

# Machine learning

By - Andrew Ng

Thu 28 Jul

week 1









## Optimization objective

$\leftarrow$  means cost function  $\rightarrow c^{(i)} = \text{index of cluster } (1 \text{ to } k) \text{ to which } x^{(i)}$  is currently assigned.

$\mu_k = \text{cluster centroid } k$

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

$$\underline{x}^{(i)} \in \underline{\mu}_{c^{(i)}}$$

Cost function:

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

Distortion

### Cost function for K-means

$$\underline{\underline{J}}(\underline{\underline{c}}^{(1)}, \dots, \underline{\underline{c}}^{(m)}, \underline{\underline{\mu}}_1, \dots, \underline{\underline{\mu}}_K) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$$

Repeat {

# Assign points to cluster centroids

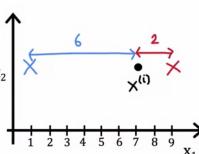
for  $i = 1$  to  $m$

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

# Move cluster centroids

for  $k = 1$  to  $K$

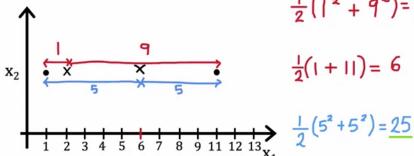
$\mu_k := \text{average of points in cluster } k$



### Moving the centroid

$$\frac{1}{2}(1^2 + 9^2) = \frac{1}{2}(1+81) = 41$$

$$\frac{1}{2}(1+11) = 6$$



$$\frac{1}{2}(5^2 + 5^2) = 25$$

## Initialization K-means ; random location

Step 0: Randomly initialize  $K$  cluster centroids  $\mu_1, \mu_2, \dots, \mu_K$

Repeat {  
Step 1: Assign pt-s to cluster centroids  
Step 2: Move cluster centroids

}

$$K < m$$

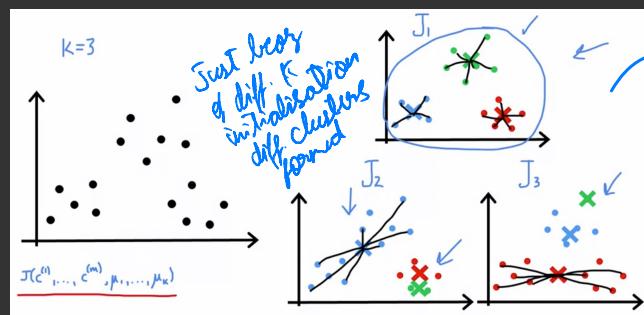
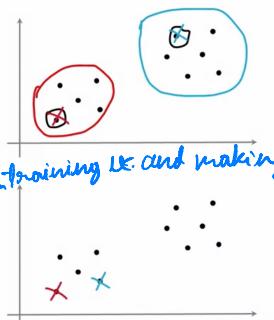
### Random initialization

Choose  $K < m$

Randomly pick  $K$  training

examples. {In this case they  
are using 2 random training LC and making  
the cluster centroid ?}

Set  $\mu_1, \mu_2, \dots, \mu_K$  equal to these  
 $K$  examples.



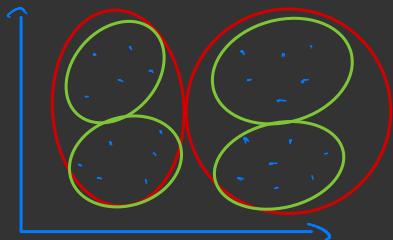
We can choose  
the best  $K$   
initialization  
by comparing  
the cost fun.

Algo :  $\text{for } i = 1 \text{ to } 100 \{$  random samples.  
Run K-means. get  $c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K$   
compute cost fun (distortion)  
 $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$  ←

} Pick set of clusters that gave the lowest cost fun.  $J$

## Choosing the no. of clusters ( $K$ )

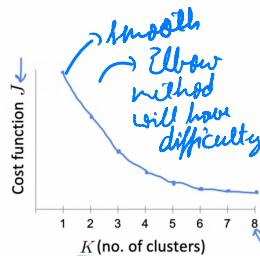
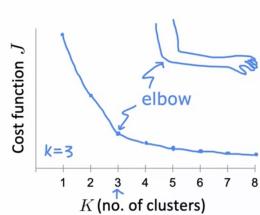
right value of  $K = ?$



Elbow method  $\Rightarrow \{ \text{not preferred} \}$

### Choosing the value of $K$

Elbow method:

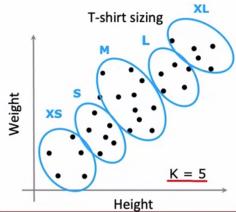
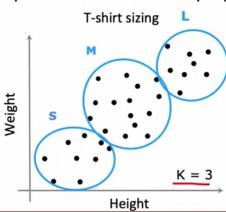


→ Don't choose  $K$  by seeing where the cost fun "is lowest"

### 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.

E.g.

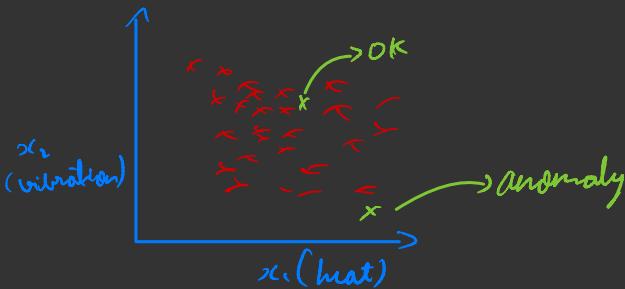


## Anomaly Detection

- Eg Aircraft engine failures:
- $x_1$  = heat generated  
 $x_2$  = vibration intensity

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

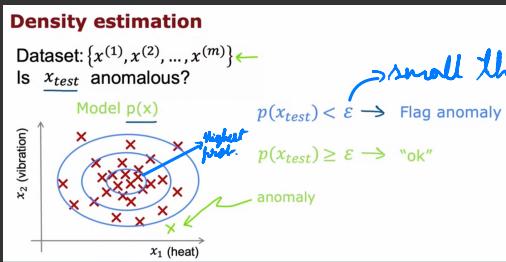
New engine =  $x_{test}$



Is  $x_{test}$  anomalous?

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

First makes a model called  $p(x)$ , which tells which  $x_1$  &  $x_2$  are more probable.



- Eg Fraud detection:

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

Model  $p(x)$  from data

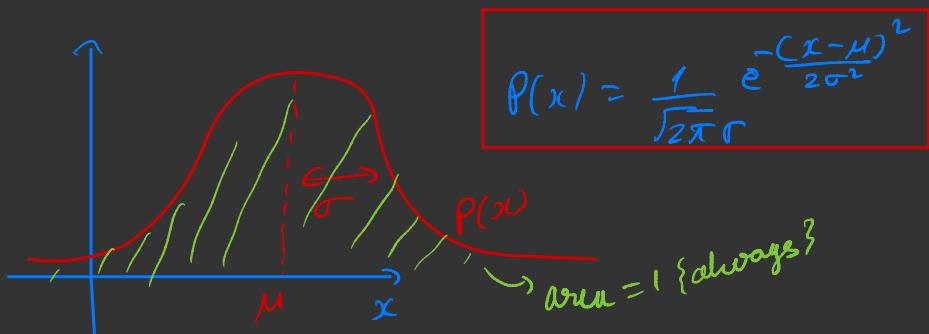
Identify unusual users by checking which have  $p(x) < \varepsilon$

- Eg Monitoring computers in data center.

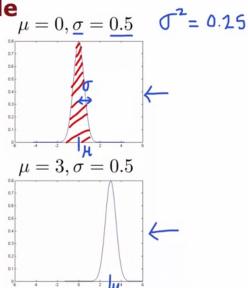
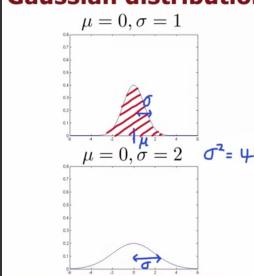
$x^{(i)}$  = features of machine i like CPU load, memory usage.

## Gaussian Distribution (Normal distribution) for P(x)

Say  $x$  is a no., If  $x$  is a Distributed Gaussian with mean  $\mu$ ,  $\sigma^2$  variance



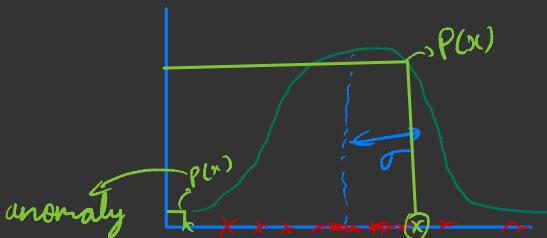
### **Gaussian distribution example**



} Now  $P(x)$  for normal distribution varies with  $\mu$  &  $\sigma$  standard deviation.

## Parameter estimation :

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



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

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

What if  $n$  no. of features? Then how would you implement gaussian distribution?

Training set:  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}_{x_1}$

Each example  $x^{(i)}$  has  $n$  features

$$x^i = \begin{bmatrix} x_1^i \\ x_2^i \\ \vdots \\ x_n^i \end{bmatrix}$$

$$p(x) = p(x_1; \mu_1, \sigma_1^2) * p(x_2; \mu_2, \sigma_2^2) \dots p(x_n; \mu_n, \sigma_n^2)$$

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

### Anomaly detection algorithm

- Choose  $n$  features  $x_i$  that you think might be indicative of anomalous examples.

- Fit parameters  $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$

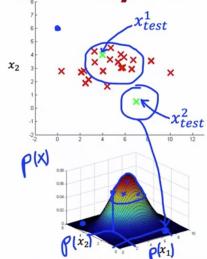
$$\mu_j = \frac{1}{m} \sum_{l=1}^m x_j^{(l)} \quad \sigma_j^2 = \frac{1}{m} \sum_{l=1}^m (x_j^{(l)} - \mu_j)^2$$

- Given new example  $x$ , compute  $p(x)$ :

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

Anomaly if  $p(x) < \epsilon$

### Anomaly detection example



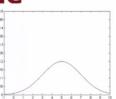
$$\mu_1 = 5, \sigma_1^2 = 2$$

$$\mu_2 = 3, \sigma_2^2 = 1$$

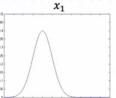
$$\epsilon = 0.02$$

$$p(x_1^{(1)}; \mu_1, \sigma_1^2) = 0.0426 \quad \text{"ok"}$$

$$p(x_2^{(2)}; \mu_2, \sigma_2^2) = 0.0021 \quad \text{"anomaly"}$$



$$p(x_1; \mu_1, \sigma_1^2)$$



$$p(x_2; \mu_2, \sigma_2^2)$$

# How do choose $\varepsilon$ ?

## The importance of real-number evaluation

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

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

Training set:  $x^{(1)}, x^{(2)}, \dots, x^{(m)}$  (assume normal examples/not anomalous)  $\leftarrow y=0$

Cross validation set:  $(x_{cv}^{(1)}, y_{cv}^{(1)}), \dots, (x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$  } include a few anomalous examples

Test set:  $(x_{test}^{(1)}, y_{test}^{(1)}), \dots, (x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$  }  $y=1$

## Aircraft engines monitoring example

10000 good (normal) engines  $\leftarrow 2-50$

$\rightarrow 20$  flawed engines (anomalous)  $\leftarrow y=1$

Training set: 6000 good engines  $\leftarrow y=0$  Train algorithm

CV: 2000 good engines ( $y = 0$ ), 10 anomalous ( $y = 1$ )  $\leftarrow \varepsilon, x_j$

Test: 2000 good engines ( $y = 0$ ), 10 anomalous ( $y = 1$ )  $\leftarrow \varepsilon, x_j$

Alternative: If we dataset like  $y$  anomaly is only 2.

$\rightarrow$  Training set: 6000 good engines  $\leftarrow 2$

$\rightarrow$  CV: 4000 good engines ( $y = 0$ ), 20 anomalous ( $y = 1$ )  $\leftarrow \varepsilon, x_j$

## Algorithm evaluation

Fit model  $p(x)$  on training set  $x^{(1)}, x^{(2)}, \dots, x^{(m)}$

On a cross validation/test example  $x$ , predict

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

## Algorithm evaluation

$\rightarrow$  Fit model  $p(x)$  on training set  $x^{(1)}, x^{(2)}, \dots, x^{(m)}$

On a cross validation/test example  $x$ , predict

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

10  
2000

Possible evaluation metrics:

- True positive, false positive, false negative, true negative
- Precision/Recall
- $F_1$ -score

Can also use cross validation set to choose parameter  $\varepsilon$ .

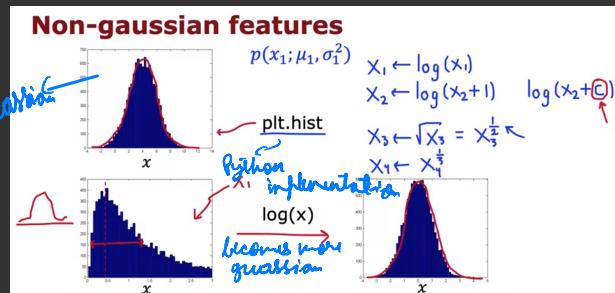
When to use anomaly detection vs when to use supervised learning?

Anomaly detection	vs.	Supervised learning
<p>Very small number of positive examples (<math>y = 1</math>). (0-20 is common). Large number of negative (<math>y = 0</math>) examples. <math>p(x)</math> <math>y=1</math></p> <p>Many different "types" of anomalies. Hard for any algorithm to learn from positive examples what the anomalies look like; <u>future anomalies may look nothing like any of the anomalous examples we've seen so far.</u></p> <p><u>Fraud</u> <math>\rightarrow</math> always new type of frauds</p>		<p>Large number of positive and negative examples.</p> <p><b>20 positive examples</b></p> <p>Enough positive examples for algorithm to get a sense of what positive examples are like, <u>future positive examples likely to be similar to ones in training set.</u></p> <p><u>Spam</u> <math>\rightarrow</math> very similar</p>

Anomaly detection	vs.	Supervised learning
<ul style="list-style-type: none"><li>→ Fraud detection</li><li>→ Manufacturing - Finding <u>new previously unseen defects</u> in manufacturing. (e.g. aircraft engines)</li><li>→ Monitoring machines in a data center</li><li>⋮</li></ul>		<ul style="list-style-type: none"><li>→ Email spam classification</li><li>→ Manufacturing - Finding known, previously <u>seen</u> defects <math>y=1</math> <u>scratches</u></li><li>→ Weather prediction (sunny/rainy/etc.)</li><li>→ Diseases classification</li><li>⋮</li></ul>

## Choosing what features to use :

Non-gaussian features  $\rightarrow$



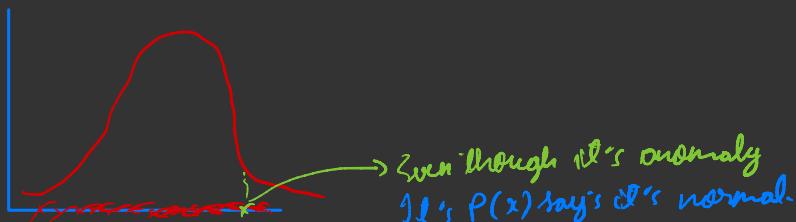
## Error analysis for anomaly detection :

want  $p(x) \geq \varepsilon$  large for normal samples  $x$ .

$p(x) < \varepsilon$  small for anomalous samples  $x$ .

most common problem:

$p(x)$  is comparable (say, both large) for normal & anomalous samples.



Eg: Monitoring computers in a data center :

$x_1$  = memory use of computer

$x_2$  = no. of disk accesses/sec.

$x_3$  = CPU load

$x_4$  = network traffic

$$x_5 = \frac{\text{CPU Load}}{\text{Network Traffic}}$$

$$x_6 = \frac{(\text{CPU Load})^2}{\text{Network Traffic}}$$

5. You are monitoring the temperature and vibration intensity on newly manufactured aircraft engines. You have measured 100 engines and fit the Gaussian model described in the video lectures to the data. The 100 examples and the resulting distributions are shown in the figure below.

The measurements on the latest engine you are testing have a temperature of 17.5 and a vibration intensity of 48. These are shown in magenta on the figure below. What is the probability of an engine having these two measurements?

