$\hookrightarrow \text{Combine: } \theta_j := \theta_j - \alpha \frac{1}{400} (\text{temp}_j^{(1)} + \text{temp}_j^{(2)} + \text{temp}_j^{(3)} + \text{temp}_j^{(4)})$$

$$(j=0, \dots, n)$$

$\rightarrow$  Many learning algorithms can be expressed as computing sums of functions over the training set.

E.g. for advanced optimization, with logistic regression, need:

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

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

week 11

## Photo OCR: problem description and pipeline

photo OCR: photo Optical character Recognition

### Photo OCR pipeline

1. Text detection



2. Character segmentation



3. Character classification

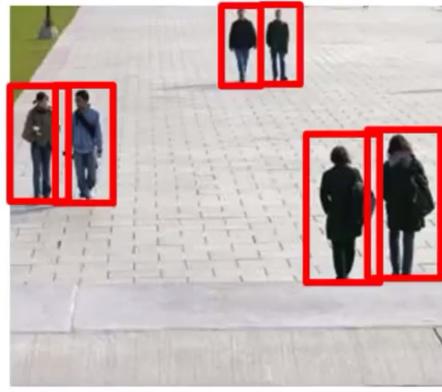


## photos OCR: sliding windows

Text detection



Pedestrian detection



### Supervised learning for pedestrian detection

$x$  = pixels in  $82 \times 36$  image patches  $\rightarrow$  aspect ratio

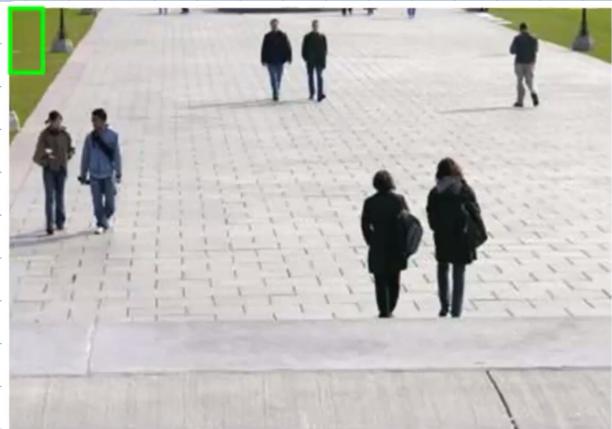


Positive examples ( $y = 1$ )

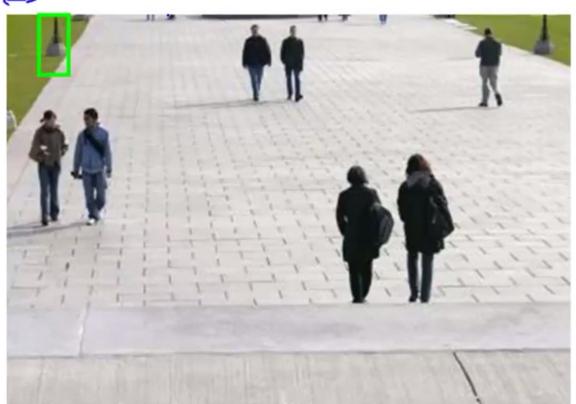


Negative examples ( $y = 0$ )

$\hookrightarrow$  collect large set of positive / negative examples



step-size / stride



Run the first patch ( $82 \times 36$  img) through classifier. (hopefully return  $y=1$ ).

Then slide the green box, and run that patch to the classifier again.

### Text detection



Positive examples ( $y = 1$ )



Negative examples ( $y = 0$ )

### Text detection



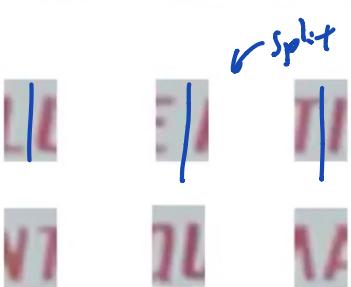
too thin box

expansion algorithm

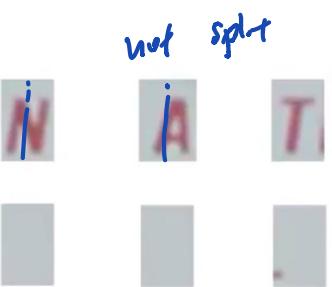
### 1D Sliding window for character segmentation



→ train classifier to  
split text



Positive examples ( $y = 1$ )



Negative examples ( $y = 0$ )

## Artificial Data Synthesis

### Artificial data synthesis for photo OCR



Real data

Abcdefg

*Abcdefg*

*Abcdefg*

Abcdefg

Abcdefg

we can take character of different font, and paste them against random backgrounds , to increase training examples.

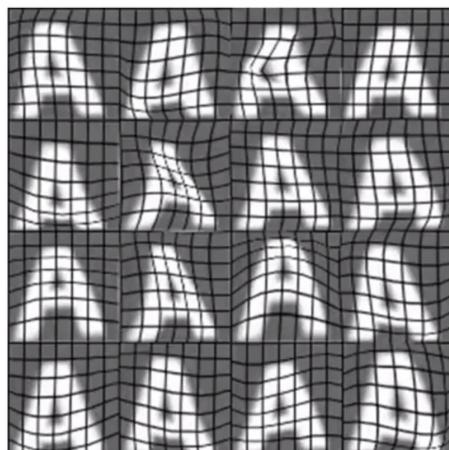
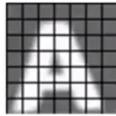


Real data



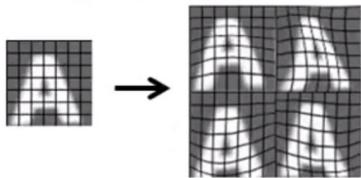
Synthetic data

### Synthesizing data by introducing distortions



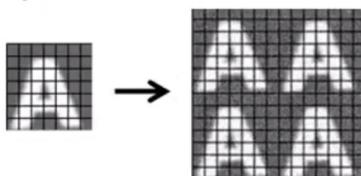
## Synthesizing data by introducing distortions

Distortion introduced should be representation of the type of noise/distortions in the test set.



Audio:  
Background noise,  
bad cellphone connection

Usually does not help to add purely random/meaningless noise to your data.



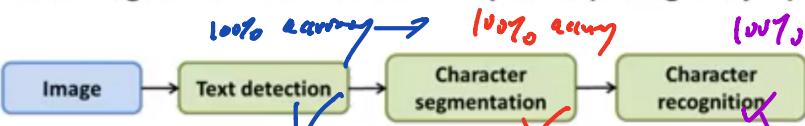
$x_i$  = intensity (brightness) of pixel  $i$   
 $x_i \leftarrow x_i + \text{random noise}$

## Discussion on getting more data

1. Make sure you have a low bias classifier before expending the effort. (Plot learning curves). E.g. keep increasing the number of features/number of hidden units in neural network until you have a low bias classifier.
2. "How much work would it be to get 10x as much data as we currently have?"
  - Artificial data synthesis
  - Collect/label it yourself
  - "Crowd source" (E.g. Amazon Mechanical Turk)

## Ceiling analysis

Estimating the errors due to each component (ceiling analysis)



What part of the pipeline should you spend the most time trying to improve?

component	accuracy
overall system	72%
Text detection	89%
char seg	90%
char rec	100%

## Summary: Main topics

$$(x^{(i)}, y^{(i)})$$

### → Supervised Learning

- Linear regression, logistic regression, neural networks, SVMs

$$x^{(i)}$$

### → Unsupervised Learning

- K-means, PCA, Anomaly detection

### → Special applications/special topics

- Recommender systems, large scale machine learning.

### → Advice on building a machine learning system

- Bias/variance, regularization; deciding what to work on next: evaluation of learning algorithms, learning curves, error analysis, ceiling analysis.