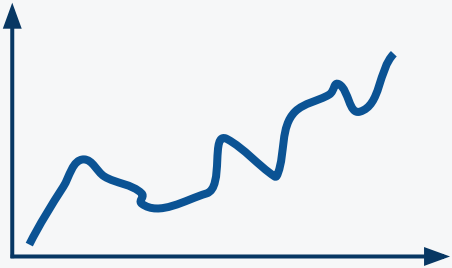


---

# Detección de anomalías

PORTAFOLIOS | EQUIPO 1  
Caso de estudio 2

Forma supervisada



VS

Forma no supervisada



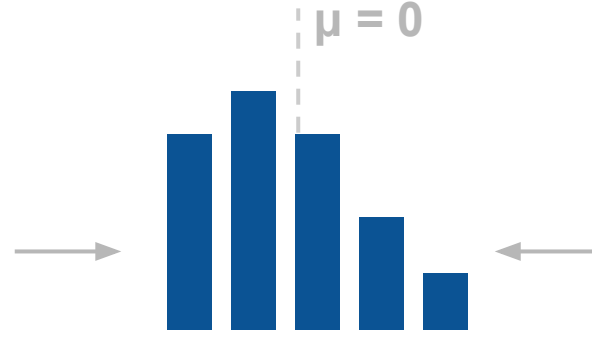


## 2

## Problemas



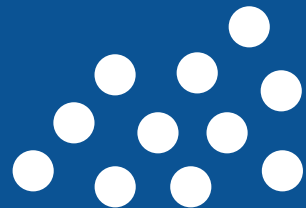
Distribución **sesgada**



Datos **normalizados**

---

# Distancias





## ¿Cómo detectar un outlier?

La detección se basa en la suposición de que los datos atípicos se encuentran separados de la mayoría de los datos del conjunto de datos.

1

Se le asigna un puntaje a cada instancia, cuyo valor es la distancia hasta su k-vecino más cercano.

2

Se seleccionan las  $n$  instancias con mayor puntaje y se dice que son probables outliers.



1

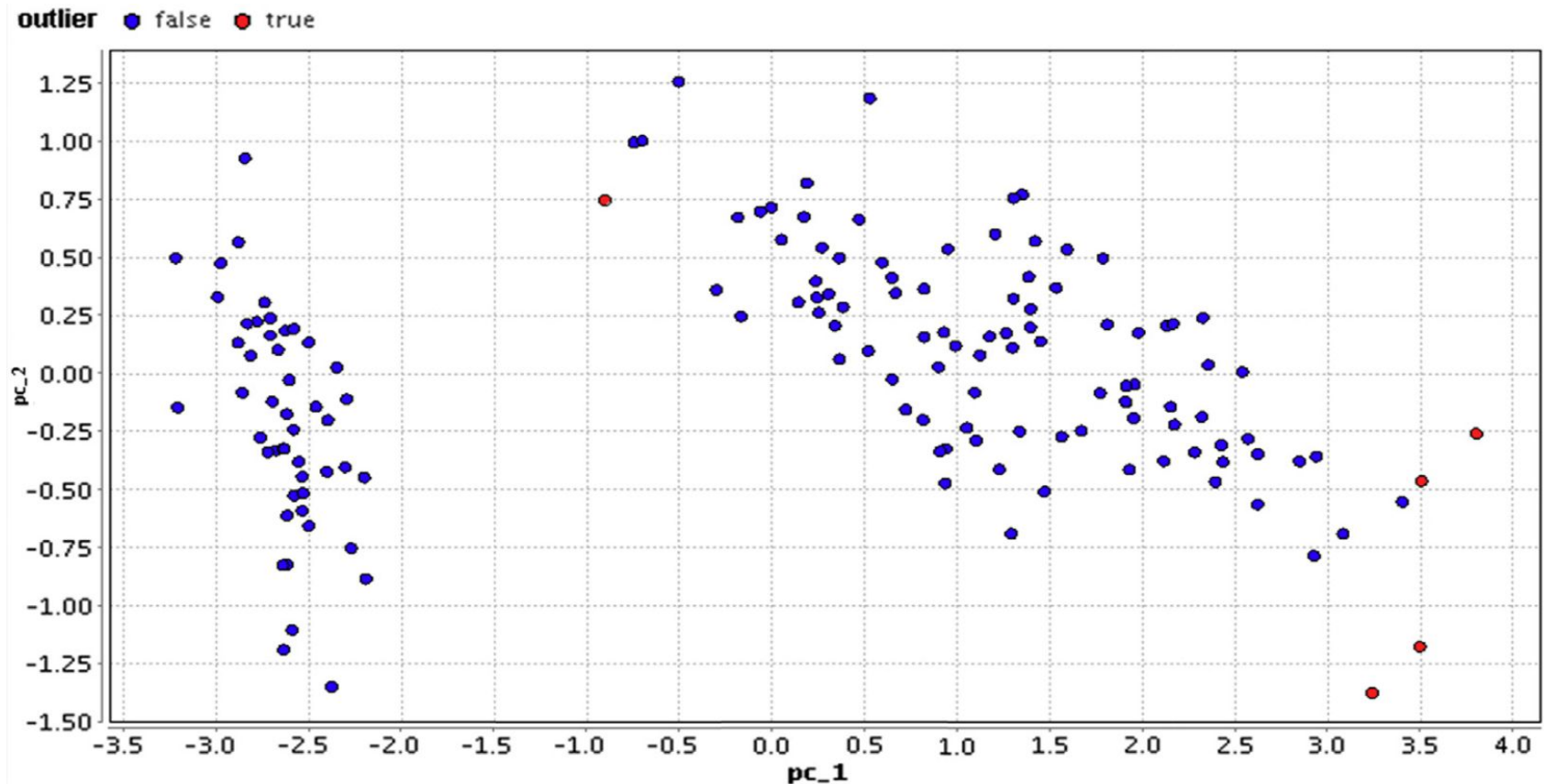
Número de **vecinos**

2

Número de **outliers**

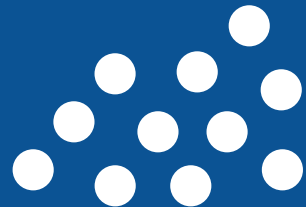
3

Función de **distancia**



---

# Densidad







## ¿Cómo detectar un outlier?

- ★ Los valores atípicos ocurren con menos frecuencia que datos normales → En el espacio de los datos atípicos ocupan áreas de baja densidad y puntos de datos normal ocupan zonas de alta densidad.
- ★ La densidad es un recuento de los puntos de datos en una unidad normalizada espacio y es inversamente proporcional a la distancia entre puntos de datos



1

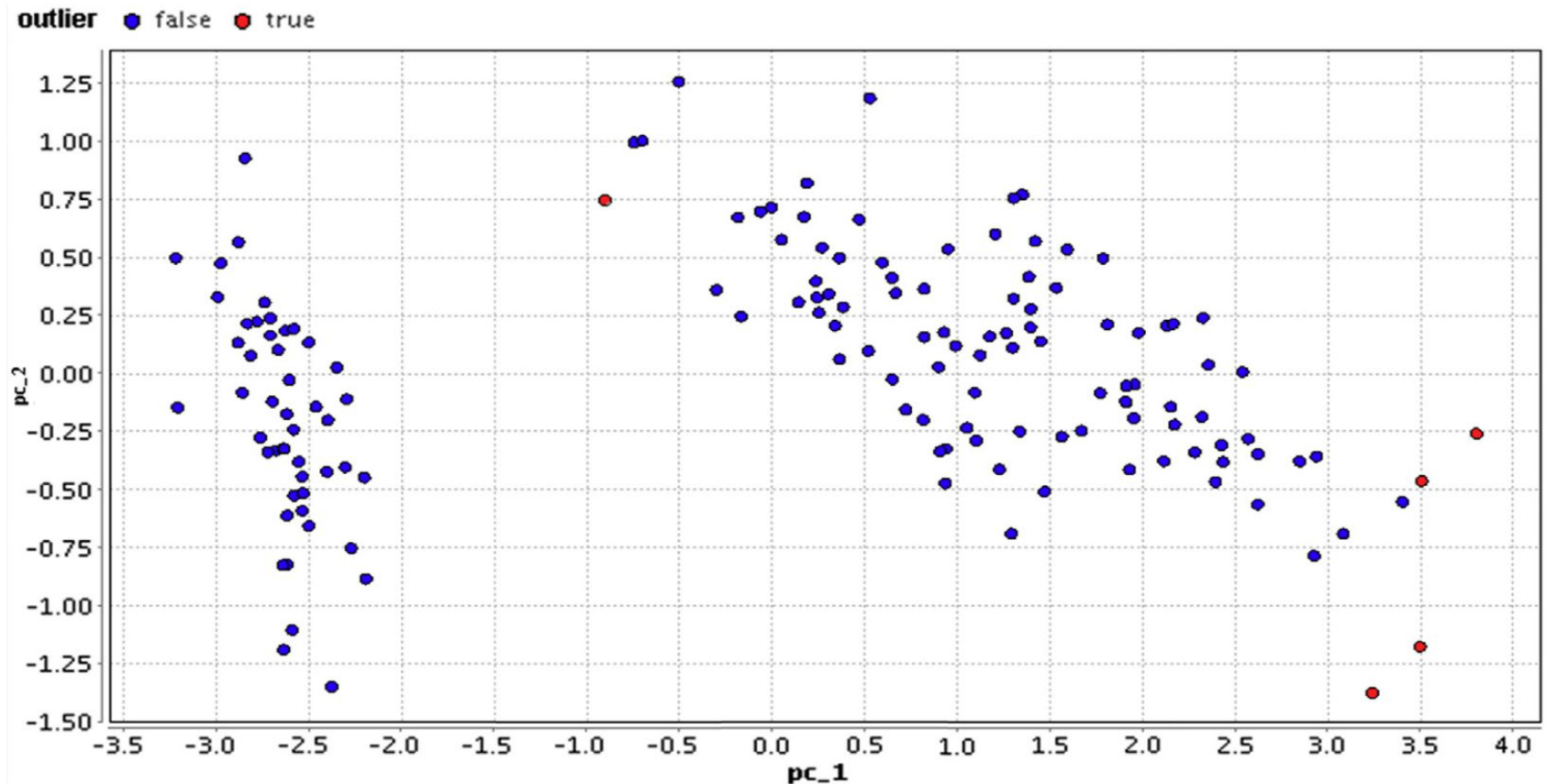
Distancia

2

Proporción

3

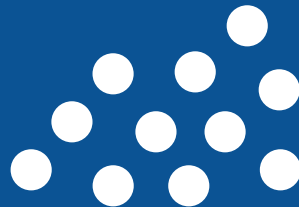
Función de **distancia**



---

**LOF**

Local Outlier Factor





## ¿Cómo detectar un outlier?

Densidad del punto



Densidad del vecindario

Densidad del punto

(Promedio de la **distancia a todos los puntos** del vecindario)<sup>-1</sup>

Densidad del vecindario

Promedio de las **densidades** de los puntos del vecindario



1

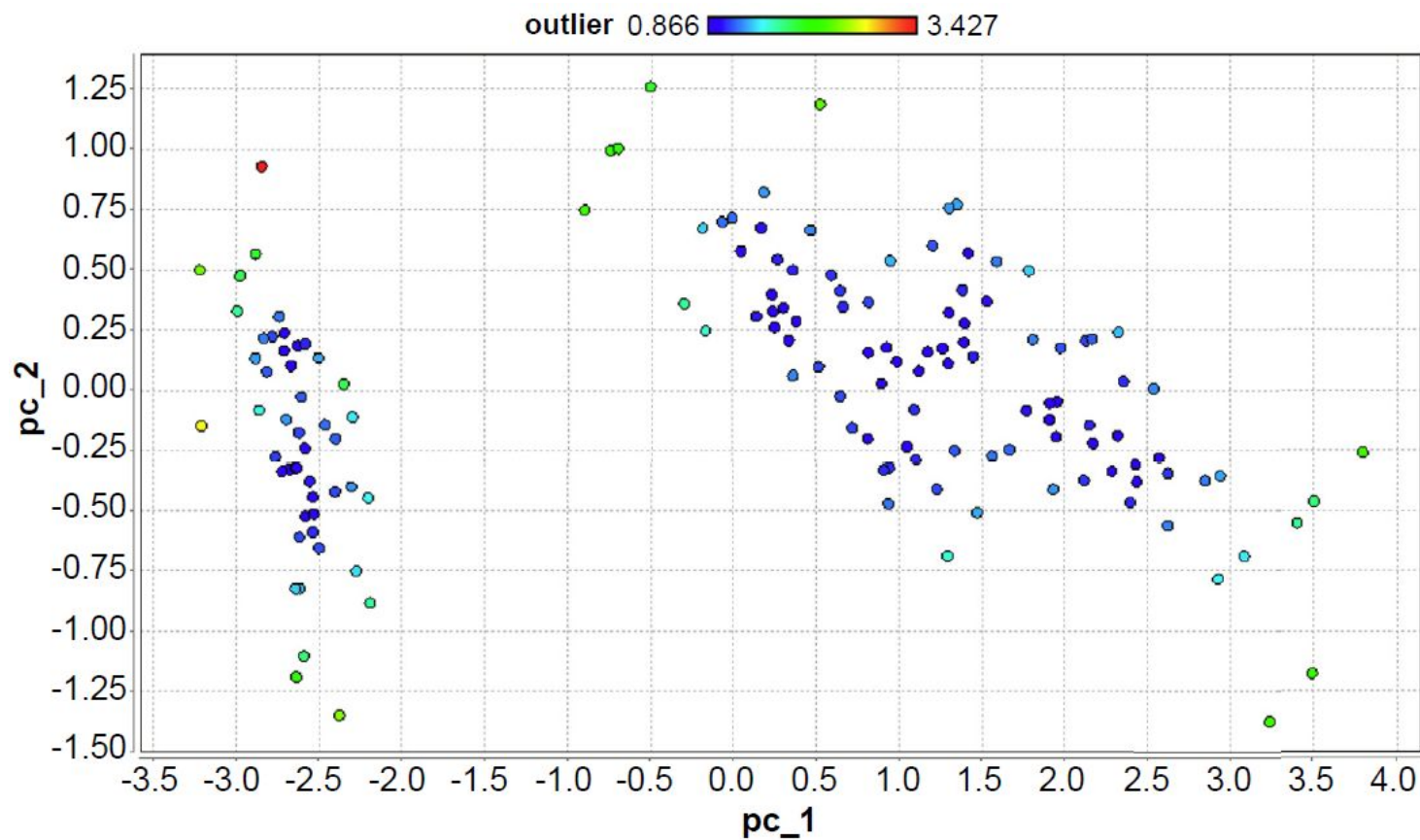
**Límite inferior** (cant. mínima de puntos)

2

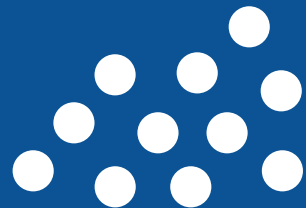
**Límite superior** (cant. mínima de puntos)

3

Función de **distancia**



—  
COF





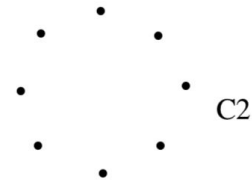


Densidad del punto



Aislamiento

C1 .....o1



---

# Ejemplo

## Hipotiroidismo





## Thyroid Disease Data Set

<http://archive.ics.uci.edu/ml/datasets/thyroid+disease>

2213 instancias

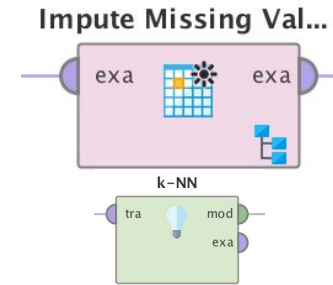
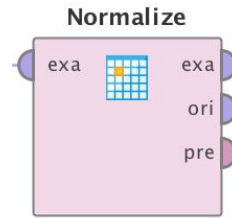
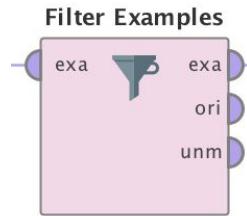
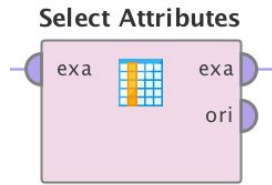
30 atributos

---

# Preprocesamiento del dataset



# Preprocesamiento



# Atributos utilizados (23)



✓ age	✓ I131 treatment	✓ TSH
✓ sex	✓ query hypothyroid	✓ T3
✓ on thyroxine	✓ query hyperthyroid	✓ TT4
✓ query on thyroxine	✓ lithium	✓ T4U
✓ on antithyroid med	✓ goitre	✓ FTI
✓ sick	✓ tumor	✓ referral source
✓ pregnant	✓ hypopituitary	✓ class
✓ thyroid surgery	✓ psych	



## Por ser irrelevantes

- ✗ TSH measured
- ✗ T3 measured
- ✗ TT4 measured
- ✗ FTI measured
- ✗ T4U measured

## Por falta de información

- ✗ TBG measured
- ✗ TBG

---

# Modelo

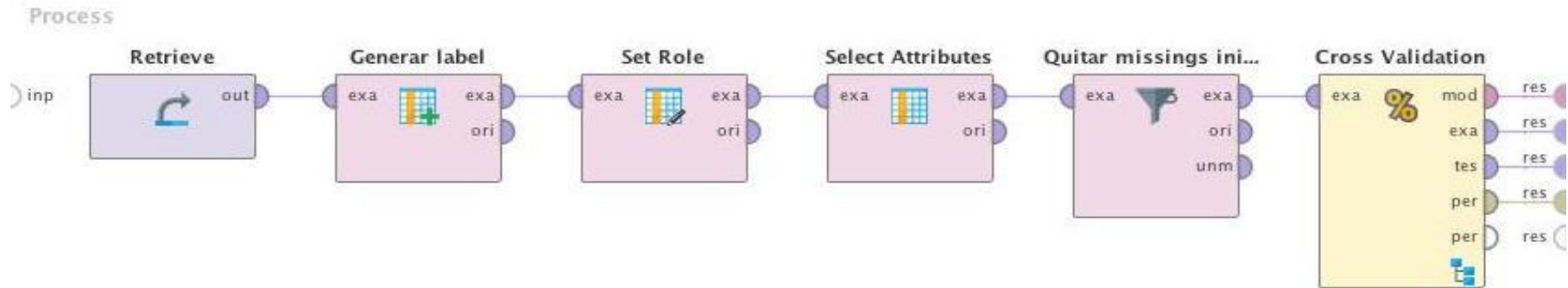
Detección de outliers





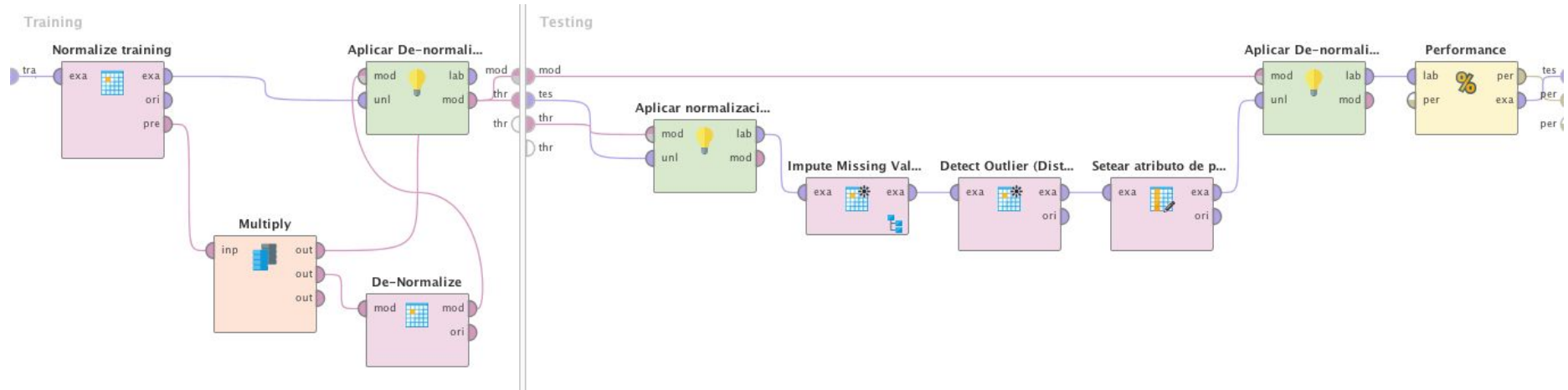


## Vista de diseño





## Vista de diseño > cross-validation



---

# Conclusiones

Detección de outliers





**Preguntas?**

—

**Muchas gracias!**

---