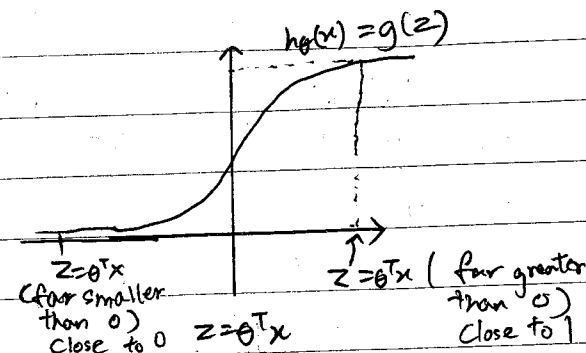


# Support Vector Machines

## Optimization objective

Alternative view of logistic regression

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



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

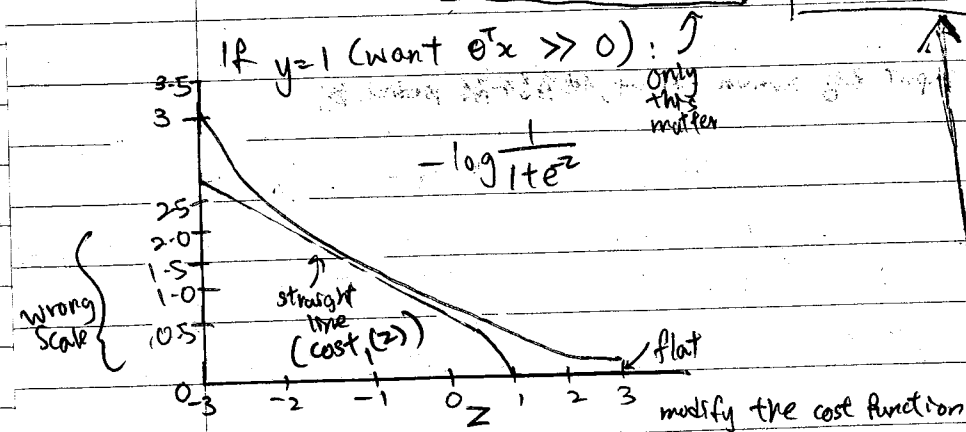
If  $y=0$ , we want  $h_{\theta}(x) \approx 0$ ,  $\theta^T x \ll 0$

Alternative view of logistic regression

single  $(x, y)$

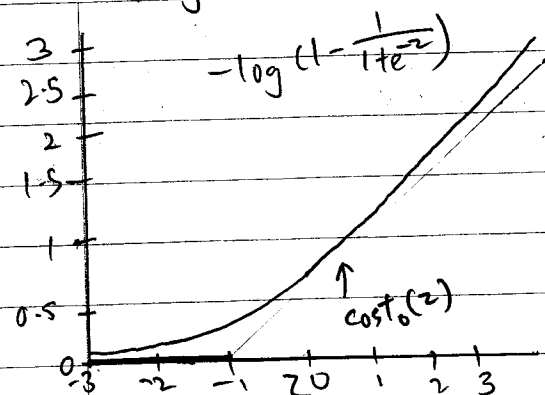
Cost of example:  $-(y \log h_{\theta}(x) + (1-y) \log(1-h_{\theta}(x)))$

$$= -y \log \frac{1}{1 + e^{-\theta^T x}} - (1-y) \log \left( 1 - \frac{1}{1 + e^{-\theta^T x}} \right)$$



when  $z \uparrow$ ,  $-\log \frac{1}{1+e^z} \downarrow$

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



## Support vector machine

Logistic regression,

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

(Note: The first term is labeled 'cost' and the second term is labeled 'B' in the original image)

Support vector machine:

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

$\Rightarrow$  cut out the constant in  $\min_u (u-5)^2 + 10 \rightarrow u=5$  (get smallest)

$\min_u 10(u-5)^2 + 10 \rightarrow u=5$  (still  $u=5$  to get smallest)

LP:

$A + \lambda B$   
in logistic regression,  $\lambda$  can act as a control to trade off the relative weight b/w how much we wanted the training set well.

So multiply something that you're minimizing over by some constant, 10 in this case, it does not change the value of  $u$  that gives us, that minimizes this function.

SVM:  $CA + B$  C playing a role as  $\frac{1}{\lambda}$

$C \downarrow, B > A, C$

overall optimization objective function for svm

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

if minimize this function, what you have is the parameters learned by SVM

SVM hypothesis

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

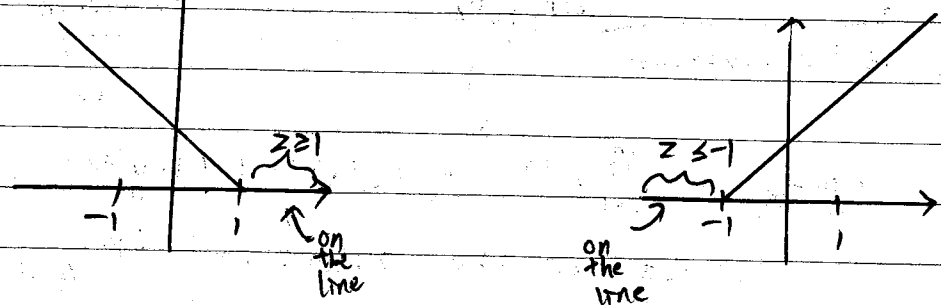
Hypothesis

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

## Large Margin Intuition

SVM:

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

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  $< 0$ )

→ This higher value to determine  $y=1/y=0$  builds in an extra safety factor or safety margin factor into SVM.

SVM Decision Boundary

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

if  $C$  very very large like 100,000 → we have to choose the value so this = 0

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

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

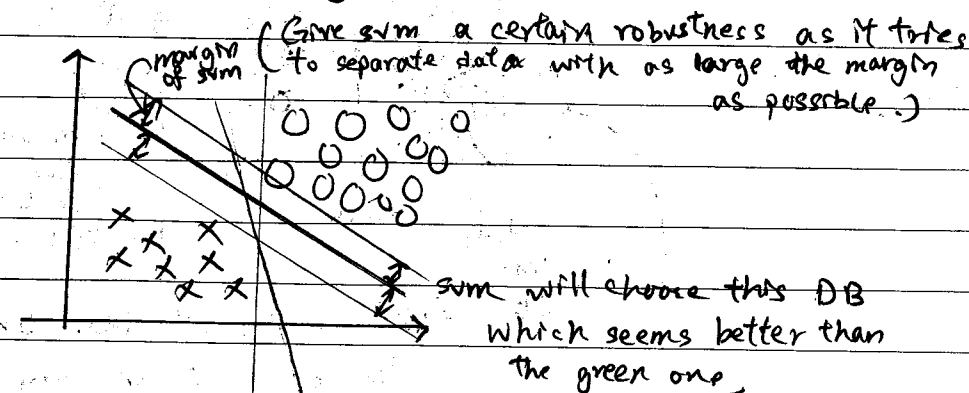
$$\min_{\theta} C \sum_{i=1}^m [y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)})] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

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

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

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

SVM Decision Boundary: Linearly separable case

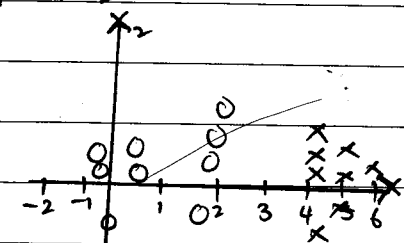


Linear DB that separate

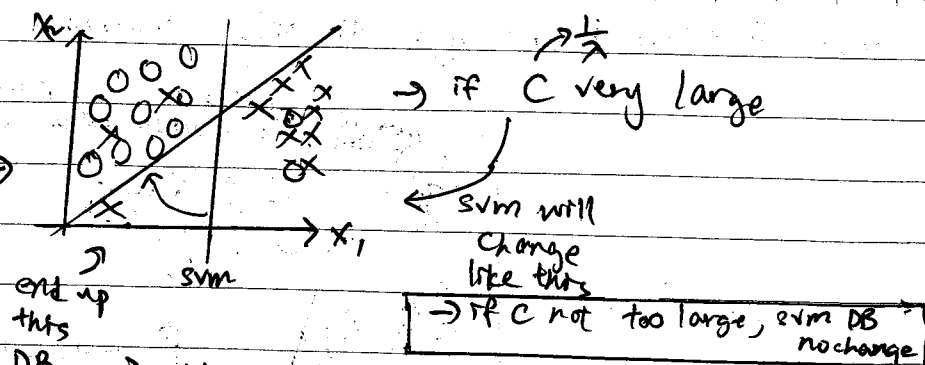
SVM sometimes called large margin classifier

this DB has a larger distance which is known as margin from any of the training examples, whereas green line is close to examples.

Quiz: Consider the training set to the right, where "x" denotes the examples ( $y=1$ ) and "o" denotes -ve examples ( $y=0$ ). Suppose you train an SVM (which will predict 1 when  $\theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$ ). What values might SVM give for  $\theta_0$ ,  $\theta_1$ , and  $\theta_2$ ?

Answer:  $\theta_0 = -3$ ,  $\theta_1 = 1$ ,  $\theta_2 = 0$

# Large margin classifier in presence of outliers



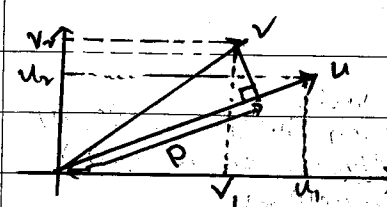
end up this DB  
 => It's really not clear based on the single outlier, and it's not a good idea to change DB.

if got training examples like green "o" and "x" appear, svm will also do the right thing

## Support vector machines

### The mathematics behind large margin classification

#### Vector Inner Product



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

$$\|u\| = \text{length of vector } u \quad u^T v \quad (\text{called as inner product b/w vectors } u \text{ and } v)$$

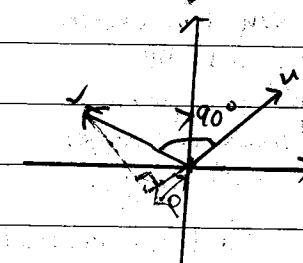
$$= \sqrt{u_1^2 + u_2^2} \in \mathbb{R}$$

$p$  = length/magnitude of the projection of vector  $v$  onto the vector  $u$

Signed (could be +ve or -ve)

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

$$= u_1 v_1 + u_2 v_2 \quad p \in \mathbb{R}$$



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

$$p < 0$$

\* In inner product, if angle b/w  $u$  and  $v$  is less than  $90^\circ$ , then  $p$  is +ve  
 if angle is greater than  $90^\circ$ ,  $p$  is -ve.

#### SVM Decision Boundary

$$\min \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

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

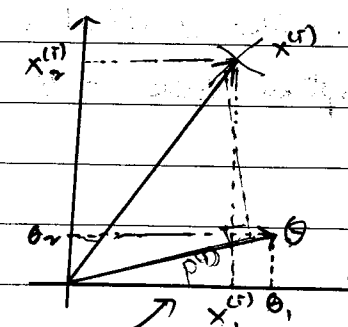
$$\theta^T x^{(i)} = ?$$

Simplification:  $\theta_0 = 0$

$n=2$

$$= \frac{1}{2} (\theta_1^2 + \theta_2^2) = \frac{1}{2} (\sqrt{\theta_1^2 + \theta_2^2})^2 = \frac{1}{2} \| \theta \|^2$$

$$\begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} \theta_0 = 0$$



Projection of the  $i$ -th training example onto parameter vector  $\theta$

$$\theta^T x^{(i)} = p^{(i)} \cdot \| \theta \|$$

$$= \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)}$$

=> continue at the back.

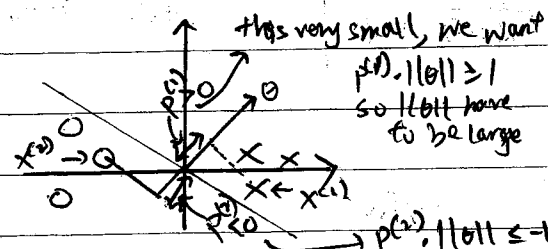
## SVM Decision Boundary

$$\min_{\theta} \frac{1}{2} \sum_{i=1}^n \theta_i^2 = \frac{1}{2} \|\theta\|^2$$

$$\text{s.t. } \left. \begin{array}{l} p^{(i)} \cdot \|\theta\| \geq 1 \text{ if } y^{(i)} = 1 \\ p^{(i)} \cdot \|\theta\| \leq -1 \text{ if } y^{(i)} = 0 \end{array} \right\} \leftarrow \text{very large}$$

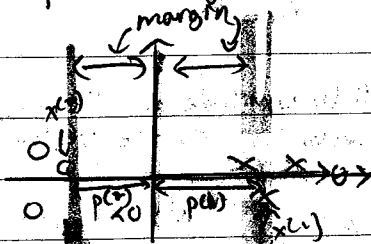
Where  $p^{(i)}$  is the projection of  $x^{(i)}$  onto the vector  $\theta$ .  
 (if  $\theta_0 \neq 0$ , DB can be not passing through origin  $(0,0)$ )  
 Simplification:  $\theta_0 = 0$  (means DB must pass through origin  $(0,0)$ )

Note:  
vector  $\theta$  should be  $90^\circ$  to DB



Let's say SVM choose green DB  
 → not a good choice, margin too small

But SVM will not do this  
 ⇒ In this,  $p^{(i)}$  will be pretty small numbers. And the  $\|\theta\|$  must be big



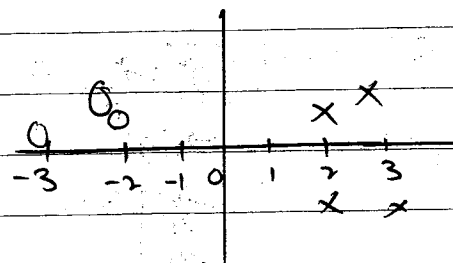
Let's say SVM choose this pink DB.  
 ⇒  $p^{(i)}$  and  $p^{(2)}$  have much more bigger

$$p^{(1)} \cdot \|\theta\| \geq 1$$

$\|\theta\|$  can be smaller

(SVM will choose this)  
 Margin is the value of  $p^{(i)}$ ,  
 $p^{(i)}$  bigger, margin bigger  
 and  $\|\theta\|$  can be smaller  
 which is what we trying to do

Quiz:



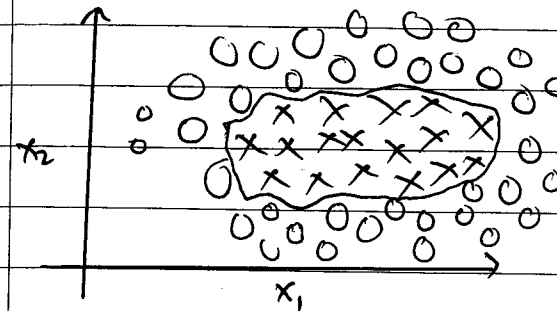
where  $p^{(i)}$  is the (signed +ve/-ve) projection of  $x^{(i)}$  onto  $\theta$ . Consider the training set above. At the optimal value of  $\theta$ . What is  $\|\theta\|$ ?

$$\text{Answer: } \frac{1}{2}$$

## Support Vector Machines

### Kernels 1

#### Non-linear Decision Boundary



Predict  $y=1$  if

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

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

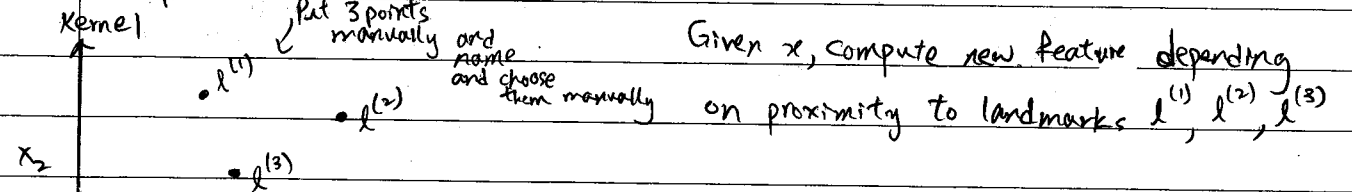
$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \dots$$

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

Is there a different/better choice of the features  $f_1, f_2, f_3, \dots$ ?

Idea on how to define new features  $f_1, f_2, f_3$

Example:



Given  $x$ :  $f_1 = \text{similarity}(x, l^{(1)})$  (1st feature)

measure of similarity between training example  $x$  and first landmark  $l^{(1)}$  =  $\exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$

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

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

Kernel function (Gaussian Kernels) ← specific example here is named this

$$\Rightarrow k(x, l^{(i)})$$

## Kernels and Similarity

$$f_i = \text{similarity}(x, l^{(i)}) = \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right) = \exp\left(-\frac{\sum_{j=1}^n (x_j - l_j^{(i)})^2}{2\sigma^2}\right)$$

If  $x \approx l^{(i)}$ :  $f_i \approx \exp\left(-\frac{0^2}{2\sigma^2}\right) \approx 1$

If  $x$  is far from  $l^{(i)}$ :

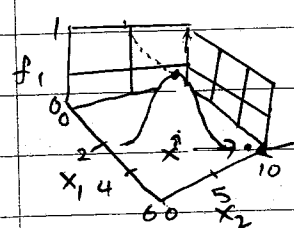
$$f_i = \exp\left(-\frac{(\text{large number})^2}{2\sigma^2}\right) \approx 0$$

Example: 2 features  $\rightarrow x_1, x_2$

landmark at location  $\begin{bmatrix} 3 \\ 5 \end{bmatrix}$   $l^{(1)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$ ,  $f_1 = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$

$\sigma^2$ : parameter of Gaussian Kernel

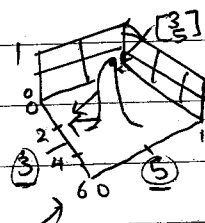
$$\sigma^2 = 1$$



when  $x = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$ ,  $f_1 = 1$

nice bump (凸起来的)

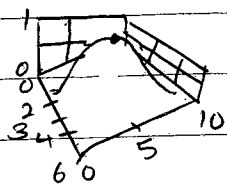
$$\sigma^2 = 0.5$$



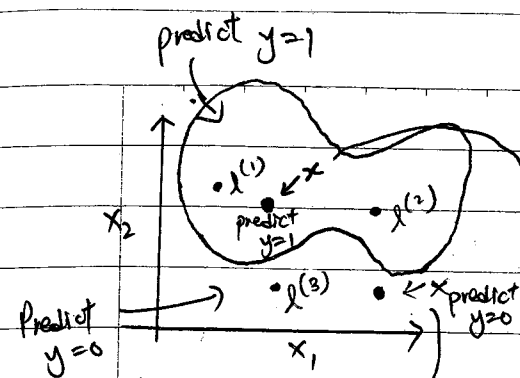
$f_1$  falls to zero much more rapidly

the bump becomes narrower

$$\sigma^2 = 3$$



as you move away from  $l^{(i)}$ , value of feature falls away much more slowly



Predict "1" when

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

Put  $x$  as

Let's say

$$\theta_0 = -0.5, \theta_1 = 1, \theta_2 = 1, \theta_3 = 0$$

Because  $x$  is near to  $l^{(1)}$ , but far from  $l^{(2)}$  and  $l^{(3)}$ ,

$$\text{so, } f_1 \approx 1, f_2 \approx 0, f_3 \approx 0$$

not exactly 1 or 0, just close to

$$\theta_0 + \theta_1 x_0 + \theta_2 x_0 + \theta_3 x_0$$

$$= \theta_0$$

$$\approx -0.5 < 0$$

Because  $15 < 0$ , for this training example  $x$ , predict  $y=0$

$$\theta_0 + \theta_1 x_1 + \theta_2 x_0 + \theta_3 x_0$$

$$= -0.5 + 1$$

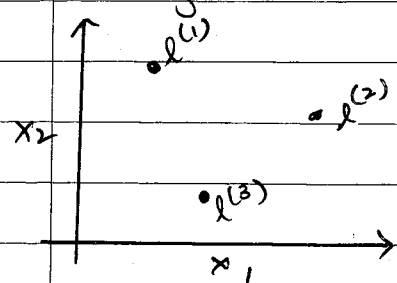
$$= 0.5 \geq 0$$

Because  $15 \geq 0$ , for training example  $x$ , predict  $y=1$

We end up predicting those example near to  $l^{(1)}, l^{(2)}$  as  $y=1$  and far from  $l^{(1)}, l^{(2)}$  as  $y=0$ . So we got a DB like this.

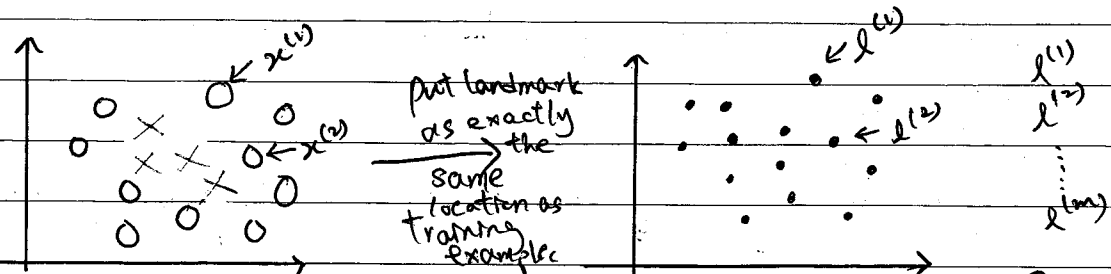
## Kernels II

Choosing the landmarks

Given  $x$ :

$$f_i = \text{similarity}(x, l^{(i)}) \\ = \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right) \leftarrow \text{Gaussian Kernel}$$

Predict  $y = 1$  if  $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$   
 where to get  $l^{(1)}, l^{(2)}, l^{(3)}, \dots$ ?



Features are going to measure how close an example is to one of the things in training set.

## SVM with kernels

Given  $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$ ,Choose  $l^{(1)} = x^{(1)}, l^{(2)} = x^{(2)}, \dots, l^{(m)} = x^{(m)}$ Given example  $x$ :  $\leftarrow$  can be in the training set / cv set / test set

$$f_1 = \text{similarity}(x, l^{(1)}) \\ f_2 = \text{similarity}(x, l^{(2)}) \\ \vdots \\ f_m = \text{similarity}(x, l^{(m)}) \Rightarrow f = \begin{bmatrix} f_0 \\ f_1 \\ f_2 \\ \vdots \\ f_m \end{bmatrix} \quad \begin{matrix} f_0 = 1 \text{ (x0, intercept)} \\ \text{(as a feature vector)} \end{matrix}$$

For training example  $(x^{(i)}, y^{(i)})$ :

Gaussian Kernel

$$x^{(i)} \rightarrow \begin{bmatrix} f_1^{(i)} \\ f_2^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix} = \begin{bmatrix} \text{sim}(x^{(i)}, l^{(1)}) \\ \text{sim}(x^{(i)}, l^{(2)}) \\ \vdots \\ \text{sim}(x^{(i)}, l^{(m)}) \end{bmatrix} \leftarrow f_i^{(i)} = \text{sim}(x^{(i)}, l^{(i)}) \quad \begin{matrix} \text{Gaussian Kernel} \\ = \exp\left(-\frac{0}{2\sigma^2}\right) = 1 \end{matrix}$$

$[f_i^{(i)}]$  is just be similarity b/w  $x$  and itself

Can take these and group them into a feature vector

new feature vector to represent training example

## SVM with kernels

Hypothesis: Given  $x$ , compute features  $f \in \mathbb{R}^{m+1}$   
 Predict " $y=1$ " if  $\theta^T f \geq 0$   $\leftarrow \theta_0 f_0 + \theta_1 f_1 + \dots + \theta_m f_m$

Training:

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

$\leftarrow \theta^T f^{(i)}$   $\leftarrow$  Not regularizing  $\theta_0$

$$\sum_{j=1}^n \theta_j^2 = \theta^T \theta \leftarrow \theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_m \end{bmatrix} \text{ (ignoring } \theta_0)$$

replace Sum Implementation  $\rightarrow$  instead of minimizing exactly  $\| \theta \|^2$

$$\text{with } \theta^T M \theta \quad \text{rescale version} \rightarrow \text{The parameter vector theta depends on kernel}$$

this  $\rightarrow$  Matrix (depends on kernel you use)

kind of mathematical detail allows svm software to run more efficiently

this modification  $\rightarrow$  allow to scale to much bigger training sets

Can use kernel idea on logistic regression, but will very slow. Because of computational tricks like that in body and how it modifies and details of how svm software is implemented, so svm and kernels tend to go particularly well together.

- Use the software/packages that developed by others is better than writing them yourself.

SVM parameters:

Bias and Variance Trade-off

$$C (= \frac{1}{\lambda})$$

Large  $C$ : Lower bias, high variance (small  $\lambda$ )

Small  $C$ : Higher bias, low variance (large  $\lambda$ )

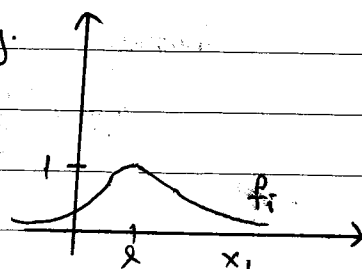
In Gaussian kernels:

$$\sigma^2$$

Large  $\sigma^2$ : Features  $f_i$  vary more smoothly.

Higher bias, lower variance.

$$\exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right)$$



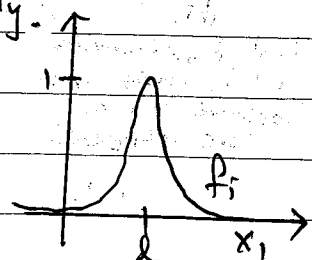
Gaussian kernel  
falls off smoothly,  
will get a hypothesis  
that varies slowly/  
varies smoothly  
as you change the  
input  $x$

\*

determine the width  
of the Gaussian kernel

Small  $\sigma^2$ : Features  $f_i$  vary less smoothly.

Lower bias, higher variance.



The features vary  
less smoothly as  
higher slopes/higher  
derivatives, end up  
fitting hypothesis  
of lower bias and  
higher variance

Quiz: Suppose you train an SVM and find it overfits your training data.

Which of these would be a reasonable next step?

- Decrease  $C$

- Increase  $\sigma^2$

## Support Vector Machines

Using an SVM

Use SVM software package (e.g. liblinear, libsvm, ...) to solve for parameters  $\theta$ .

Need to specify:

Choice of parameter  $C$

Choice of kernel (similarity function):

E.g. No kernel ("linear kernel")

Predict " $y=1$ " if  $\theta^T x \geq 0$  ( $\theta_0 + \theta_1 x_1 + \dots + \theta_n x_n \geq 0$ )

Why want to do this?

$n$  large,  $m$  small

( $\Rightarrow$ ) maybe want to just fit a linear DB  
and not try to fit a very complicated  
nonlinear function, because might not  
have enough data.

Gaussian kernel:

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

Need to choose  $\sigma^2$

When to use?

$x \in \mathbb{R}^n$ ,  $n$  small, and/or  $m$  large

and want to use kernel to fit more

kernel (similarity) functions:  $x^{(1)}, x^{(2)}, \dots, x^{(m)}$

function  $f = \text{kernel}(x_1, x_2)$

$$f = \exp\left(-\frac{\|x_1 - x_2\|^2}{2\sigma^2}\right)$$

return

Note: Do perform feature scaling before using Gaussian kernel.

$$\|x - l\|^2 \rightarrow v = x - l \quad x \in \mathbb{R}^n$$

$$\|v\|^2 = v_1^2 + v_2^2 + \dots + v_n^2$$

$$= (x_1 - l_1)^2 + (x_2 - l_2)^2 + \dots + (x_n - l_n)^2$$

Use housing problem

if this is 1000 ft<sup>2</sup>

this is #bedrooms (1-5)

The difference will  
very large, so need  
feature scaling

## Other choices of kernel

Note: Not all similarity functions  $\text{similarity}(x, l)$  make valid kernels. (Need to satisfy technical condition called "Mercer's Theorem" to make sure SVM packages' optimizations run correctly, and do not diverge).

Many off-the-shelf kernels available:

- Polynomial kernel:  $k(x, l) = (x^T l)^2$

one measure how similar  $x$  and  $l$  are. If  $x$  and  $l$  very close to each other, the inner product tend to be large.

Another version:  $(x^T l)^3, (x^T l + 1)^3, (x^T l + 5)^3$

$\Rightarrow$  has 2 parameters: degree of polynomial, number you add here

In general:  $(x^T l + \text{constant})^{\text{degree}}$

But almost always perform worse than Gaussian kernel

Usually used only for data where  $x$  and  $l$  are all strictly non negative, so that ensure these inner products are never -ve

- More esoteric: String kernel, chi-square kernel, histogram intersection kernel, ...

to measure similarity between different objects

use if your input data is text strings or other types of strings.

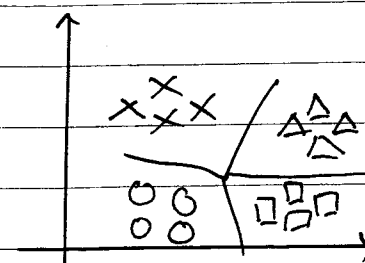
The more similar two strings/two objects are, the higher the value of  $K(a, b)$  will be.

## Quiz:

Suppose you are trying to decide among a few different choices of kernel and are also choosing parameters such as  $C, \sigma^2$ , etc. How should you make the choice?

Answer: Choose whatever performs best on cross-validation data.

## Multi-class classification



$y \in \{1, 2, 3, \dots, K\}$

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

## Logistic regression vs. SVMs

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

$\Rightarrow$  if  $n$  is large (relative to  $m$ ): (e.g.  $n \geq m$ ,  $n=10,000$ ,  $m=10 \dots 1000$ )

Use logistic regression, or SVM without a kernel ("linear kernel")

$\Rightarrow$  if  $n$  is small,  $m$  is intermediate: ( $n=1-1000$ ,  $m=10-10,000$ )

Use SVM with Gaussian kernel (Create/Add more polynomial features)

$\Rightarrow$  if  $n$  is small,  $m$  is large: ( $n=1-1000$ ,  $m=50,000$ ++)

Create/add more features, then use logistic regression or SVM without a kernel

logistic regression and SVM without a kernel usually do pretty similar things and give pretty similar performance, but depending on the implementational details, one may be more efficient than the other.

$\Rightarrow$  Neural Network likely to work well for most of these settings, but may be slower to train.

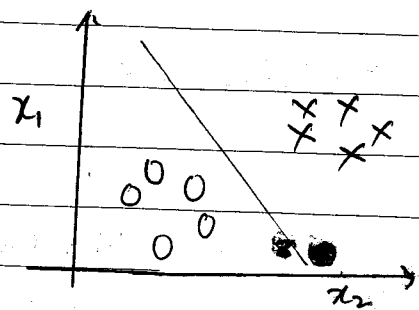
SVM has a convex optimization problem, and good SVM optimization software packages will always find the global min. or something close to it. So for SVM, no need to worry about local optima.



# Clustering

## Unsupervised learning introduction (learned from unlabeled data)

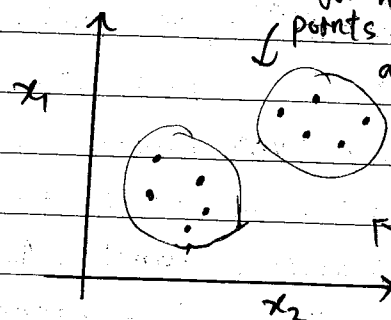
### Supervised learning



Training set =  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}), \dots, (x^{(n)}, y^{(n)})\}$

Goal: To find DB that separates the label examples and -ve label examples

### Unsupervised learning



algorithm find, looks like this data set has points grouped into 2 separate clusters, algorithm that find clusters like the circles one is called a clustering algorithm.

Data that has no labels

Training set =  $\{x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(n)}\}$  ← don't get label y

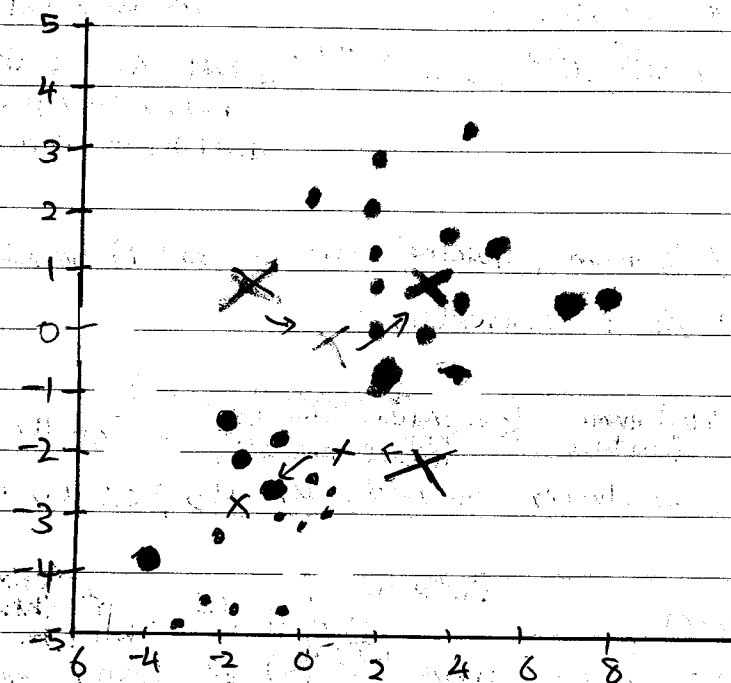
Give this sort of unlabeled training set to an algorithm and ask algorithm find some structure in the data for us.

### Application of clustering

- Market segmentation
- social network analysis
- Organize computing clusters
- Astronomical data analysis

# Clustering

K-means algorithm  $\Rightarrow$  to automatically group unlabeled data into coherent subsets/coherent clusters for us.



First, randomly initialize (X) two points, called cluster centroids.

Got 2 cluster centroids, because going to group data into 2 clusters.

K-Mean  $\Rightarrow$  iterative algorithm does 2 things — cluster assignment step — move centroid step

Going through each of the example (dots), depend which cluster centroid is closer to, black or blue, assign each of the data points to one of it.

(first thing) — Those blue belong to blue cluster centroid, while black belong to black cluster centroid

(second thing) — Look at the mean of blue dot, move the blue cluster centroid there, and look at the mean of black dot, move the black cluster centroid there. (arrow indicate moving)

(first thing again) — After moving the cluster centroid, do the cluster assignment step again, which colour them either blue/black, so the colour of some dots change again

(first thing again) — Move to the mean of new training data again, and recolour those data which are closer to either blue/black cluster centroid again.

Keep running additional iterations of K-mean until a final mean, it will not change any further, and had done finding 2 clusters in data.

# K-mean algorithm

Input:

- K (number of clusters)  $\leftarrow$  we want certain number of cluster, and we tell how many we want in data set.
- Training set  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$   $\leftarrow$  just Xs, as this is unsupervised learning, we don't have labels Y.

For unsupervised learning, of K-mean, training examples are n-dimensional rather than nt)  
 $x^{(i)} \in \mathbb{R}^n$  (drop  $x_0 = 1$  convention)

What it does

$K \Rightarrow$  total num. of centroids  $\Rightarrow$  index (between 1 to K) into different centroid location of the cross on last page

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

Repeat {

Cluster assignment step that colour the points  $\leftarrow$  for  $i=1$  to  $m$   
 $c^{(i)} := \text{index (from 1 to K) of cluster centroid closest to } x^{(i)}$   $\leftarrow$  Distance  $\min_k \|x^{(i)} - \mu_k\|^2$  [value of k that minimizes this is what to set in  $c^{(i)}$ ]

Move centroid step

for  $k=1$  to  $K$

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

let's say  $x^{(1)}, x^{(5)}, x^{(6)}, x^{(10)}$  assigned to cluster 2,  $c^{(1)}=2, c^{(5)}=2, c^{(6)}=2, c^{(10)}=2$

moving  $\mu_2 \Rightarrow$  (location)  $\mu_2 = \frac{1}{4} [x^{(1)} + x^{(5)} + x^{(6)} + x^{(10)}] \in \mathbb{R}^n$   $\leftarrow$  because all  $x^{(1)}, x^{(5)}, x^{(6)}, x^{(10)}$  are n-dimensional vector (location of cluster centroid k)

Quiz 1

Suppose you run K-means and after the algorithm converges, you have  $c^{(1)}=3, c^{(2)}=3, c^{(3)}=5, \dots$

True statement: - Third example  $x^{(3)}$  has been assigned to cluster 5

- First and second training examples  $x^{(1)}$  and  $x^{(2)}$  have been assigned to the same cluster.

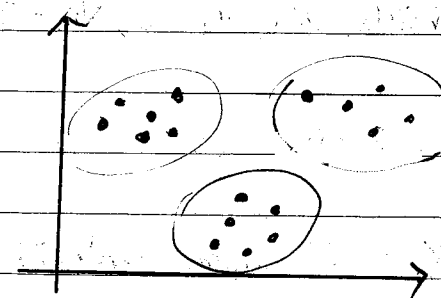
- Out of all the possible values of  $k \in \{1, 2, \dots, K\}$  the value  $k=3$  minimizes  $\|x^{(2)} - \mu_k\|^2$ .

Note:

- Think of  $c^{(i)}$  is picking the cluster centroid with the smallest squared distance to training example  $x^{(i)}$
- Cluster centroid with no points / zero points assigned to it, can just eliminate this cluster centroid.

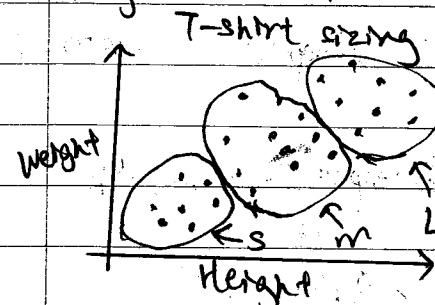
(Sometimes if you really need K cluster, can randomly reinitialize the cluster that don't have points, but more common to just eliminate it)

## K-means for non-separated clusters



$\leftarrow$  3 well separated clusters, we can just let algorithm to find 3 clusters for us.

very often k-mean also apply to data set like this



For example, want to sell T-shirt in 3 sizes  $\Rightarrow$  S, M, L

K-Mean will group these data into 3 clusters

$\leftarrow$  Kind of market segmentation

-Tries to fit the needs of each of the 3 separate sub-populations well.

Clustering - Optimization objective (knowing this help use in debugging learning algorithm, make sure it is running correctly and find better costs and avoid local optima)

K-means optimization objective

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

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

$\mu_{c^{(i)}}$  = cluster centroid of cluster to which example  $x^{(i)}$  has been assigned

let's say  $x^{(1)}$  assigned to 5,  $c^{(1)}=5$

Optimization objective:

$$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$$

$\leftarrow$  location of training example  $x^{(i)}$  distance  $\mu_{c^{(i)}}$  location of cluster centroid to which example  $x^{(i)}$  has been assigned

Try to find c and  $\mu$  to try to minimize cost function J  $\rightarrow \min_{c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K} J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$

$\uparrow$  Distortion cost function / Distortion of the K-means algorithm

## K-means algorithm

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

Repeat {

for  $i=1$  to  $m$

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

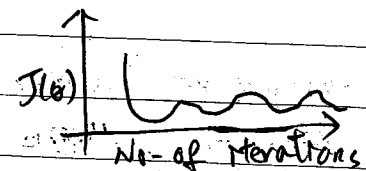
for  $k=1$  to  $K$

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

} move centroid step — minimize  $J(\dots)$  with respect to  $\mu_1, \dots, \mu_K$

Note: First, minimize  $J(\dots)$  with respect to  $c^{(1)}, c^{(2)}, \dots, c^{(m)}$   
Then, minimize  $J(\dots)$  with respect to  $\mu_1, \dots, \mu_K$

Quiz: Suppose you have implemented k-means and to check that it is running correctly, you plot the cost function  $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$  as a function of the number of iterations. Your plot looks like this:



What does this mean?

← It is not possible for the cost function to sometimes increase. There must be a bug in the code.

## Clustering

Random initialization (initialize k-mean)

First step in k-means algorithm

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

Random initialization

Should have  $K < m$

( $K=2$ )

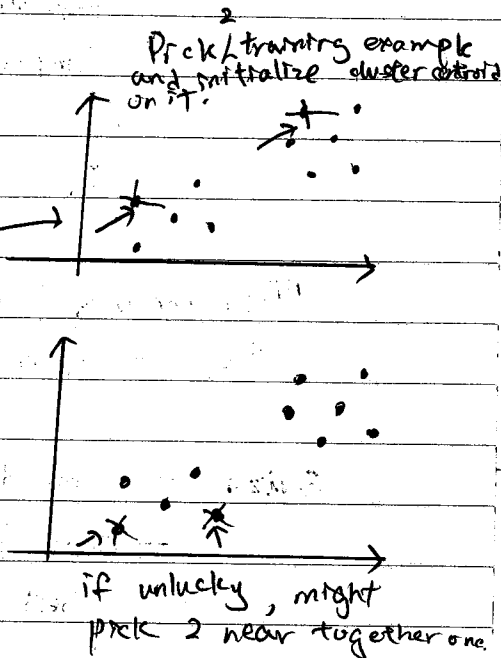
Randomly pick  $K$  training examples.

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

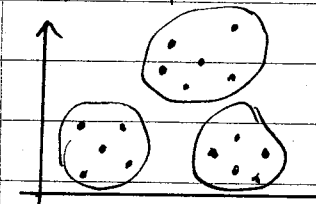
$\mu_1 = x^{(i)}$  randomly choose

$\mu_2 = x^{(j)}$  2 value  $i$  and  $j$

and go on  
if you have more  
clusters centroids.



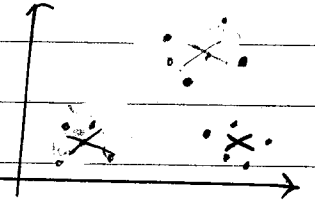
Local optima



Run K-mean

nicely

ends up at a good local optima, might be really global optima, might end up with that clustering.

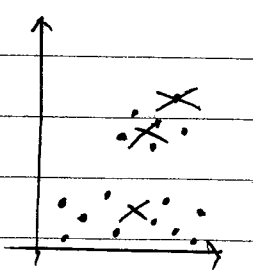
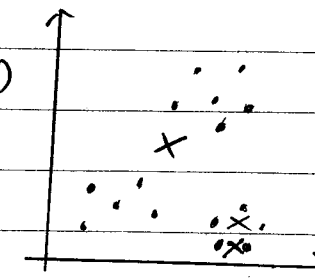


if unlucky random initialization, K-mean also will get stuck at different local optima.

Term local optima refers to local optima of distortion function  $J$

$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$

Solutions where K-means got stuck at local optima and not doing good sub minimizing distortion function.



Note: Want to increase the odd of K-means finding the best possible clustering. One way to do this is to run multiple, random initialization, to keep running K-mean to try, so we can get the good local/global optima as possible.

Random initialization may range from 50-1000

For  $i=1$  to 100 {

Randomly initialize K-means.

Run K-means. Get  $c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K$ .

Compute cost function (distortion)

$$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$$

}

After running 100 times, will have 100 different ways of clustering data

- Pick clustering that gave lowest cost  $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$

Quiz: Recommended way to initialize k-means.

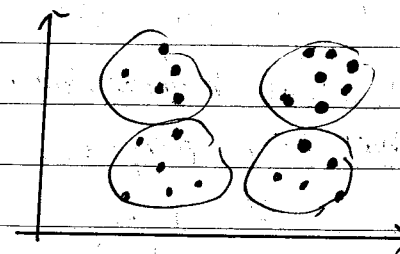
- Pick  $k$  distinct random integers  $i_1, \dots, i_k$  from  $\{1, \dots, m\}$

Set  $\mu_1 = x^{(i_1)}, \mu_2 = x^{(i_2)}, \dots, \mu_k = x^{(i_k)}$

## Clustering

Choosing the number of clusters

What is the right value of  $K$ ?



Some people might say 4, and some might say 2, might say 3 as well. So it is quite ambiguous.

Choosing number of cluster

=> Elbow method:

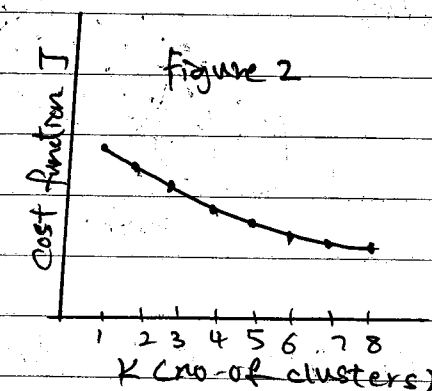
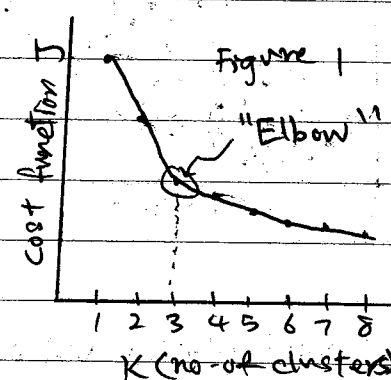


Figure 1 {  
- vary the value of  $K$  and run K-mean to see the value of cost function  $J$   
- go down rapidly from  $1 \rightarrow 2 \rightarrow 3$  and goes very slowly after that  
- So maybe using 3 clusters is the right number of clusters, because it is the elbow of the curve.

Figure 2 {  
- Fairly often will end up with curve that much more ambiguous  
- might no clear elbow, looks like distortion continuously goes down.  
- Harder to choose number of cluster

Note: Not really good method to choose the number of cluster

Quiz: Suppose you run k-means using  $K=3$  and  $K=5$ . You find that cost function  $J$  is much higher for  $K=5$  than  $K=3$ . What can you conclude?

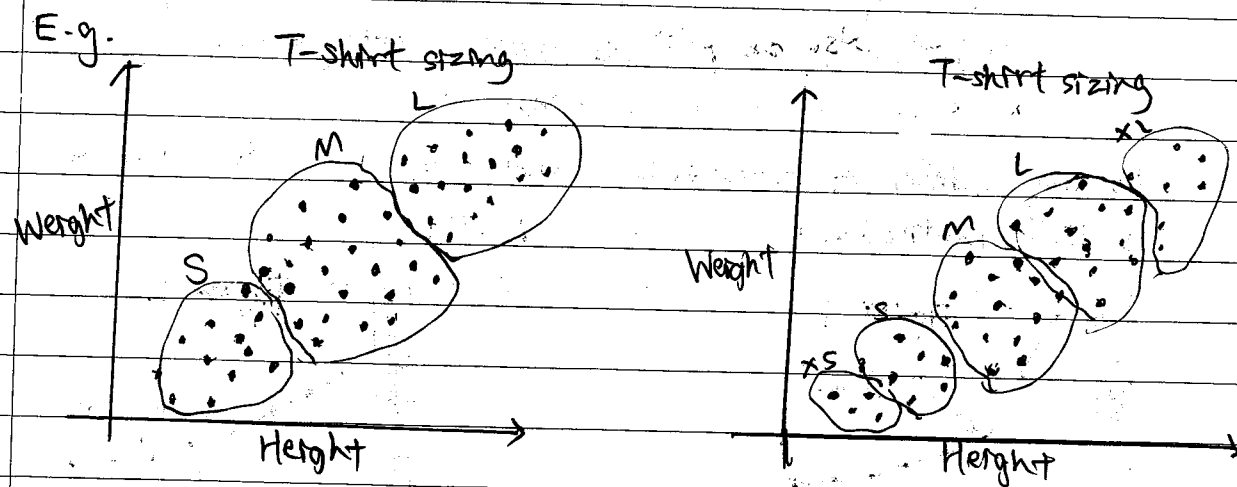
- In the run with  $K=5$ , k-means got stuck in a bad local minimum. You should try re-running k-means with multiple random initializations.

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.

If choose  
 $k=3 \Rightarrow S, M, L$

If choose  
 $k=5 \Rightarrow XS, S, M, L, XL$  Redundant  $\Rightarrow$



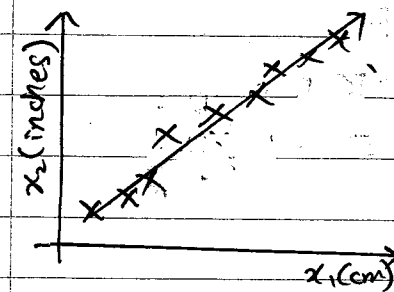
- Can think of this from perspective of T-shirt business and ask
- If has 5 segments, how well will T-shirt fit the customers
  - How many T-shirts can I sell?
  - How happy will my customers be?
  - Want have more/fewer T-shirt size to fit customers better?
  - Or make fewer size of T-shirt to sell customers cheaply.

Better way to choose value of  $K$  is to ask yourself, what purpose are you running K-means? And think, what is the number of cluster  $K$  that serve that what is the later purpose that you actually run the K-means for.

## Dimensionality Reduction

Motivation 1: Data Compression

Data Compression

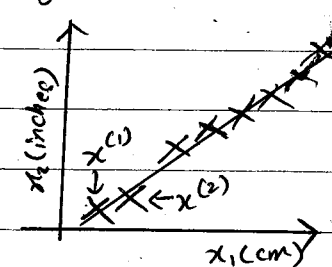


Reduce data from 2D to 1D (Reduce redundancy)

$x_1$  can also be pilot skill  
 $x_2$  can also be pilot employment } 2 features are highly correlated  
 $\Rightarrow$  What you need to care about is that this different feature really measures pilot aptitude

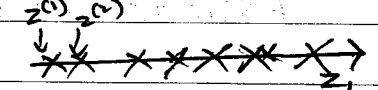
If the features are really highly correlated, need to reduce dimension.

Change to different colour



Draw a straight line  $\Rightarrow$

Measure position of each of the examples on line



Come out with new feature,  $z_1$

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

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

$\vdots$

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

$\Rightarrow$  This halve the memory/space to store data. Data Compression can also allow us to run learning algorithm more quickly. This more interesting in Data compression.

## Data Compression

1000D  $\rightarrow$  100D  
Reduce Data from 3D to 2D

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

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

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

Project the 3D training set onto 2D, and use  $z_1$  and  $z_2$  as the axis.

Quiz Suppose we apply dimensionality reduction to a dataset of  $m$  examples  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ , where  $x^{(i)} \in \mathbb{R}^n$ .

As a result of this, we will get out:

Answer: A lower dimensional dataset  $\{z^{(1)}, z^{(2)}, \dots, z^{(m)}\}$  of  $m$  examples where  $z^{(i)} \in \mathbb{R}^k$  for some value of  $k$  and  $k \leq n$ .

## Dimensionality Reduction

### Motivation II: Data Visualization

Dimensionality Reduction algorithms can reduce 50D (which is 50 features) into

2D ( $\mathbb{R}^2$ ). The output of this algorithm, doesn't ascribe a particular meaning to new features ( $z_1$  and  $z_2$ ).

But it will up to us to figure out what does this features mean

$\Rightarrow$  will find that, in plotting graph, axes  $z_1$  and  $z_2$  can help us to most succinctly capture really what are the 2 main dimensions of the variations amongst different

(example) countries/people/other things that being measured and plot.

Quiz: Suppose you have a dataset  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$  where  $x^{(i)} \in \mathbb{R}^n$ . In order to visualize it, we apply dimensional reduction and get  $\{z^{(1)}, z^{(2)}, \dots, z^{(m)}\}$  where  $z^{(i)} \in \mathbb{R}^k$  is  $k$ -dimensional. In typical setting, which is true?

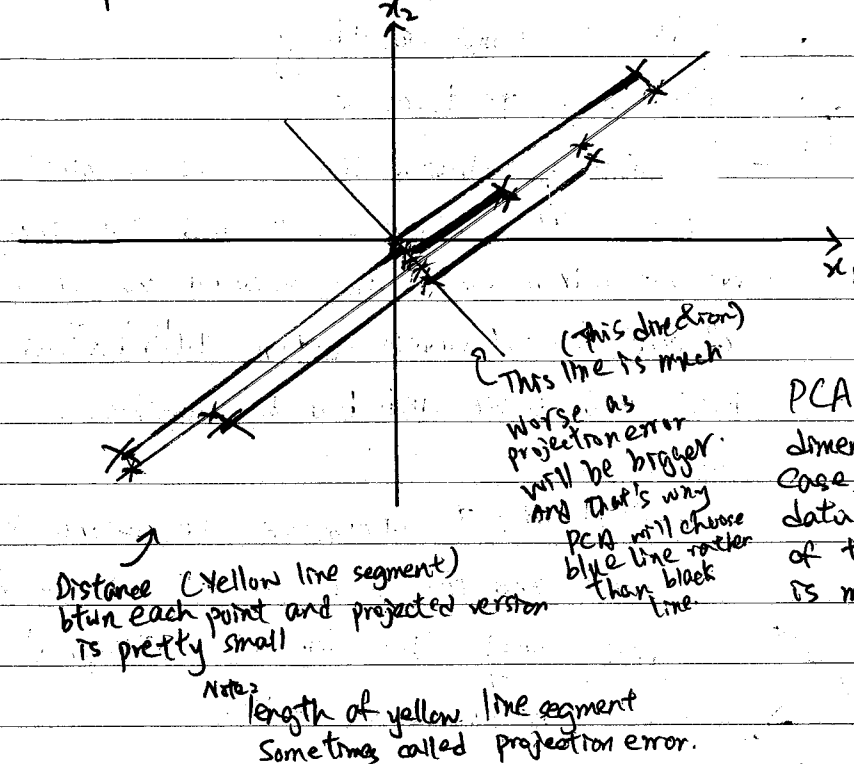
-  $k \leq n$

-  $k=2$  or  $k=3$  (since we can plot 2D or 3D data but don't have ways to visualize higher dimensional data)

## Dimensionality Reduction

### Principal Component Analysis problem formulation

#### PCA problem formulation

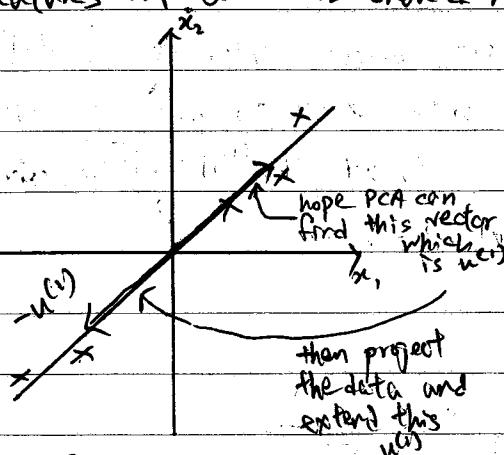


$$x \in \mathbb{R}^2$$

$\leftarrow$  Want reduce from 2D to 1D (to find a line onto which to project the data)

PCA tries to find a lower dimensional surface (a line in this case) onto which to project the data so that the sum of squares of these little yellow line segments is minimized

Before applying PCA, perform mean normalization at feature scaling so that features  $x_1$  and  $x_2$  should have zero mean and comparable ranges of value



Reduce from 2-D to 1-D = Find a direction (a vector  $u^{(1)} \in \mathbb{R}^n$ ) onto which to project the data so as to minimize the projection error

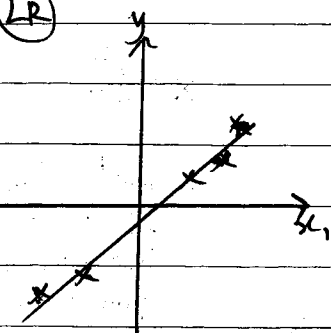
In general: Reduce from  $n$ -dimension to  $k$ -dimension: Find  $k$  vectors  $u^{(1)}, u^{(2)}, \dots, u^{(k)}$  onto which to project the data, so as to minimize the projection error.

Positive or negative doesn't matter as it is only the direction, it can be extended.



PCA is not linear regression (totally different algorithms)

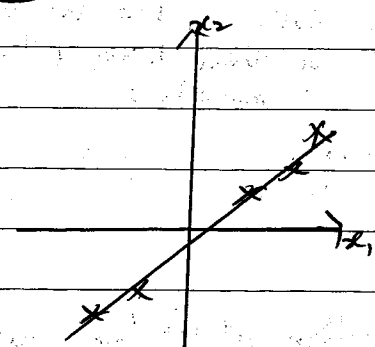
LR



- linear regression is to predict value of some variable  $y$  given info features  $x$
- we fit a straight line, so to minimize square error b/w point and the straight line
- we are minimizing the squared magnitude of purple lines (which are vertical distance between point and value predicted by hypothesis)

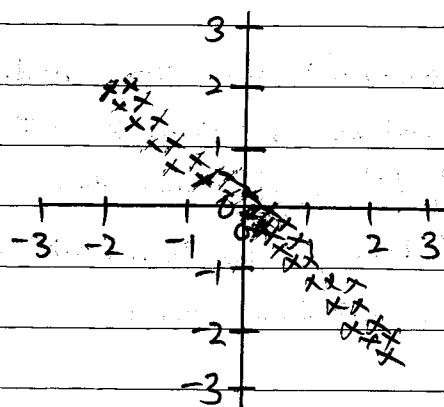
make each feature has exactly zero mean

PCA



- PCA tries to minimize the distance between points and the line, no matter in which angle, just want the shortest orthogonal distances.
- not necessary vertical

Suppose you run PCA on the dataset below. Which of the following would be a reasonable vector  $u^{(1)}$  onto which to project the data? (By convention, we choose  $u^{(1)}$  so that  $\|u^{(1)}\| = \sqrt{(u_1^{(1)})^2 + (u_2^{(1)})^2}$ , the length of the vector  $u^{(1)}$ , equals 1.)



$$u^{(1)} = \begin{bmatrix} -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix}$$

## Dimensionality Reduction

### Principal Component Analysis algorithm

Before applying PCA, have a step to do first.

Data Preprocessing

Training set:  $x^{(1)}, x^{(2)}, \dots, x^{(m)}$  (unlabeled)

Preprocessing (feature scaling/mean normalization):

Compute mean of each feature  $\left\{ \begin{array}{l} M_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)} \\ \text{Replace each } x_j^{(i)} \text{ with } x_j - M_j \end{array} \right.$

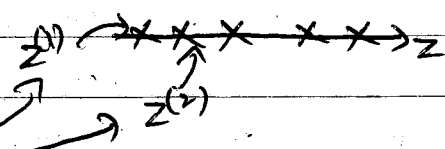
Scale  $\left\{ \begin{array}{l} \text{If different features on different scales (e.g. } x_1 = \text{size of house, } x_2 = \text{number of bedrooms), scale features to have comparable range of values.} \end{array} \right.$

$\frac{x_j^{(i)} - M_j}{S_j}$   
 $S_j \leftarrow$  Measure of range of value of features  $j$  / max-min value /  $\sigma$  of  $j$   
 standard deviation

Reduce data from 2D to 1D  
 $x^{(i)} \in \mathbb{R}^2 \rightarrow z^{(i)} \in \mathbb{R}$

Reduce data from 3D to 2D  
 $x^{(i)} \in \mathbb{R}^3 \rightarrow z^{(i)} \in \mathbb{R}^2$   
 $z = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$

Initially they were  $x^{(1)}, x^{(2)}$  in 2D



Procedure in Principal Component Analysis (PCA) algorithm

Reduce data from  $n$ -dimensions to  $k$ -dimensions

Compute "covariance matrix":  $\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)})(x^{(i)})^T$  (This is summation symbol)  
 Just an alphabet  $\rightarrow$   $n \times n$  matrix

Compute "eigenvectors" of matrix  $\Sigma$ :

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

svd outputs 3 matrices

Thing that SVD really need is  $U$  matrix

If not using octave/MATLAB to compute this, can download a numerical linear algebra library that can compute SVD

$U = \begin{bmatrix} | & | & | & | & | \\ u^{(1)} & u^{(2)} & u^{(3)} & \dots & u^{(m)} \\ | & | & | & | & | \end{bmatrix} \in \mathbb{R}^{n \times n}$   
 Take first  $k$  vectors  $\rightarrow$   $k$  direction onto which we want to project data

SVD  $\rightarrow$  Singular value decomposition

# Principal Component Analysis (PCA) algorithm

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

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

$\underbrace{\hspace{10em}}_k$

Now to make  $x \in \mathbb{R}^n \rightarrow z \in \mathbb{R}^k$

$$\begin{bmatrix} | & | & \dots & | \\ u^{(1)} & u^{(2)} & \dots & u^{(k)} \\ | & | & \dots & | \end{bmatrix}$$

$\underbrace{\hspace{10em}}_{n \times k}$

Matrix  $U_{\text{reduce}}$  (as it is reduced version of  $U$  matrix)

Use this to reduce dimension of data

$$Z^{(i)} = \begin{bmatrix} | & | & \dots & | \\ u^{(1)} & u^{(2)} & \dots & u^{(k)} \\ | & | & \dots & | \end{bmatrix}^T \cdot X^{(i)} = \begin{bmatrix} -(u^{(1)})^T \\ \vdots \\ -(u^{(k)})^T \end{bmatrix} \cdot \underbrace{X^{(i)}}_{n \times 1}$$

$\underbrace{\hspace{10em}}_{k \times n} \quad \underbrace{\hspace{10em}}_{k \times 1}$

The  $x$  can  $\Rightarrow$   
be examples  
in training set / cv set / test set

$z \in \mathbb{R}^k$

## PCA algorithm summary

- After mean normalization (ensure every feature has zero mean) and optionally feature scaling:

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

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

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

$$z = U_{\text{reduce}}^T * x;$$

$$X = \begin{bmatrix} -(x^{(1)})^T \\ \vdots \\ -(x^{(m)})^T \end{bmatrix}$$

$\downarrow$   
Vectorized implementation  
 $\text{Sigma} = (1/m) * X' * X;$

You apply PCA with  $x \in \mathbb{R}^n$

so  ~~$x \in \mathbb{R}^n$~~

Quiz: In PCA, we obtain  $z \in \mathbb{R}^k$  from  $x \in \mathbb{R}^n$  as follows:

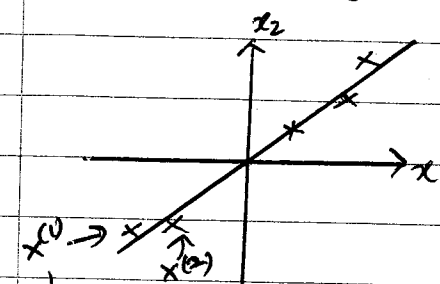
$$z = \begin{bmatrix} | & | & \dots & | \\ u^{(1)} & u^{(2)} & \dots & u^{(k)} \\ | & | & \dots & | \end{bmatrix}^T x = \begin{bmatrix} -(u^{(1)})^T \\ -(u^{(2)})^T \\ \vdots \\ -(u^{(k)})^T \end{bmatrix} x$$

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

Dimensionality Reduction — Reconstruction from compressed representation

PCA is kind of compression algorithm

- So how do we get back to the form that are not compressed?



$$z \in \mathbb{R}^k \rightarrow x \in \mathbb{R}^n$$

$z = U_{\text{reduce}}^T x$

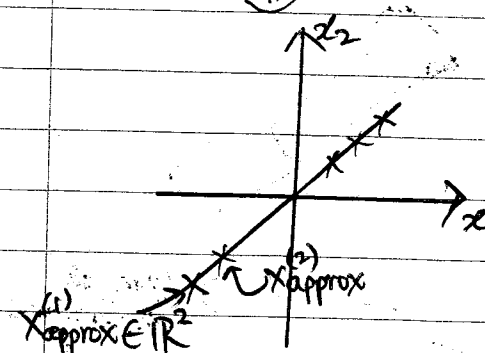
if squared projection error is not too big  $\approx x$  (will close to original value of  $x$  that have used to derive  $z$  in the first place)

Equation:  $X_{\text{approx}} = U_{\text{reduce}} * z$

$\underbrace{\hspace{10em}}_{n \times k} \quad \underbrace{\hspace{10em}}_{k \times 1}$

$\underbrace{\hspace{10em}}_{n \times 1}$

Back to this



Quiz: Suppose we run PCA with  $k=n$ , so that dimension of data is not reduced at all. (This is not useful in practice but is good thought exercise.) Recall that the percent/fraction of variance retained is given by:  $\frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}}$  which ~~are~~ are true?

-  $U_{\text{reduce}}$  will be an  $n \times n$  matrix

-  $X_{\text{approx}} = x$  for every example  $x$ .

- The percentage of variance retained will be 100%



# Dimensionality Reduction

## Choosing the number of principal components

PCA - We take  $n$ -dimensional features and reduce to  $k$ -dimensional feature representation.

number  $k$  is the parameter of PCA algorithm.

↳ also called number of principal components

So, how people choose parameter  $k$  for PCA?

Choosing  $k$  (number of principal components)

PCA tries to minimize  $\Rightarrow$

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

Total variation in the data:  $\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2 \Rightarrow$  on average, how far are the examples from vector or just zero (origin)

Typically, choose  $k$  to be smallest value so that

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

$$(0.05) \quad (5\%)$$

$$(0.10) \quad (10\%)$$

Choose  $k$  so that, "99% of variance is retained"

common range of value is 95-99%

Choosing  $k$  (number of principal components)

Algorithm:

Try PCA with  $k=1$

Compute  $U_{\text{reduce}}, z^{(1)}, z^{(2)}, \dots, z^{(m)}, x_{\text{approx}}^{(1)}, \dots, x_{\text{approx}}^{(m)}$

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 it is good, then use  $k=1$ .  
If not good then next try  $k=2$  and run through procedure and check again, maybe keep trying until a satisfied  $k$  value which 99% of data is retained

For given  $k$ ,

$$1 - \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \leq 0.01$$

↑  
1 minus this ratio can give you quantity of this

$$\text{OR } \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \geq 0.99$$

slowly increase  $k$ , and see which  $k$  ensures 99% of variance is retained.

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

$S = \begin{bmatrix} S_{11} & & \\ & S_{22} & \\ & & S_{33} & \\ & & & \ddots & \\ & & & & S_{nn} \end{bmatrix} \in \mathbb{R}^{n \times n}$

everything except diagonal is zero

Let's say  $k=3$

Sum of  $S_{11}, S_{22}$  and  $S_{33}$

Sum of  $S_{11}$  to  $S_{nn}$

Quiz:

PCA chooses a direction  $u^{(1)}$  (or  $k$  direction  $u^{(1)}, \dots, u^{(k)}$ ) onto which to project the data so as to minimize the (squared) projection error.

Another way is to say the same is the PCA tries to minimize:

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

## Dimensionality Reduction — Advice for applying PCA

Most common use of PCA  $\rightarrow$

Supervised learning speedup

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

Extract inputs: (Put aside the  $y$ 's)

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

↓ PCA

$z^{(1)}, z^{(2)}, \dots, z^{(m)} \in \mathbb{R}^{100}$

New training set:

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

and feed to learning algorithm  $\Rightarrow h_{\theta}(z) = \frac{1}{1 + e^{-\theta^T z}}$

If Get a new example  $x$ , take it and map it to the same mapping that found by PCA to get  $z$ . and feed to hypothesis

Unreduce, after computing all features on training set, just apply same mapping to cv and test set.

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

## Application of PCA

- Compression

- Reduce memory / disk needed to store data
- Speed up learning algorithm
- choose  $k$  by % of variance retained

- Visualization

- plot out the data that is compressed
- $k=2$  or  $k=3$  as we can plot only 2D and 3D data sets.

Bad use of PCA: To prevent overfitting

Use  $z^{(i)}$  instead of  $x^{(i)}$  to reduce the number of features to  $k < n$

Thus, fewer features, less likely to overfit.

And this will work BAD!

Cause PCA only use  $x$  and don't use  $y$ . It will reduce and throw away some information without knowing value of  $y$ .

This might work OK, but isn't a good way to address overfitting. Use regularization instead.

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

PCA is sometimes used where it shouldn't be

Design of ML system:

- 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  $\{(z_{\text{test}}^{(1)}, y_{\text{test}}^{(1)}), \dots, (z_{\text{test}}^{(m)}, y_{\text{test}}^{(m)})\}$

Before run this, ask yourself → How about doing the whole thing without using PCA?

Before implementing PCA, first try running whatever you want to do with the original/raw data  $x^{(i)}$ . Only if that doesn't do what you want, then implement PCA and consider using  $z^{(i)}$ .

only if  $x^{(i)}$  don't work, just consider using compressed representation ( $z^{(i)}$ ).

Note: PCA

- Use to speed up running time of algorithm
- Compress data
- to reduce memory/disk space requirements
- to use it to visualize data

Anomaly Detection

Problem motivation (Density Estimation)

Anomaly detection example

Aircraft engine features:

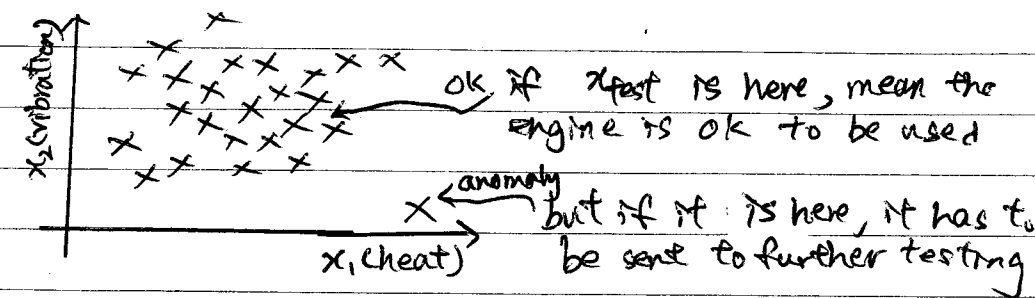
$x_1$  = heat generated

$x_2$  = vibration intensity

...

Dataset:  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

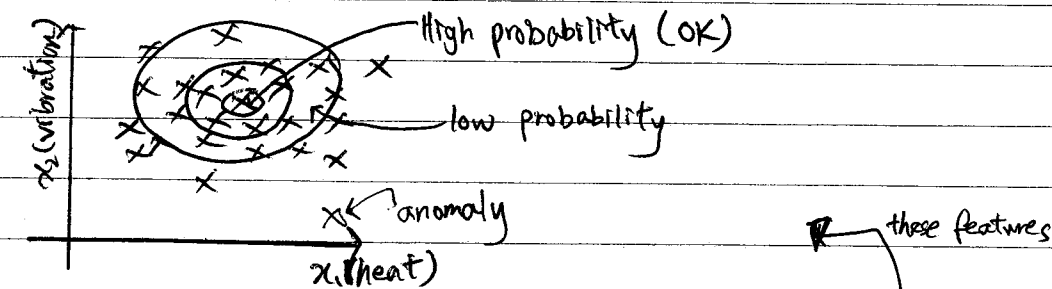
New engine:  $x_{\text{test}}$



Density estimation

Dataset:  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

Is  $x_{\text{test}}$  anomalous?



Given this training set, we build a model for  $p(x)$

$p(x_{\text{test}}) < \epsilon \rightarrow$  flag anomaly

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

Anomaly detection example

① Fraud detection:

$x^{(i)}$  = features of user  $i$ 's activities

Model  $p(x)$  from (data)

maybe  $x_1$  = how often does user log in

$x_2$  = number of web pages visited.

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

If it is unusual of users behavior, send profile of users for further review / demand additional identification from users

② Manufacturing (aircraft engine)

③ Monitoring computers in a data center.

$x^{(i)}$  = features of machine  $i$   
 $x_1$  = memory use,  $x_2$  = number of disk accesses/sec,  $x_3$  = CPU load,  $x_4$  = CPU load / network traffic...

Quiz:

Your anomaly detection system flags  $x$  as anomalous whenever  $p(x) \leq \epsilon$ . Suppose your system is flagging too many things as anomalous that are not actually so (similar to supervised learning, these mistakes are called false positives). What should you do?  
 - Try decreasing  $\epsilon$ .

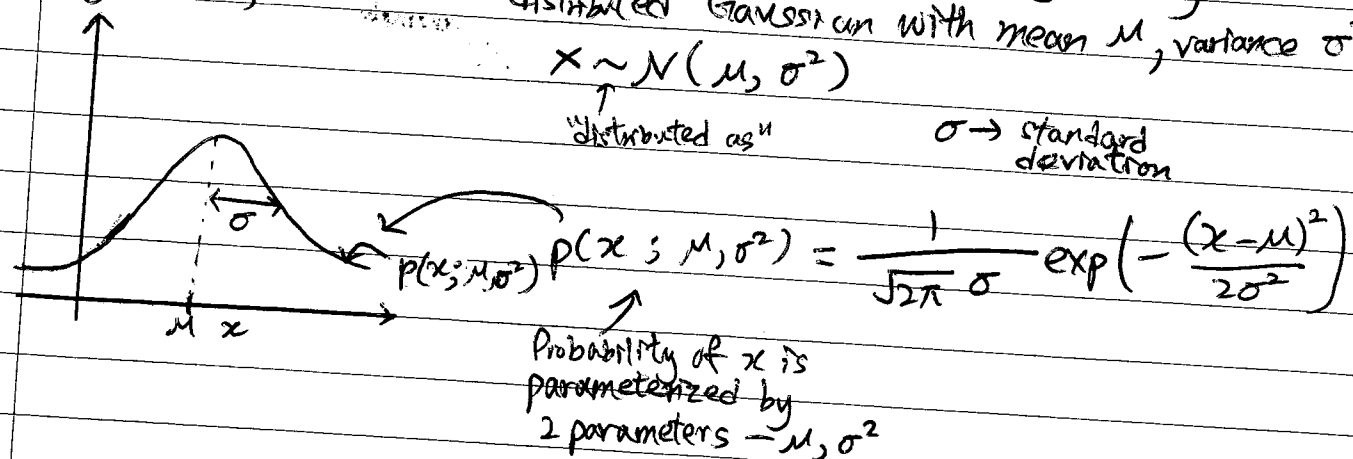
Anomaly detection — Gaussian distribution (normal distribution)  
 Gaussian (Normal) distribution / Gaussian probability density

Say  $x \in \mathbb{R}$ , if  $x$  is distributed Gaussian with mean  $\mu$ , variance  $\sigma^2$ .

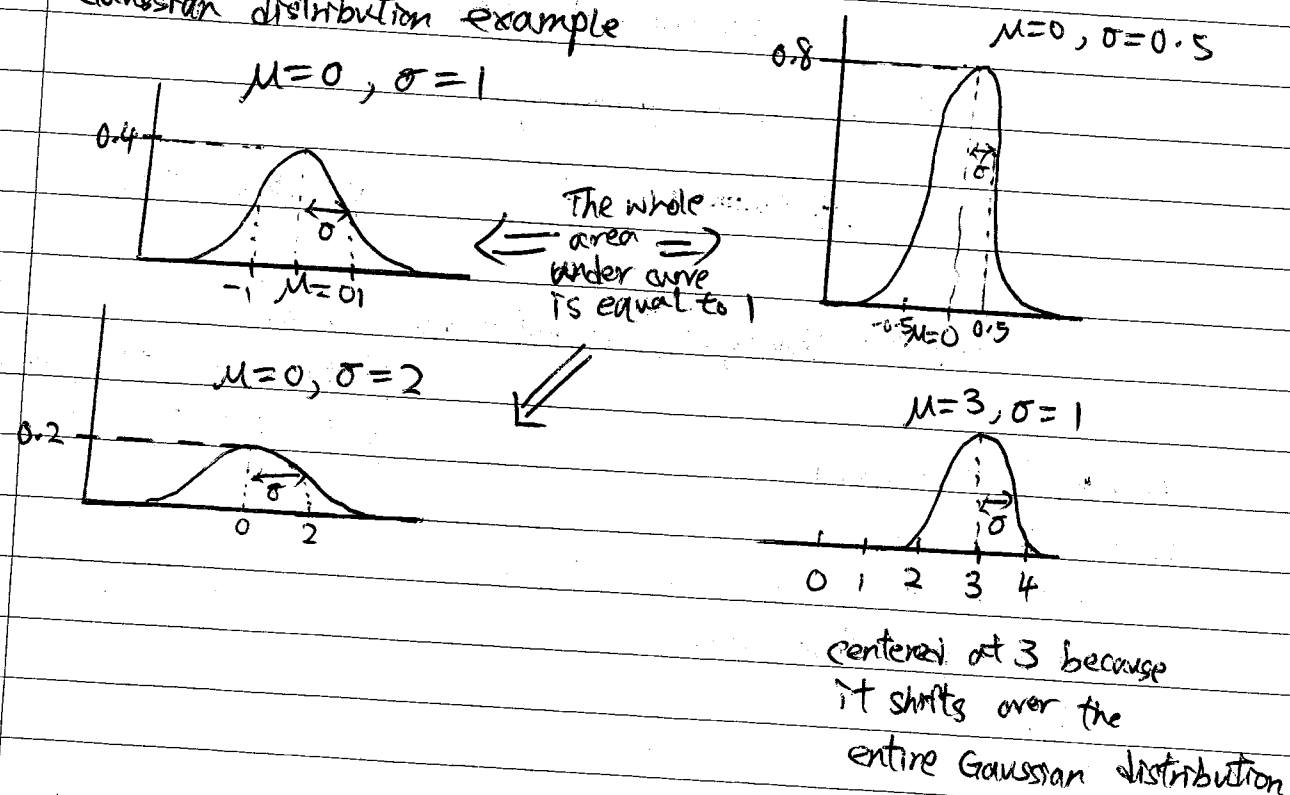
$$x \sim N(\mu, \sigma^2)$$

"distributed as"

$\sigma \rightarrow$  standard deviation

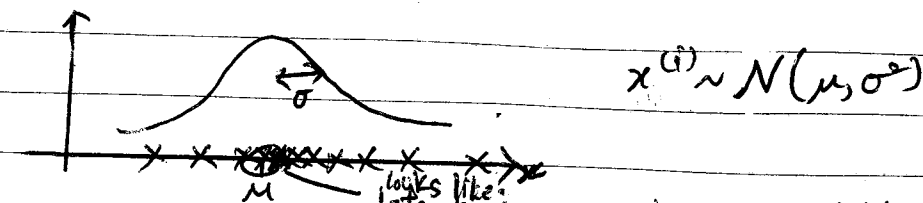


Gaussian distribution example



Parameter estimation

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



Problem of parameter estimation

Is given data set, try to figure out, estimate what are the values of  $\mu$  and  $\sigma^2$

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