

## Anomaly Detection

e.g. aircraft engine features

$x_1 = \text{heat generated}$

$x_2 = \text{vibration intensity}$

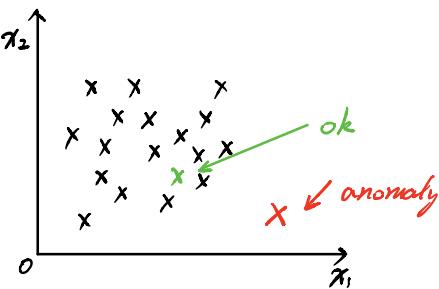
.....

$p(x) \geq q$  ok

$p(x) < q$  anomaly

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

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



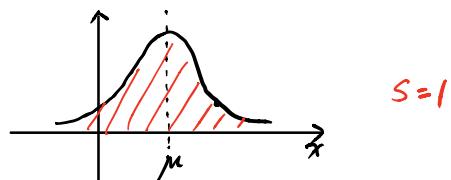
## Anomaly detection example

- fraud detection
- manufacturing (example above)
- monitor computers in data center

## Gaussian distribution

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

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



$$P(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

parameter estimation:

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

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

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

square error

anomaly detection algorithm:

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

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

$$x_j \in N(\mu_j, \sigma_j^2) \quad j=1, 2, \dots, n$$

2. fit parameters  $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)} \quad j=1, 2, \dots, n$$

$$\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2 \quad j=1, 2, \dots, n$$

3. give new example  $x$ , compute  $p(x)$

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

if  $p(x) < \epsilon \Leftrightarrow \text{anomaly}$

algorithm evaluation (how to choose  $\epsilon$ ?)

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

on cross validation/test example  $x$ , predict

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

evaluation metrics:

- true positive, false positive, true negative, false negative
- precision / recall
- $F_1$  - score

anomaly detection  
 mainly normal cases, few anomaly cases.  
 many different type of cases hard for algorithm to learn.

vs  
 supervised learning  
 large number of positive and negative cases  
 enough positive cases for algorithms to learn.

- fraud detection
- manufacturing
- monitoring machines

- spam email classification
- weather prediction
- cancer classification

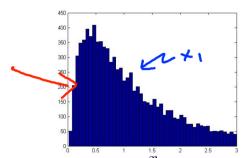
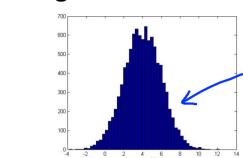
:

:

choosing what features to use.

check each dimension of data manually to see if the feature is gaussian or not. If they are non-gaussian, we need to make them gaussian by using:  $x \leftarrow \log(x+c)$  or  $x \leftarrow x^a$

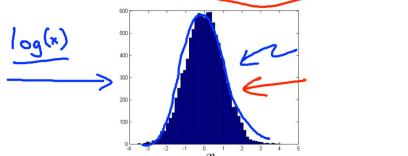
Non-gaussian features



$$p(x_i | \mu_i, \sigma_i^2)$$

$$x_1 \leftarrow \log(x_i)$$

$$\log(x)$$



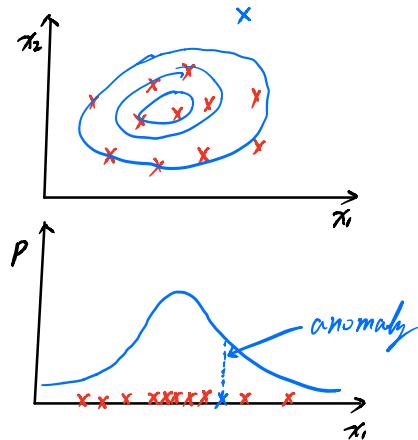
error analysis for anomaly detection

want  $p(x)$  large for normal examples  $x$

want  $p(x)$  small for anomalous examples  $x$

most common problem:

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

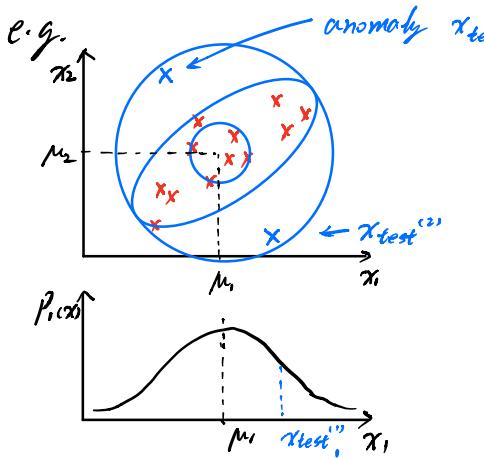


$x_1$  = memory use of computer

$x_2$  = CPU load

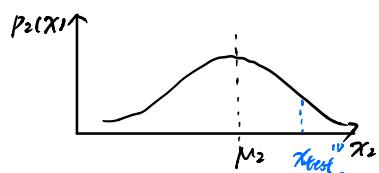
we now need to create a new feature

$$x_3 = \frac{x_1}{x_2} \text{ or } \frac{x_1^2}{x_2} \dots$$



$$P(x \text{ test''}) = P_1(x_{\text{test}}^{(1)}) \times P_2(x_{\text{test}}^{(2)}) > \epsilon$$

$\times$  predict normal but actual anomaly  
for the boundary is  $\circ$  rather than  $\square$



can add new feature  $x_3 = \frac{x_1}{x_2}$   
or use "multivariate gaussian distribution"

## Multivariate Gaussian Distribution

$x \in \mathbb{R}^n$ , model  $p(x)$  all in one instead of modelling  $p(x_1), p(x_2) \dots$  separately

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

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

parameter fitting:

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

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)} \in \mathbb{R}^n$$

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)(x^{(i)} - \mu)^T \in \mathbb{R}^{n \times n}$$

relation between the origin model and new model:

$$\text{origin } p(x) = p(x_1; \mu_1, \sigma_1^2) \times p(x_2; \mu_2, \sigma_2^2) \times \dots \times p(x_n; \mu_n, \sigma_n^2)$$

in new model:

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)(x^{(i)} - \mu)^T = \begin{bmatrix} \sigma_1^2 & & \\ & \sigma_2^2 & \\ & & \ddots \\ & & \sigma_n^2 \end{bmatrix}_{n \times n}$$

origin model

manually create features to capture anomalies  
where  $x_1, x_2$  take unusual combination of  
values.

$$x_3 \leftarrow \frac{x_1}{x_2}$$

multivariate Gaussian distribution

automatically captures correlations  
between features.

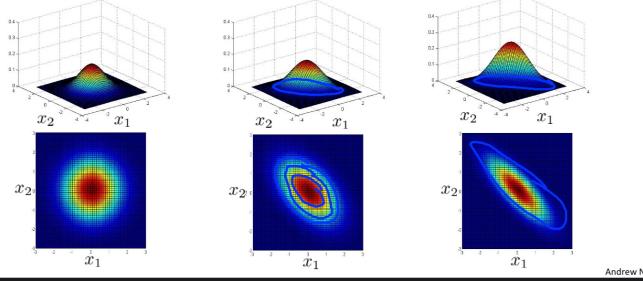
computationally cheaper, scale better to large  
ok if  $m$  is small

computationally more expensive  
must have  $m \gg n$  ( $m > n$ )  
 $\Sigma$  be invertible

when the graph like  $O, \circ, \emptyset$ , use origin model directly  
 when the graph like  $\partial, \odot$ , use multivariate gauss distribution

### Multivariate Gaussian (Normal) examples

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



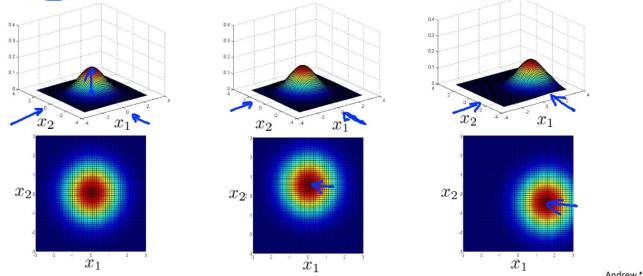
Andrew N

### Multivariate Gaussian (Normal) examples

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

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

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



Andrew N