

# **SEMANA 9 - DETECCIÓN DE ANOMALÍAS**

# 1. ESTIMACIÓN DE DENSIDAD

- Motivación del problema

Principalmente en problemas de aprendizaje no supervisado.

Supuesto de ejemplos normales.

Con los datos de entrenamiento sin etiqueta generamos un modelo de probabilidad  $P(x)$ .

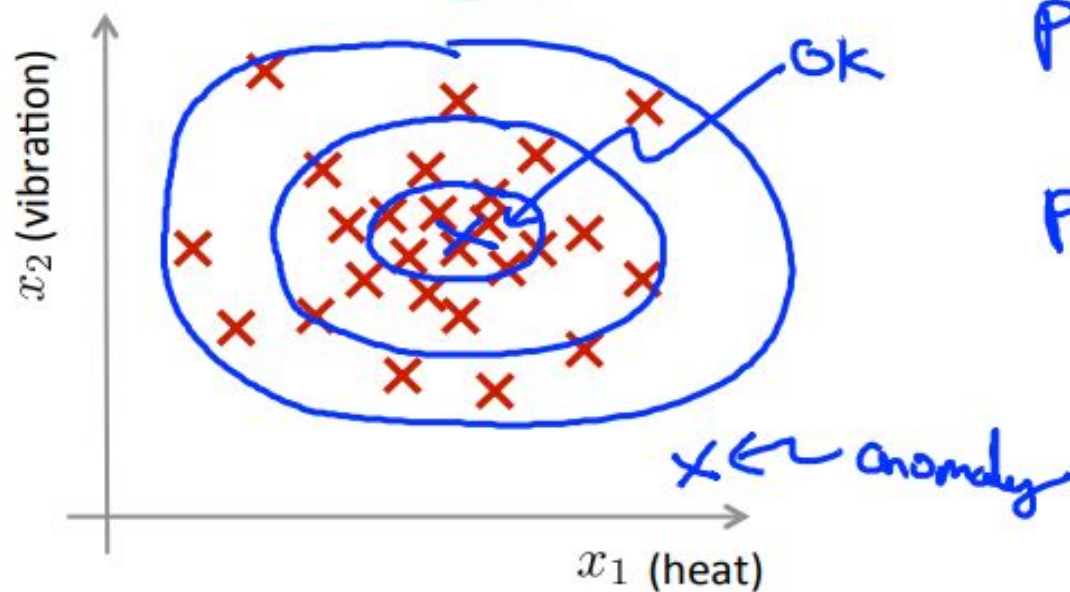
Introducimos nuevos datos (test) en el modelo de probabilidad y lo comparamos con un valor 'exilon' para er si es anormal o no.

# Density estimation

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

→ Is  $x_{test}$  anomalous?

Model  $\underline{p(x)}$ .



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

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

- Distribución Gaussiana

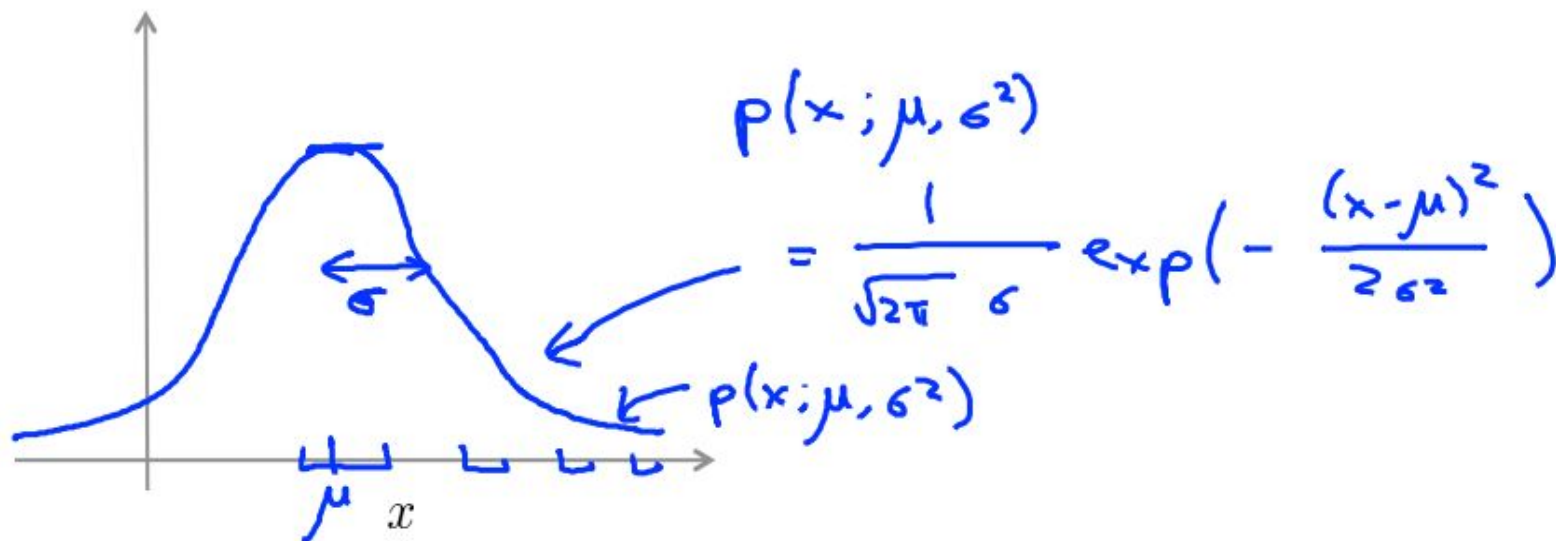
**Gaussian (Normal) distribution**

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

$$x \sim \mathcal{N}(\overset{\downarrow}{\mu}, \overset{\downarrow}{\sigma^2})$$

$\nwarrow$  "distributed as"

$\sigma$  standard deviation



- Algoritmo

## Anomaly detection algorithm

→ 1. Choose features  $x_i$  that you think might be indicative of anomalous examples.  $\{x^{(1)}, \dots, x^{(m)}\}$

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

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

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

$$p(x_j; \mu_j, \sigma_j^2)$$

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix} = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

→ 3. 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  $p(x) < \varepsilon$

# CREACIÓN DE UN SISTEMA DE DETECCIÓN DE ANOMALÍAS

- Desarrollo y Evaluación de un sistema de detección de anomalías

Tomar decisiones es mucho más fácil si tenemos una forma de evaluar nuestro algoritmo de aprendizaje.

Training Set. Asumimos datos normales ( $y=0$ )

## Aircraft engines motivating example

- 10000 good (normal) engines
- 20 flawed engines (anomalous)  $\frac{2-50}{y=1}$
- Training set: 6000 good engines ( $y=0$ )  $\mu_1, \sigma_1^2, \dots, \mu_n, \sigma_n^2$   $p(x) = p(x_1; \mu_1, \sigma_1^2) \dots p(x_n; \mu_n, \sigma_n^2)$
- CV: 2000 good engines ( $y=0$ ), 10 anomalous ( $y=1$ )
- Test: 2000 good engines ( $y=0$ ), 10 anomalous ( $y=1$ )

Alternative:

Training set: 6000 good engines

- CV: 4000 good engines ( $y=0$ ), 10 anomalous ( $y=1$ )
- Test: 4000 good engines ( $y=0$ ), 10 anomalous ( $y=1$ )



## Algorithm evaluation

- Fit model  $\underline{p(x)}$  on training set  $\{x^{(1)}, \dots, x^{(m)}\}$
- On a cross validation/test example  $\underline{x}$ , predict

$$(x_{\text{test}}^{(i)}, y_{\text{test}}^{(i)})$$

↑

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

$$\underline{y=0}$$

Possible evaluation metrics:

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

CV

Test set

Can also use cross validation set to choose parameter  $\underline{\epsilon}$  ←



## - Detección de anomalías VS Aprendizaje Supervisado

### Anomaly detection

- Very small number of positive examples ( $y = 1$ ). (0-20 is common).
- Large number of negative ( $y = 0$ ) examples.  $p(x)$  ←
- Many different “types” of anomalies. Hard for any algorithm to learn from positive examples what the anomalies look like;
- future anomalies may look nothing like any of the anomalous examples we've seen so far.

vs.

### Supervised learning

Large number of positive and negative examples. ←

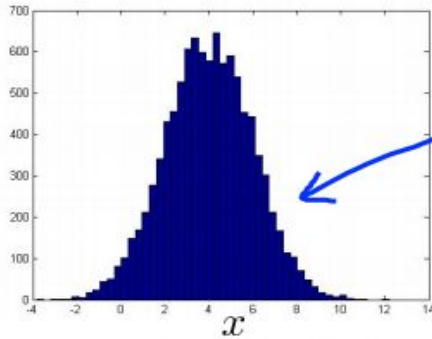
Enough positive examples for algorithm to get a sense of what positive examples are like, future positive examples likely to be similar to ones in training set. ←

Spam ←

- Elegir qué variables usar

Realizar un histograma de las variables, si presenta asimetría transformar hasta obtener la forma de Campana de Gauss.

## Non-gaussian features



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

hist

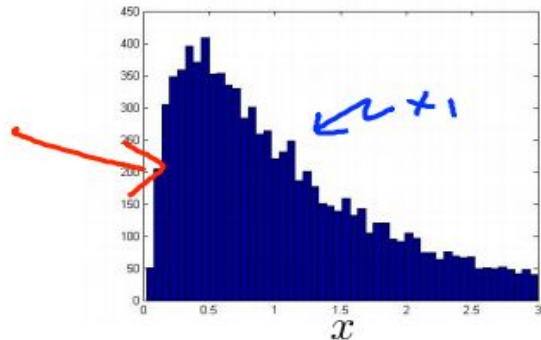
$$x_1 \leftarrow \frac{\log(x_i)}{x_i}$$

$$x_2 \leftarrow \log(x_2 + 1)$$

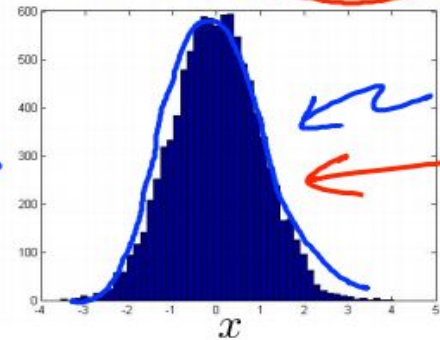
$$x_3 \leftarrow \sqrt{x_3} = x_3^{\frac{1}{2}}$$

$$x_4 \leftarrow x_4^{\frac{1}{3}}$$

$$\log(x_2 + 1)$$



$$\frac{\log(x)}{x}$$

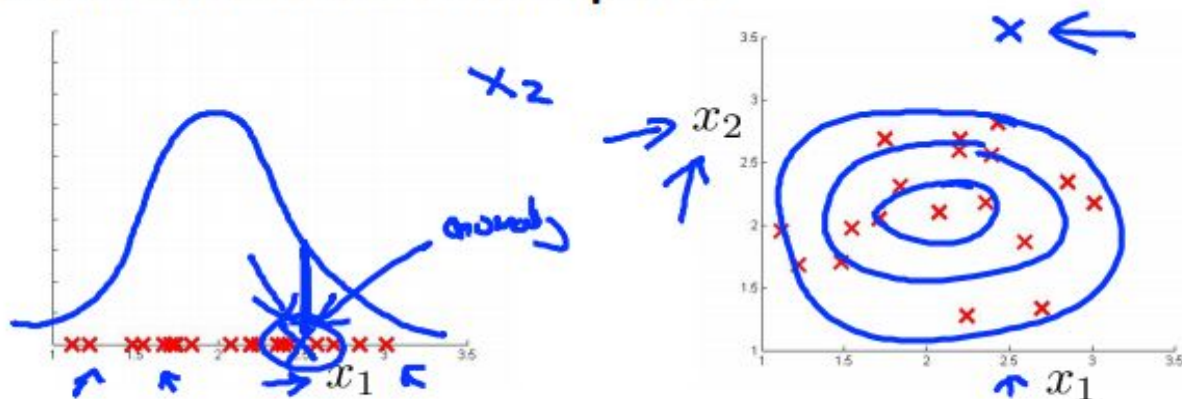


## → Error analysis for anomaly detection

Want  $p(x)$  large for normal examples  $x$ .  
 $p(x)$  small for anomalous examples  $x$ .

Most common problem:

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

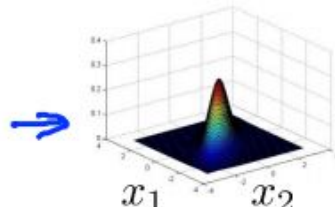
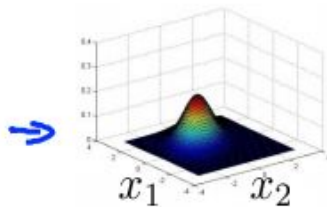
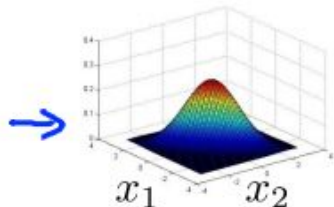


# Distribución Gaussiana Multivariante

## Multivariate Gaussian (Normal) distribution

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

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



Parameter fitting:

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

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



# Anomaly detection with the multivariate Gaussian

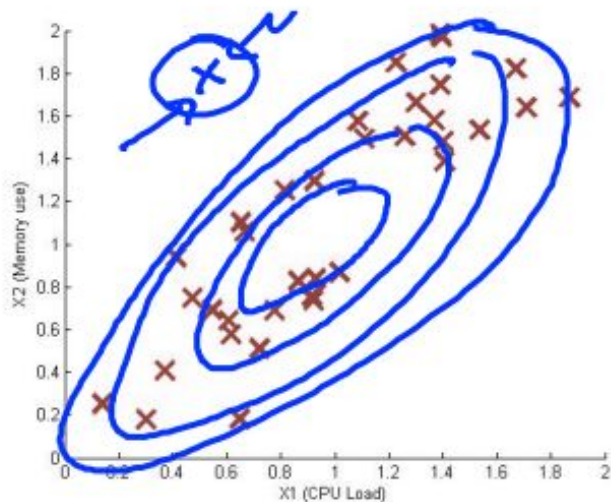
1. Fit model  $p(x)$  by setting

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

2. Given a new example  $x$ , compute

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

Flag an anomaly if  $\underline{p(x) < \varepsilon}$



## → Original model

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

Manually create features to capture anomalies where  $x_1, x_2$  take unusual combinations of values.

$$\rightarrow X_3 = \frac{x_1}{x_2} = \frac{\text{CPU load}}{\text{memory}}$$

→ Computationally cheaper (alternatively, scales better to large  $n=10,000, n=100,000$ )

OK even if  $m$  (training set size) is small

## vs. → Multivariate Gaussian

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

→ Automatically captures correlations between features

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

$$\Sigma^{-1}$$

Computationally more expensive

$$\rightarrow \Sigma \sim \frac{n^2}{2}$$

Must have  $m > n$  or else  $\Sigma$  is non-invertible. →  $m \geq 10n$

$$\left[ \begin{array}{l} \rightarrow X_1 = \cancel{X_2} \\ \cancel{X_3} = X_4 + X_5 \end{array} \right]$$



# SISTEMAS DE RECOMENDACIÓN

## Content-based recommender systems

$n_u = 4, n_m = 5$

$x_0 = 1$

$x^{(i)} = \begin{bmatrix} 1 \\ 0.9 \\ 0 \end{bmatrix}$

$n = 2$

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

- > For each user  $j$ , learn a parameter  $\theta^{(j)} \in \mathbb{R}^3$ . Predict user  $j$  as rating movie  $i$  with stars.  $\theta^{(j)} \in \mathbb{R}^{n+1}$

$$x^{(3)} = \begin{bmatrix} 1 \\ 0.99 \\ 0 \end{bmatrix} \leftrightarrow \theta^{(1)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix} \quad (\theta^{(1)})^T x^{(3)} = 5 \times 0.99 = 4.95$$

## Optimization algorithm:

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

Gradient descent update:

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

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

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

## Collaborative filtering algorithm

~~$x_0 = 1$~~   $x \in \mathbb{R}^n$ ,  $\theta \in \mathbb{R}^n$   
 ~~$\theta_0$~~   
 $\theta_1$   
 $\vdots$   
 $\theta_n$

- 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  $j = 1, \dots, n_u, i = 1, \dots, n_m$  :

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

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

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

# Collaborative filtering

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

Predicted ratings:

$$\begin{bmatrix} (\theta^{(1)})^T(x^{(1)}) & (\theta^{(2)})^T(x^{(1)}) & \dots & (\theta^{(n_u)})^T(x^{(1)}) \\ (\theta^{(1)})^T(x^{(2)}) & (\theta^{(2)})^T(x^{(2)}) & \dots & (\theta^{(n_u)})^T(x^{(2)}) \\ \vdots & \vdots & \vdots & \vdots \\ (\theta^{(1)})^T(x^{(n_m)}) & (\theta^{(2)})^T(x^{(n_m)}) & \dots & (\theta^{(n_u)})^T(x^{(n_m)}) \end{bmatrix}$$

$$X = \begin{bmatrix} -(x^{(1)})^T \\ -(x^{(2)})^T \\ \vdots \\ -(x^{(n_m)})^T \end{bmatrix}$$

$$L = \begin{bmatrix} -(\theta^{(1)})^T \\ -(\theta^{(2)})^T \\ \vdots \\ -(\theta^{(n_u)})^T \end{bmatrix}$$

→ Low rank matrix factorization