

# Introducción al Machine Learning

Ignacio Díaz Blanco

# Análisis de datos en procesos industriales ¡un problema *Big Data*!

- Cantidades ingentes de datos
  - Infinidad de sensores
  - Ubicuidad de la información
  - Heterogeneidad de la información
- Sistemas complejos
  - Procesos dinámicos
  - Conexión, interacción y acoplamiento
  - Muchas variables, problemas multivía
  - Interacción humana
  - Interacción con el tejido productivo

**Problema:**  
**Aprendizaje de modelos  
a partir de datos**

**Buscar estructuras  
en los datos  
y modelarlas**

## Referencias:

*Big data: The next frontier for innovation, competition, and productivity.*

[http://www.mckinsey.com/insights/business\\_technology/big\\_data\\_the\\_next\\_frontier\\_for\\_innovation.](http://www.mckinsey.com/insights/business_technology/big_data_the_next_frontier_for_innovation)



# ¿qué es el Aprendizaje Automático?

## “Machine Learning”

### Aprendizaje (definición formal)

*Un programa de computador se dice que “aprende” de la experiencia  $E$  con respecto a una serie de tareas  $T$  y un índice de rendimiento  $P$ , si dicho rendimiento  $P$  para las tareas  $T$  mejora con la experiencia  $E$ .*

*Tom Mitchell*

$T \rightarrow$  ej diagnóstico fallos (clasificación),  
estimación de variables como la temperatura en un motor (regresión),  
detección de anomalías, predicción de posición de objetos móviles en robótica etc.

$E \rightarrow$  datos de entrenamiento  $\{ \dots, (\mathbf{x}_i, y_i), \dots \}$

$P \rightarrow$  función de coste, ej. suma de errores cuadráticos

$$\sum_i (y_i - \underbrace{f(\mathbf{x}_i, \mathbf{W})}_{\hat{y}_i})^2$$

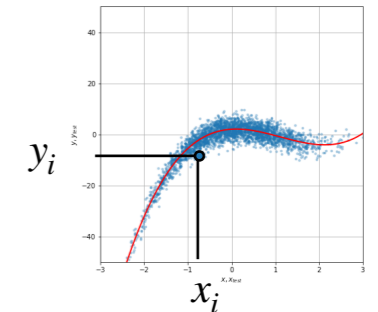


## Aprendizaje Supervisado

utiliza datos etiquetados

con entradas  $\mathbf{x}_i$  y sus correspondientes respuestas o *etiquetas*  $y_i$

$$\{\dots, (\mathbf{x}_i, y_i), \dots\}$$

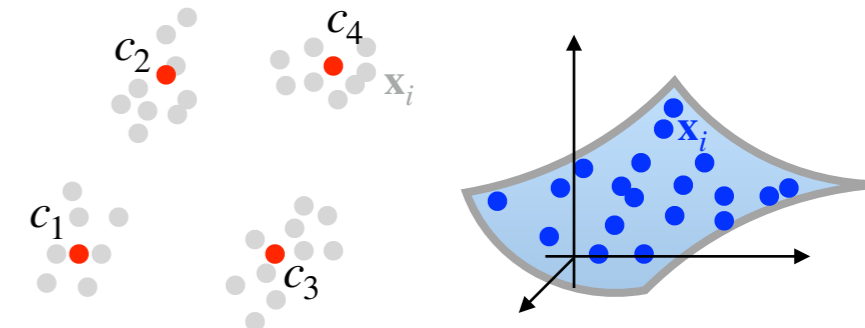


## Aprendizaje No Supervisado

utiliza datos sin etiquetar, solo las entradas:

$$\{\dots, \mathbf{x}_i, \dots\}$$

debe ser capaz de detectar patrones “interesantes”  
sin darle guías... solo con a partir de los datos



## Aprendizaje por Refuerzo

Un sistema de aprendizaje (agente) interactúa con el mundo  
y aprende a maximizar una señal de “recompensa”

# Tipos de datos de entrada

$n$  = muestra,  $k$  = timestep, descriptor,  $i, j$  = coordenadas en imagen,  $c$  = canal

Escalares  
( $n$ )

$x$

vectores  
( $n, c$ )

$(\dots, x_c, \dots)$

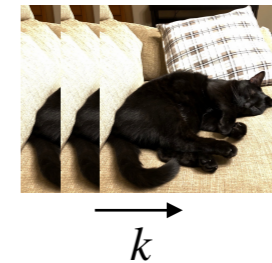
imágenes  
( $n, i, j$ )



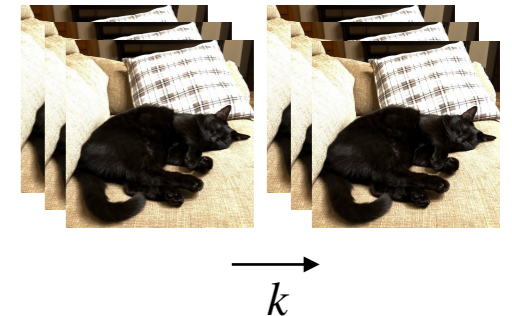
imágenes multicanal  
( $n, i, j, c$ )



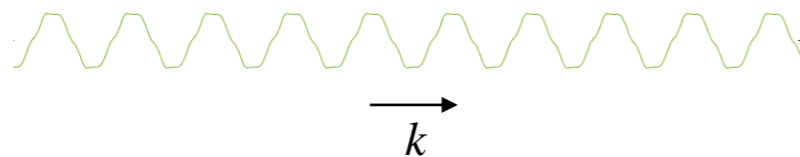
vídeos  
( $n, i, j, k$ )



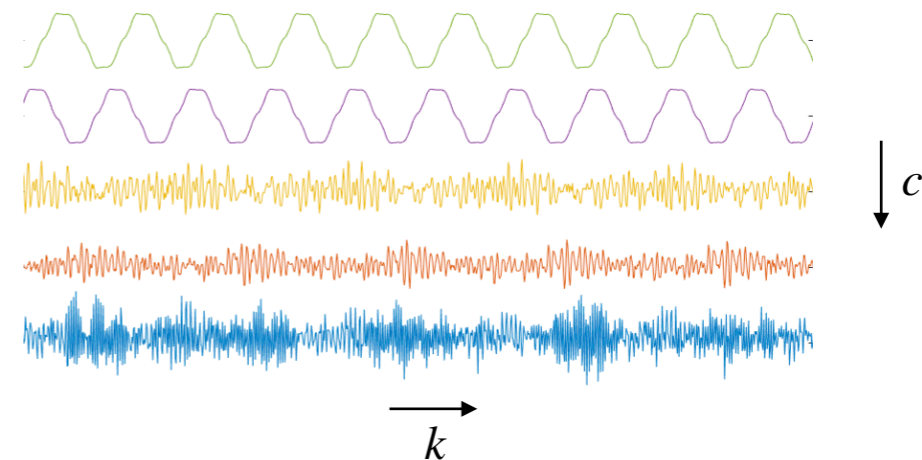
vídeos multicanal  
( $n, i, j, k, c$ )



series temporales  
( $n, k$ )



series temporales (multicanal)  
( $n, k, c$ )



## 1. Planteamiento

- ¿qué tareas  $T$  debo resolver?
- ¿qué datos  $E$  uso?
- adquirir conocimiento del problema
- visualización de datos, análisis exploratorio

## 2. Organizar los datos

- gestión de datos faltantes (eliminación, imputación)
- reducción de muestras:
  - filtrado
  - agregación
- reducción de atributos:
  - reducción de la dimensionalidad
  - selección de características relevantes
  - extracción de características (descriptores)

## 3. Elegir modelo, función de coste, regularización

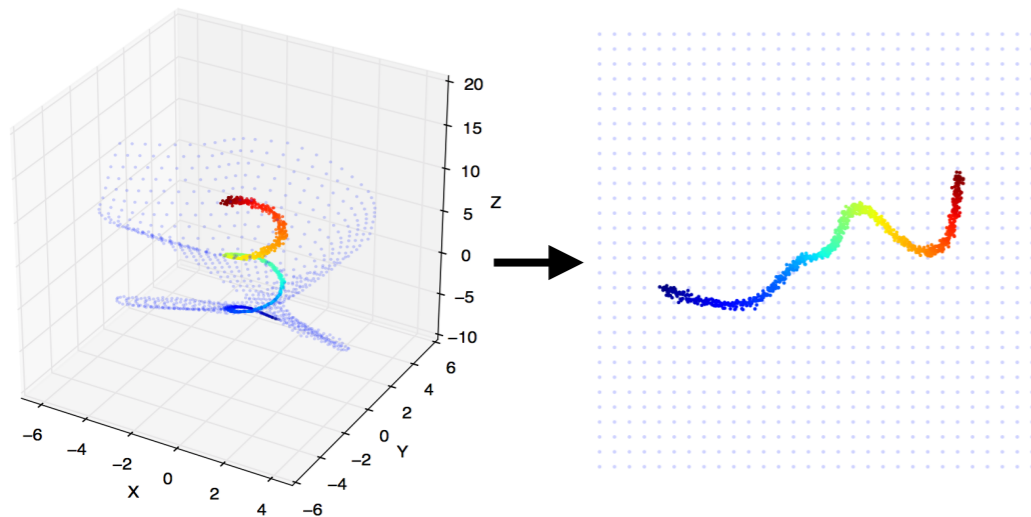
## 4. Optimización

## 5. Búsqueda de hiperparámetros

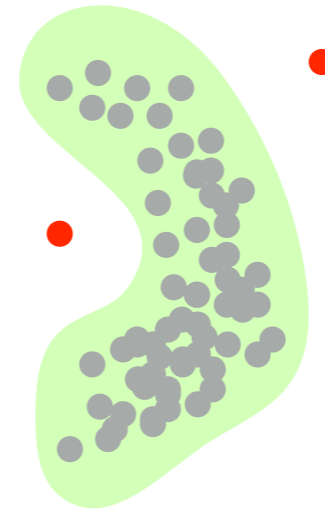
# Machine learning

detectar estructura en los datos  
desenredo de información (“disentangling”)

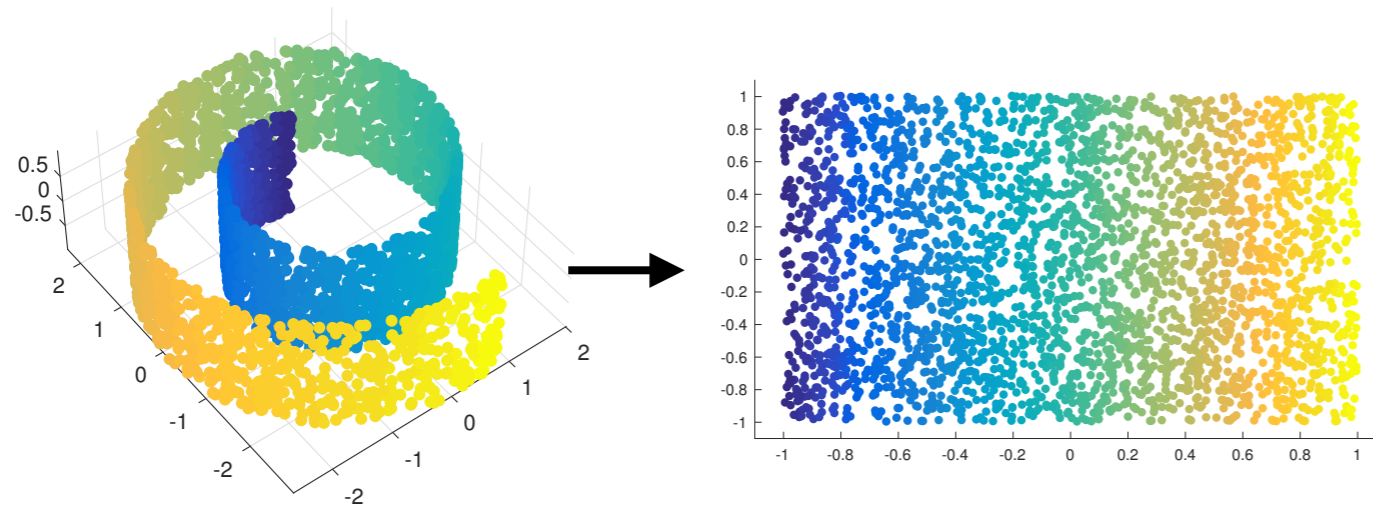
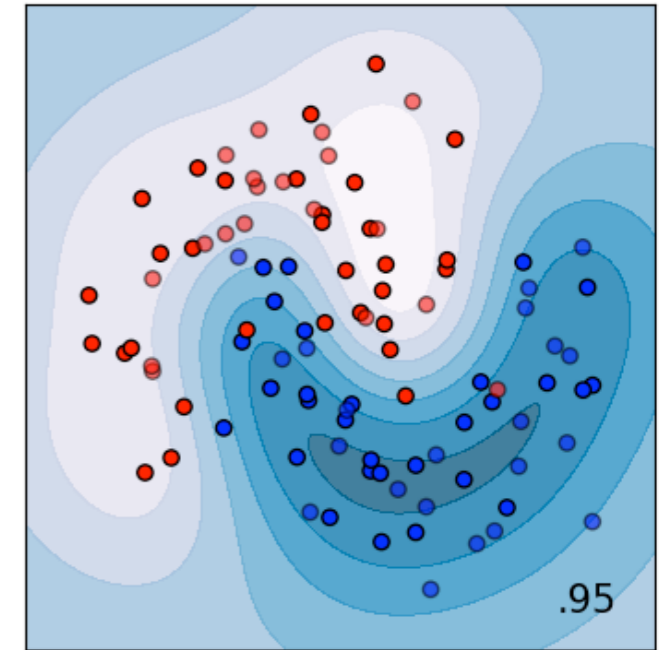
reducción de la dimensionalidad



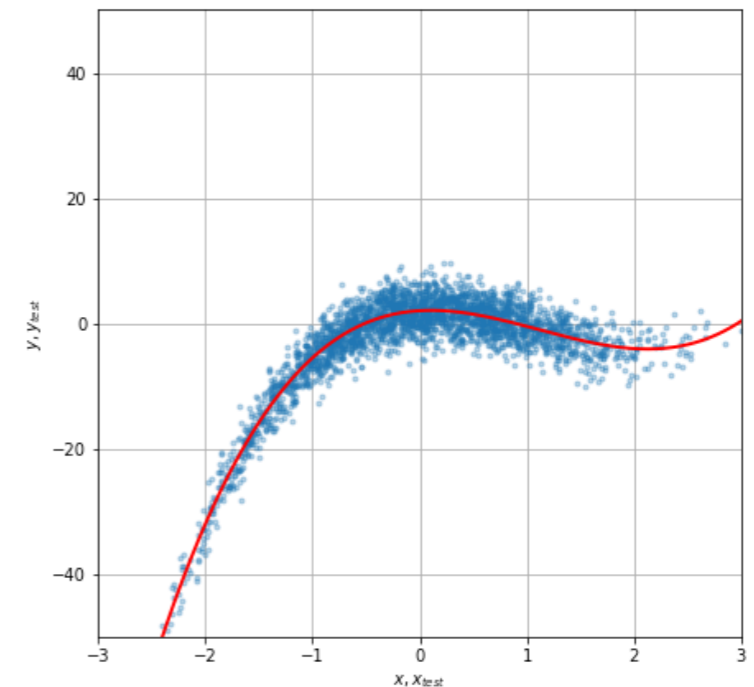
detección de anomalías



clasificación



regresión



# tareas de ML

Regresión

Clasificación

Reducción de la dimensionalidad

# Regresión



# Regresión

## Modelo y entrenamiento

$$\mathbf{y} = F(\mathbf{x}, \mathbf{W})$$

Conjunto de datos de entrenamiento

$$\{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)\}$$

clasificación:  $\mathbf{y}_i = \{(0.3,1), 0.5, 1), (-0.2, 0), \dots, (-0.3, 0)\}$

regresión:  $\mathbf{y}_i = \{(0.2, 1.1), 1.5, -1.2), (-0.23, 0.12), \dots, (-0.12, 0.43)\}$

objetivo del entrenamiento

$$W \text{ t.q. } \underbrace{\sum_i \|y_i - F(\mathbf{x}_i, \mathbf{W})\|^2}_{E(W)} \text{ sea mínimo}$$

$$\hat{\mathbf{W}} = \arg \min_{\mathbf{W}} E(\mathbf{W})$$



# Regresión

## Descenso del gradiente

Si definimos el error de un modelo como

$$E = \sum_i \| \mathbf{y}_i - F(\mathbf{x}_i, \mathbf{W}) \|^2$$

dado que  $\mathbf{x}_i, \mathbf{y}_i$  son conocidos y fijos (son los ejemplos) tenemos que

$$E = E(\mathbf{W})$$

¡El gradiente del error

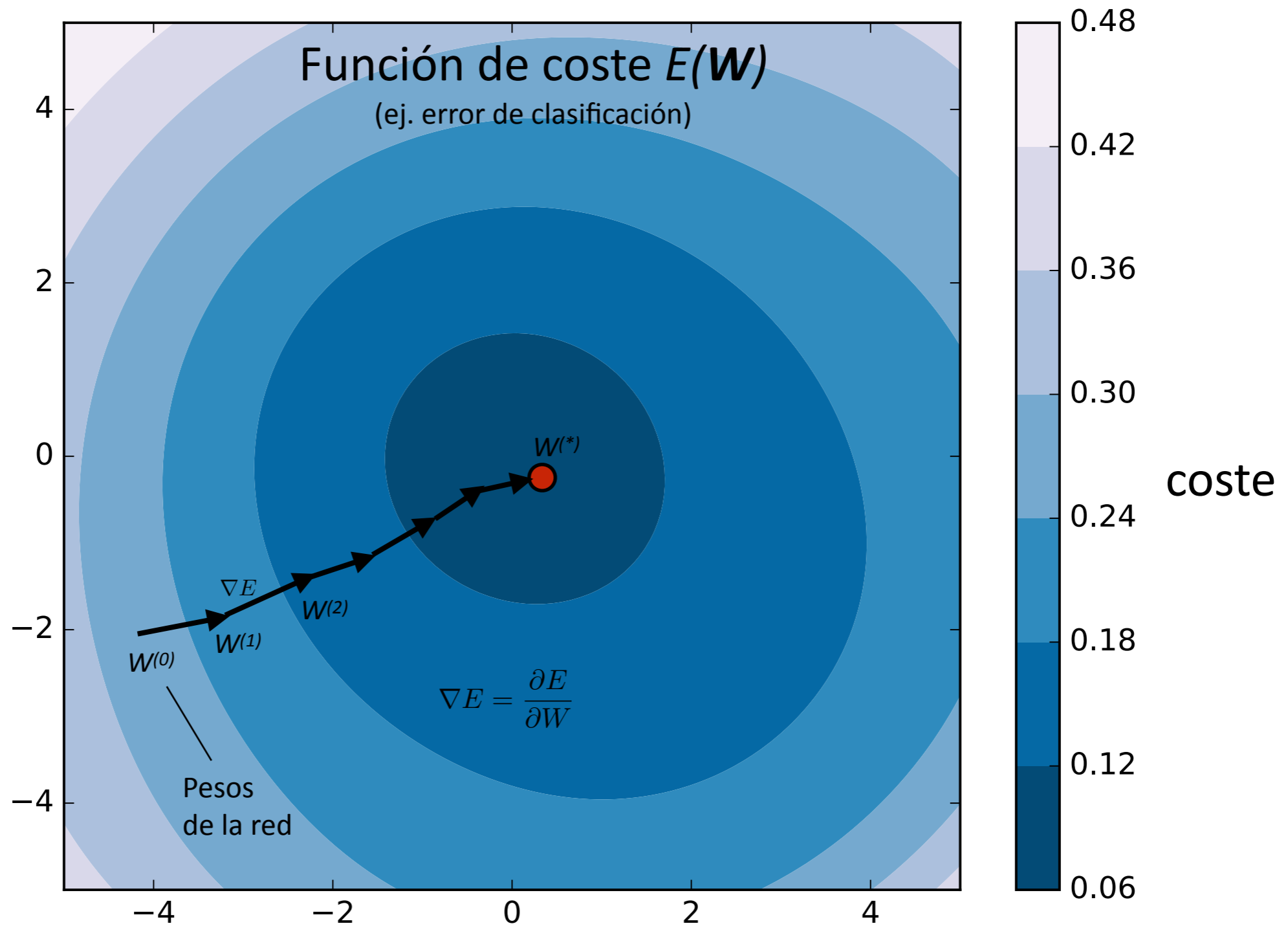
$$\nabla E = \frac{\partial E}{\partial \mathbf{W}}$$

nos dice en qué dirección debemos mover los pesos para disminuir el error!

$$\mathbf{W}^{(k)} \leftarrow \mathbf{W}^{(k-1)} - \mu \frac{\partial E}{\partial \mathbf{W}}$$

# Regresión

## Descenso del gradiente



# Regresión

## Optimización directa

Modelo lineal genérico (multivariable)

$$\mathbf{Y} = \mathbf{X}\mathbf{W}$$

datos de salida (targets)

$$\mathbf{Y} = \begin{pmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_n \end{pmatrix} \in \mathbb{R}^{nq}$$

pesos (coeficientes del modelo)

$$\mathbf{W} \in \mathbb{R}^{pq}$$

datos de entrada

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{pmatrix} \in \mathbb{R}^{np}$$

### Ridge Regression

$$\hat{\mathbf{W}} = (\mathbf{X}^T \mathbf{X} - \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{Y}$$

$$\hat{\mathbf{W}} = \arg \min_{\mathbf{W}} \|\mathbf{X}\mathbf{W} - \mathbf{Y}\|^2 + \lambda \|\mathbf{W}\|^2$$

↑  
regularización  
de Tikhonov

ejemplo en Python ↓

```
from sklearn.linear_model import Ridge
import numpy as np
n_samples, n_features = 10, 5
rng = np.random.RandomState(0)
y = rng.randn(n_samples)
X = rng.randn(n_samples, n_features)
clf = Ridge(alpha=1.0)
clf.fit(X, y)
```

# Regresión

## Optimización directa (regresión dispersa)

### Lasso

$$\hat{\mathbf{W}} = \arg \min_{\mathbf{W}} \|\mathbf{XW} - \mathbf{Y}\|^2 + \alpha \|\mathbf{W}\|_1$$

```
from sklearn import linear_model
clf = linear_model.Lasso(alpha=0.1)
clf.fit([[0,0], [1, 1], [2, 2]], [0, 1, 2])
print(clf.coef_)
print(clf.intercept_)
```

← ejemplo en Python

### Elastic Net

$$\hat{\mathbf{W}} = \arg \min_{\mathbf{W}} \|\mathbf{XW} - \mathbf{Y}\|^2 + \lambda_1 \|\mathbf{W}\|_1 + \lambda_2 \|\mathbf{W}\|^2$$

```
X, y = make_regression(n_features=2, random_state=0)
regr = ElasticNet(random_state=0)
regr.fit(X, y)
print(regr.coef_)
print(regr.intercept_)
print(regr.predict([[0, 0]]))
```

← ejemplo en Python

# Regresión

## modelos multivariable

Modelo multivariable  
y depende de varios factores  $x_1, \dots, x_p$  (regresores)

$$y = \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_n x_n + \beta$$

$$\mathbf{Y} = \mathbf{XW}$$

$$\mathbf{Y} = \begin{pmatrix} y^1 \\ y^2 \\ \vdots \\ y^n \end{pmatrix}$$

salidas

$$\mathbf{X} = \begin{pmatrix} x_1^1 & x_2^1 & \dots & x_p^1 & 1 \\ x_1^2 & x_2^2 & \dots & x_p^2 & 1 \\ \vdots & & & & \\ x_1^n & x_2^n & \dots & x_p^n & 1 \end{pmatrix}$$

$p$  regresores

↑  
término  
independiente  
(bias)

$$\mathbf{W} = \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_p \\ \beta \end{pmatrix}$$

$p + 1$  coeficientes

# Regresión

## modelos no lineales

Modelo no lineal (lineal en los parámetros)

y como combinación lineal de funciones no lineales de  $x_1, \dots, x_p$

$$y = \alpha_1 f_1(\mathbf{x}) + \alpha_2 f_2(\mathbf{x}) + \dots + \alpha_n f_n(\mathbf{x}) + \beta$$

$$\mathbf{Y} = \mathbf{XW}$$

regresores no lineales  
 $f_i(\mathbf{x}) = f_i(x_1, x_2, \dots, x_p)$

ejemplos de  $f_i(\mathbf{x})$  :

$x_3 x_5, \cos(x_3), \cos(5x_4),$

$\sqrt{(x_1 x_4)}, x_1^2 + x_3^2,$

$x_1 x_2 x_3^2 x_4 x_5$

$$\mathbf{Y} = \begin{pmatrix} y^1 \\ y^2 \\ \vdots \\ y^n \end{pmatrix}$$

salidas

$$\mathbf{X} = \begin{pmatrix} f_1(\mathbf{x}) & f_2(\mathbf{x}) & \dots & f_p(\mathbf{x}) & 1 \\ f_1(\mathbf{x}) & f_2(\mathbf{x}) & \dots & f_p(\mathbf{x}) & 1 \\ \vdots & & & & \\ f_1(\mathbf{x}) & f_2(\mathbf{x}) & \dots & f_p(\mathbf{x}) & 1 \end{pmatrix}$$

$p$  regresores

↑  
término  
independiente  
(bias)

$$\mathbf{W} = \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_p \\ \beta \end{pmatrix}$$

$p + 1$  coeficientes

# Identificación de Sistemas

## Modelos dinámicos lineales (discretos)

Función de transferencia (SISO discreto)

$$G(z) = \frac{Y(z)}{U(z)} = \frac{b_0 + b_1z^{-1} + \dots + b_mz^{-m}}{1 - a_1z^{-1} - \dots - a_nz^{-n}}$$

ecuación en diferencias

$$y(k) = a_1y(k-1) + \dots + a_ny(k-n) + b_0u(k) + \dots + b_mu(k-m)$$

modelo de regresión multivariable

$$\begin{bmatrix} \vdots \\ y(k) \\ \vdots \end{bmatrix} = \begin{bmatrix} \vdots & & & & & \\ y(k-1) & \cdots & y(k-n) & u(k) & \cdots & u(k-m) \\ \vdots & & & & & \end{bmatrix} \begin{bmatrix} -a_1 \\ \vdots \\ -a_n \\ b_0 \\ \vdots \\ b_m \end{bmatrix}$$

→  
se suponen  
disponibles  
datos para  
muchos  
valores de  $k$

$$\mathbf{Y} = \mathbf{XW}$$

# Identificación de Sistemas

## Modelos dinámicos lineales (continuos)

si se conocen, utilizar  
las derivadas como regresores

sistema SISO continuo

$\{x(t)\} \rightarrow \{y(t)\}$

$$y^{(m)}(t) = \alpha_{m-1}y^{(m-1)}(t) + \dots + \alpha_1\dot{x}(t) + \beta_r u^{(r)}(t) + \dots + b_1\dot{u}(t) + \beta_0 u(t) + \gamma$$

regresores con  
las derivadas

**Problema práctico**  
muy sensible al ruido por las derivadas



# Identificación de Sistemas

## Modelos dinámicos no lineales (NARX)

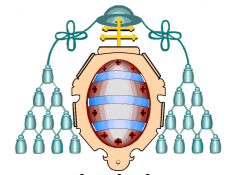
modelo NARX (*Non-linear Autoregressive eXogenous*)

$$y(k) = f(y(k-1), \dots, y(k-n_y), u(k), \dots, u(k-n_u))$$

en general, puede aproximarse con regresores  
de uno o varios retardos  
una o varias variables de salida / entrada

$$\dots + \alpha_1 y(k-2)y(k-5) + \alpha_2 y(k-6)u(k) + \alpha_3 \cos(y(k)^2 u(k)) + \dots$$

# Identificación de Sistemas espacio de estados (SINDy)



**Discovering governing equations from data by sparse identification of nonlinear dynamical systems**

Steven L. Brunton<sup>1</sup>, Joshua L. Proctor<sup>2</sup>, and J. Nathan Kutz<sup>3</sup>

<sup>1</sup>Department of Mechanical Engineering, University of Washington, Seattle, WA 98195; <sup>2</sup>Institute for Disease Modeling, Bellevue, WA 98005; and <sup>3</sup>Department of Applied Mathematics, University of Washington, Seattle, WA 98195

Edited by William Bielik, Princeton University, Princeton, NJ, and approved March 1, 2016 (received for review August 31, 2015)

Extracting governing equations from data is a central challenge in many diverse areas of science and engineering. Data are abundant whereas models often remain elusive, as in climate science, neuroscience, ecology, finance, and epidemiology, to name only a few examples. In this work, we combine sparsity-promoting techniques and machine learning with nonlinear dynamical systems to discover governing equations from noisy measurement data. The only assumption about the structure of the model is that there are only a few important terms that govern the dynamics, so that the equations are sparse in the space of possible functions; this assumption holds for many physical systems in an appropriate basis. In particular, we use sparse regression to determine the fewest terms in the dynamic governing equations required to accurately represent the data. This results in parsimonious models that balance accuracy with model complexity to avoid overfitting. We demonstrate the algorithm on a wide range of problems, from simple canonical systems, including linear and nonlinear oscillators and the chaotic Lorenz system, to the fluid vortex shedding behind an obstacle. The fluid example illustrates the ability of this method to discover the underlying dynamics of a system that took experts in the community nearly 30 years to resolve. We also show that this method generalizes to parameterized systems and systems that are time-varying or have external forcing.

**Significance**

Understanding dynamic constraints and balances in nature has facilitated rapid development of knowledge and enabled technology, including aircraft, construction engines, satellites, and electrical power. This work develops a novel framework to discover governing equations underlying a dynamical system simply from data measurements, leveraging advances in sparsity techniques and machine learning. The resulting models are parsimonious, balancing model complexity with descriptive ability while avoiding overfitting. There are many critical data-driven problems, such as understanding cognition from neural recordings, inferring climate patterns, determining stability of financial markets, predicting and suppressing the spread of disease, and controlling turbulence for greener transportation and energy. With abundant data and elusive laws, data-driven discovery of dynamics will continue to play an important role in these efforts.

Author contributions: S.L.B., J.L.P., and J.N.K. designed research; S.L.B. performed research; S.L.B., J.L.P., and J.N.K. analyzed data; and S.L.B. wrote the paper. The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

Freely available online through the PNAS open access option.

To whom correspondence should be addressed: steve@brunton.wa.edu.

This article contains supporting information online at www.pnas.org/lookup/suppl/doi:10.1073/pnas.1517384113/-DCSupplemental.

3932-3937 | PNAS | April 12, 2016 | vol. 113 | no. 15 | www.pnas.org/cgi/doi/10.1073/pnas.1517384113

$$\frac{d}{dt}\mathbf{x} = f(\mathbf{x})$$

sistema  
no lineal  
genérico

$$f(\mathbf{x}) \approx \sum_k \theta_k(\mathbf{x}) \xi_k = \Theta(\mathbf{x}) \xi$$

aproximación  
mediante  
combinación lineal  
de funciones no lineales

Brunton, S. L., Proctor, J. L., & Kutz, J. N. (2016). Discovering governing equations from data by sparse identification of nonlinear dynamical systems. *Proceedings of the national academy of sciences*, 113(15), 3932-3937.

datos y targets  
para  $m$  muestras  
en formato matricial

expansión no lineal  
de los datos  
en formato matricial

nos lleva a un  
problema de regresión  
lineal multivariable

$$\mathbf{X} = [\mathbf{x}(t_1), \mathbf{x}(t_2), \dots, \mathbf{x}(t_m)]$$

$$\dot{\mathbf{X}} = [\dot{\mathbf{x}}(t_1), \dot{\mathbf{x}}(t_2), \dots, \dot{\mathbf{x}}(t_m)]$$

$$\Theta(\mathbf{X}) = [\mathbf{1}, \mathbf{X}, \mathbf{X}^2, \dots, \cos(\mathbf{X}) \dots]$$

$$\dot{\mathbf{X}} = \Theta(\mathbf{X}) \mathbf{\Xi}$$

# Identificación de Sistemas

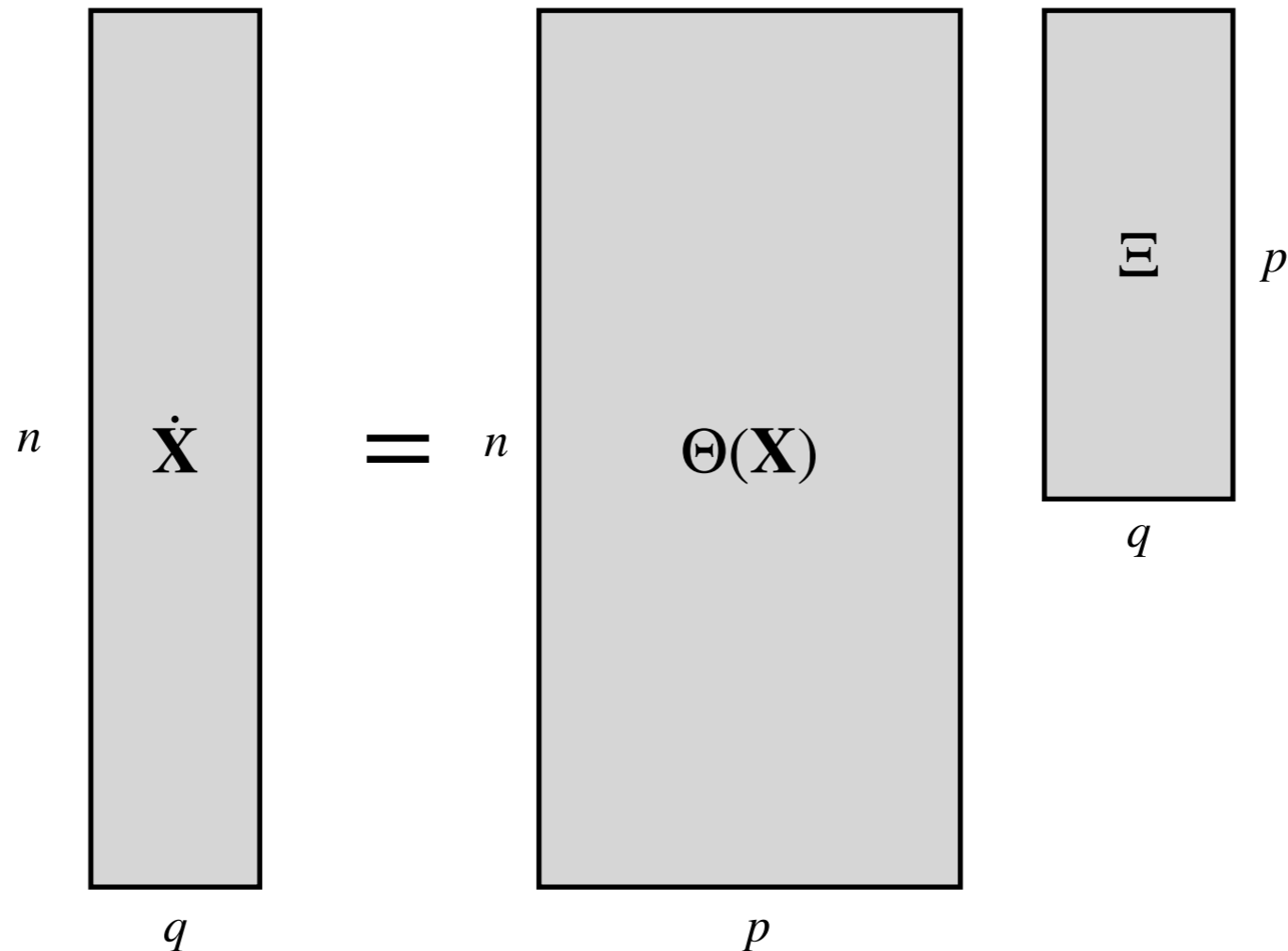
## espacio de estados (SINDy)

(Brunton & Kutz, 2022)

$p =$   
número de  
funciones no lineales  
de la expansión

$q =$   
número de  
variables  
de estado

$n =$   
muestras  
disponibles



# Identificación de Sistemas

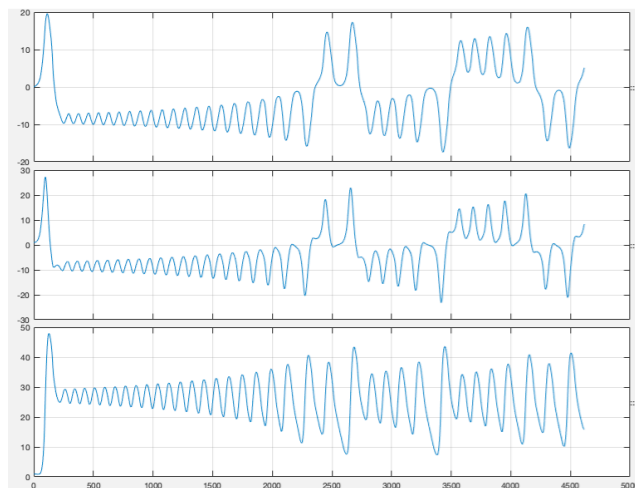
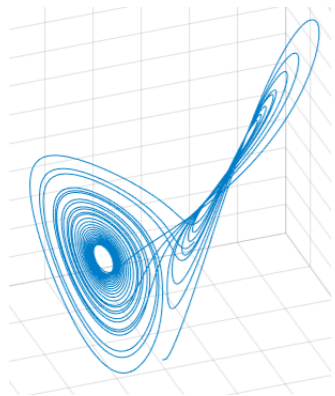
## espacio de estados (SINDy)

Ecuaciones de Lorenz

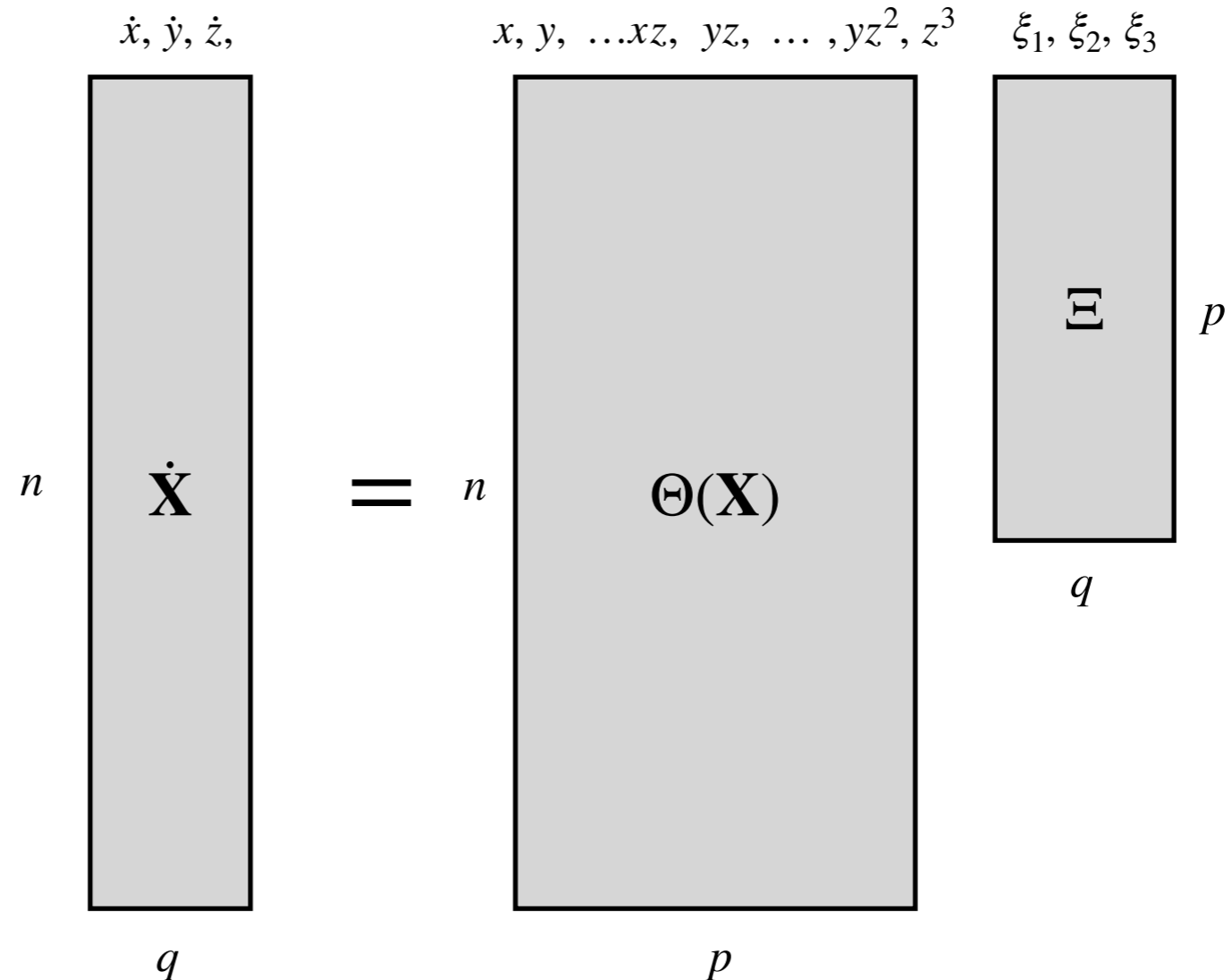
$$\frac{dx}{dt} = \sigma(y - x)$$

$$\frac{dy}{dt} = x(\rho - z) - y$$

$$\frac{dz}{dt} = xy - \beta z$$

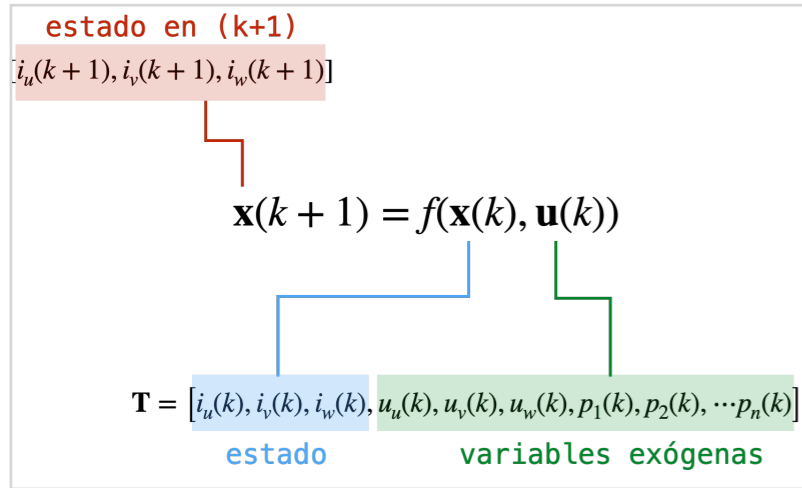


$$\Theta(\mathbf{x}) = x, y, z, x^2, xy, xz, y^2, yz, z^2, x^3, x^2y, x^2z, xy^2, xyz, xz^2, y^3, y^2z, yz^2, z^3$$

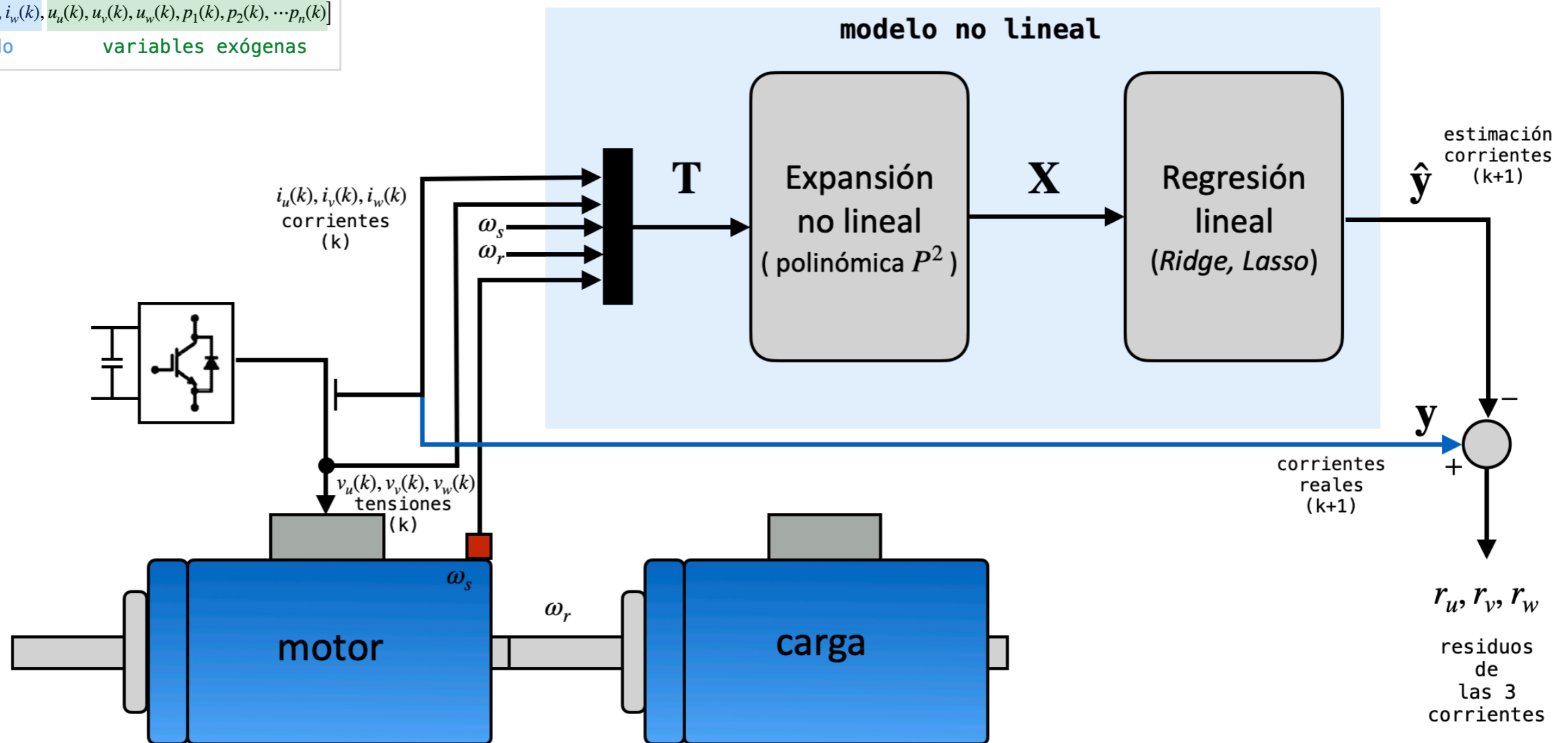


# Identificación de Sistemas

## espacio de estados (SINDy)

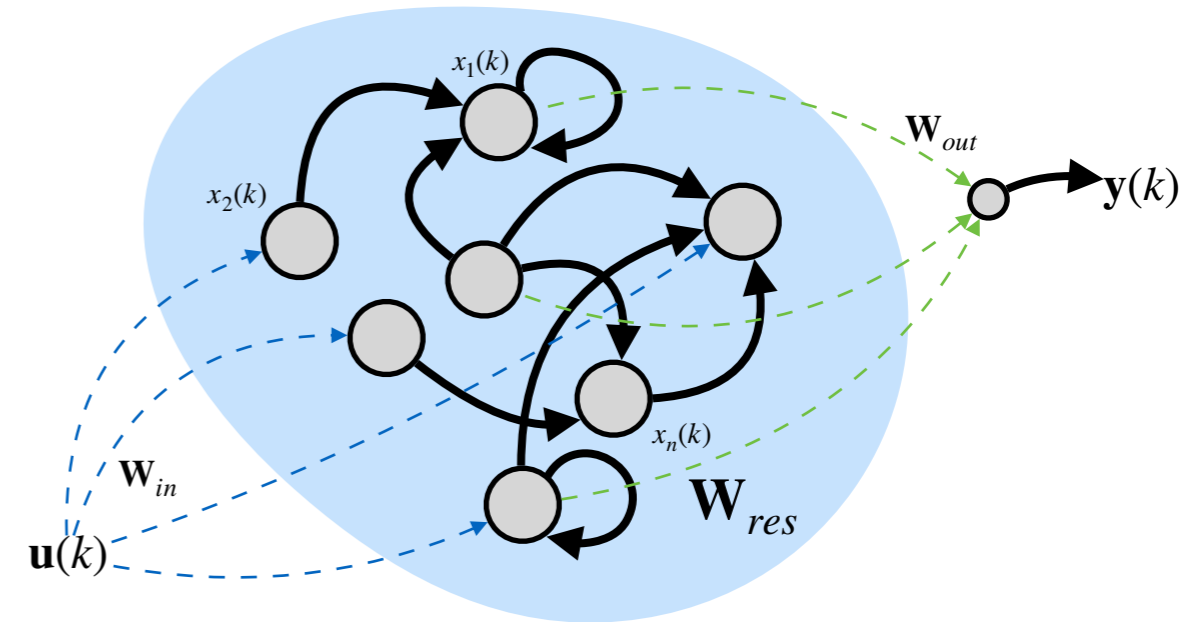


### Ejemplo: detección de anomalías



$$\mathbf{x}(k) = \sigma(\mathbf{W}_{res}\mathbf{x}(k-1) + \mathbf{W}_{in}\mathbf{u}(k))$$

$$\mathbf{y}(k) = \mathbf{W}_{out}\mathbf{x}(k)$$

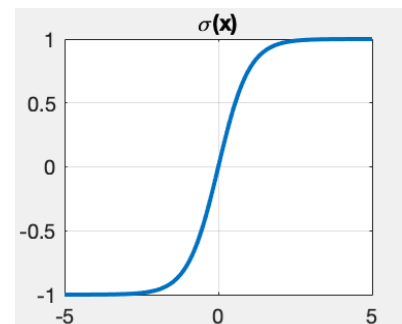
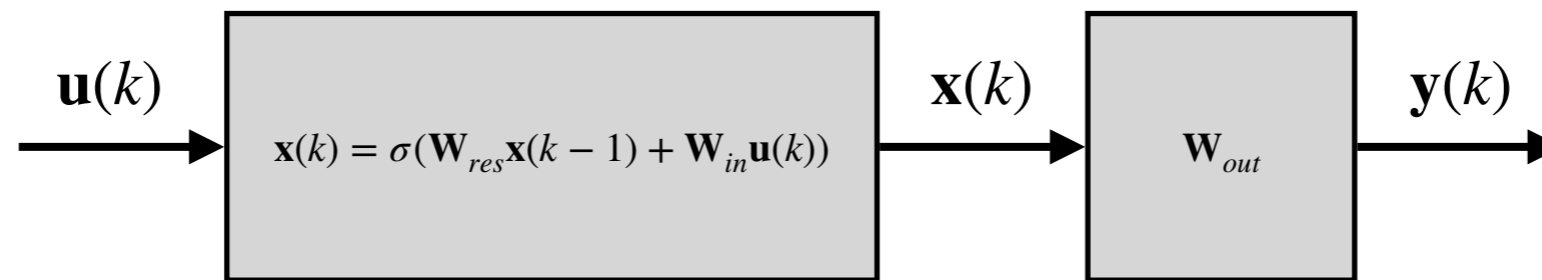


estado  
interno (de alta  
dimensionalidad)

ecuación de estado

$\mathbb{R}^n$

ecuación  
de observación

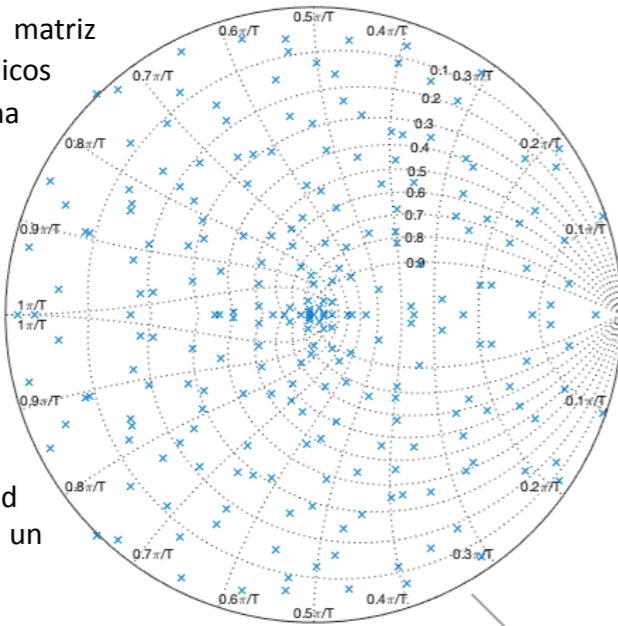




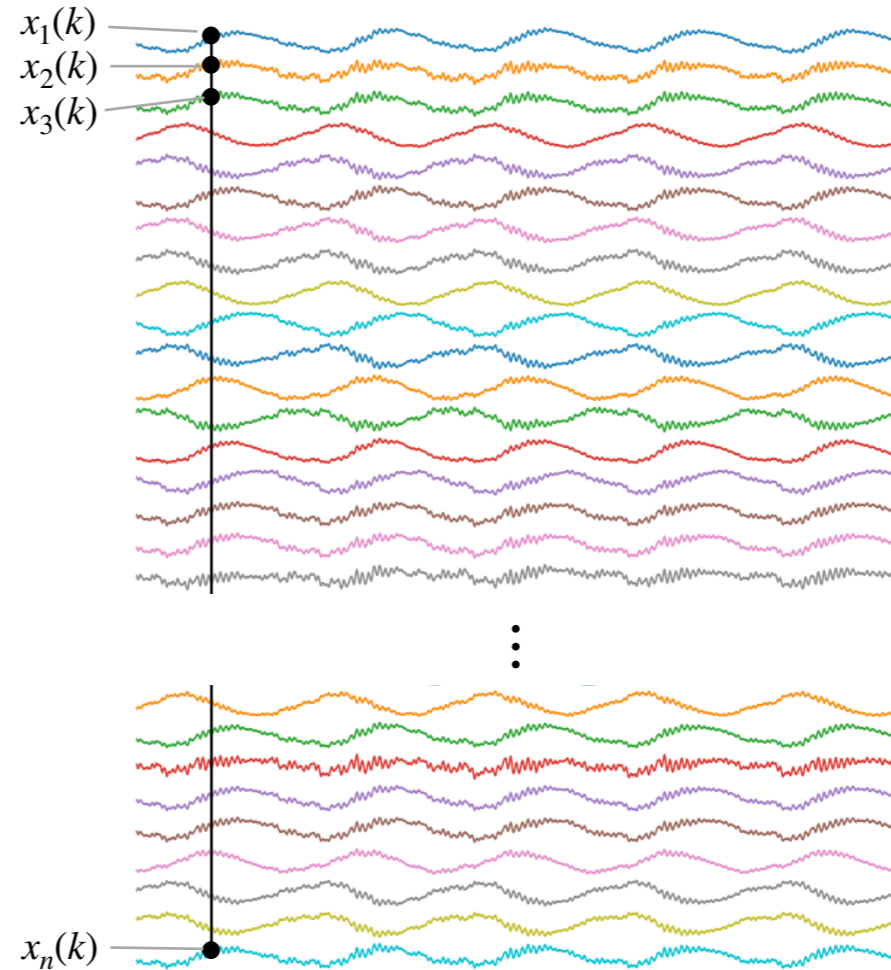
## Echo State Networks

los valores propios de la matriz  $\mathbf{W}_{res}$  son los modos dinámicos del sistema si  $\sigma(\cdot)$  es una función lineal. Influyen en la naturaleza de los modos dinámicos no lineales cuando  $\sigma$  es sigmoideal.

Para generar un amplio espectro de dinámicas se hace  $\mathbf{W}_{res}$  dispersa (mayoría de ceros) y se lleva al borde de estabilidad escalándola para alcanzar un radio espectral  $\rho \approx 1$



matriz de reservorio,  $\mathbf{X}$  (estados  $\times$  muestras)



### Reservorio

Contiene un número elevado  $n$  de modos dinámicos

$$x_1(k), x_2(k), \dots, x_n(k)$$

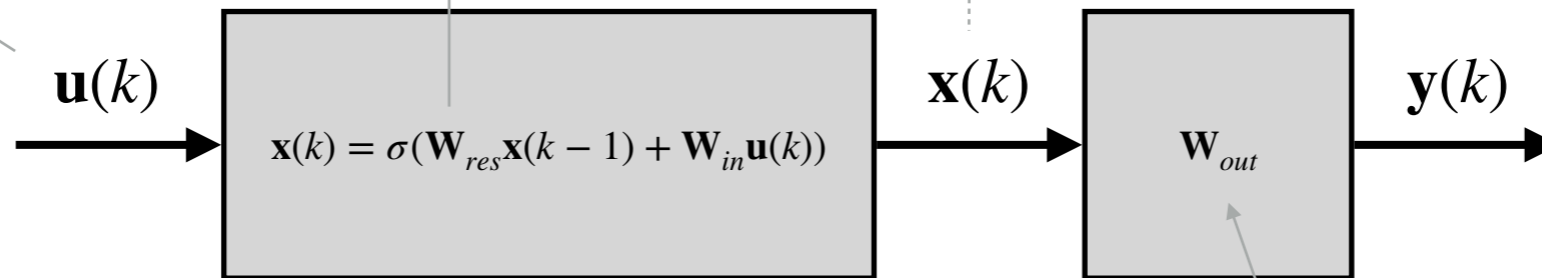
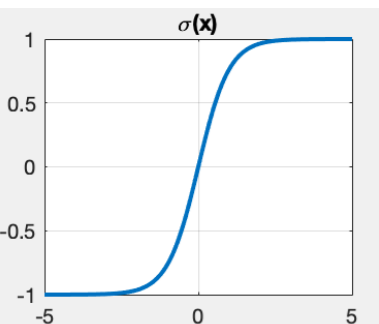
generados no linealmente a partir de la entrada  $\mathbf{u}(k)$

puede verse como una expansión no lineal de  $\mathbf{u}(k)$

que la "disecciona" en muchas componentes no lineales

### Entrada (excitación) $\mathbf{u}(k)$

señal exógena (puede ser multivariable)



### Salida

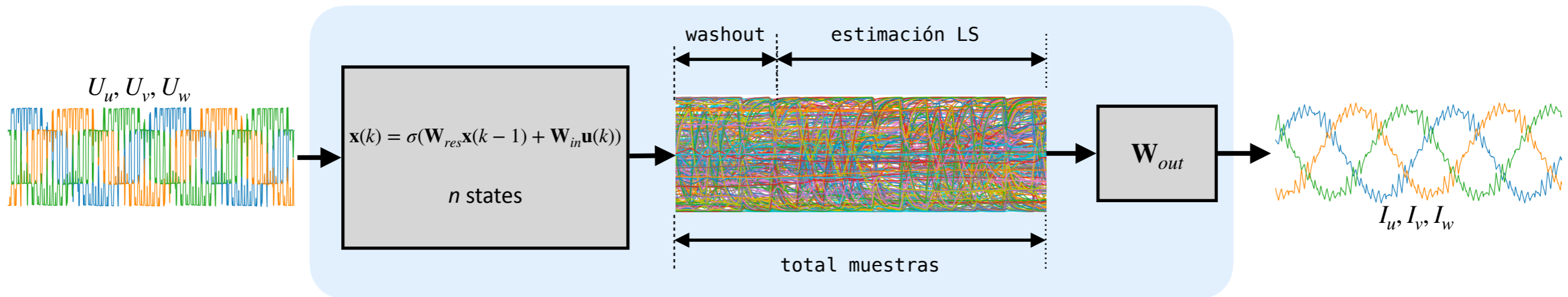
Cuando es dato (problema de identificación) permite estimar  $\mathbf{W}_{out}$

$$\mathbf{W}_{out} = (\mathbf{X}\mathbf{X}^T + \lambda\mathbf{I})^{-1}\mathbf{X}\mathbf{Y}^T$$

# Identificación de Sistemas

## Echo State Networks

Ejemplo: modelado tensión corriente





# Identificación de Sistemas

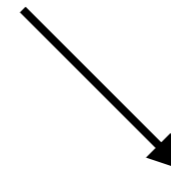
## Dificultad de la generalización

de todos modos, ¡no suele ser tan fácil!

funciona bien en un punto de trabajo,  
pero estos modelos suelen ser muy sensibles  
si se cambia el punto de operación



usar modelos  
más generales



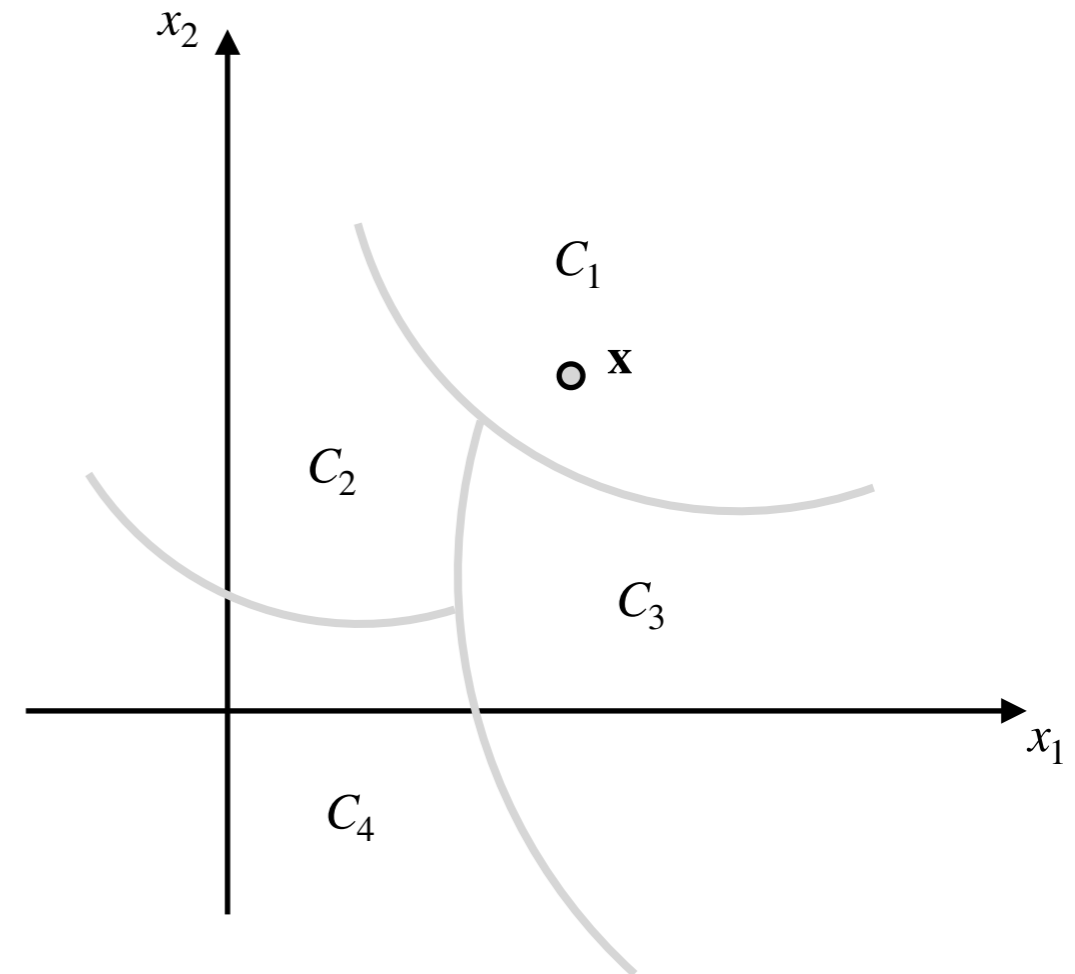
arquitecturas con modelos  
en varios puntos de trabajo

# Clasificación

Objetivo:

asignar una nueva observación  $\mathbf{x}$  a una clase  $C_i$  de entre un conjunto de clases  $\{C_1, C_2, \dots, C_n\}$

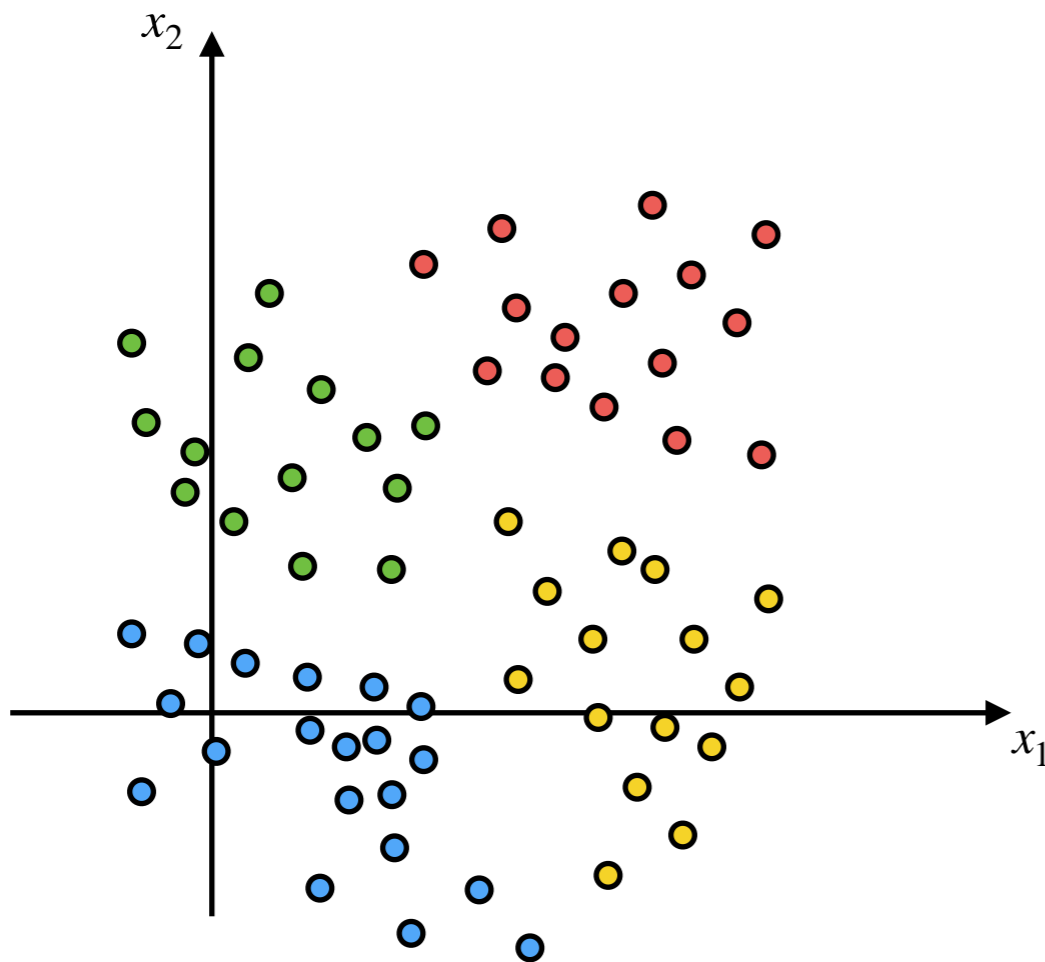
En general,  $\mathbf{x} = (x_1, x_2, \dots, x_p)$  suele ser un vector de características (*feature vector*) con descriptores  $x_i$  sensibles a las condiciones que producen el fallo



Se parte de datos etiquetados

$$\{(\mathbf{x}_1, c_1), \dots, (\mathbf{x}_n, c_n)\}$$

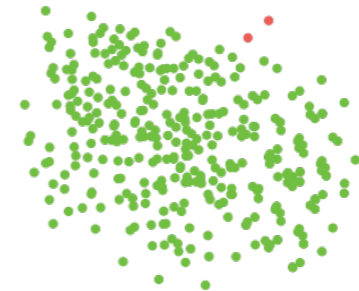
$$c_i = \{\text{rojo, verde, amarillo, azul}\}$$



¡ requerir datos etiquetados  
es una **exigencia “fuerte”** !

### A. Desbalanceo

Puede haber pocos datos etiquetados de una clase y muchos de otra



métodos  
no supervisados

### B. Falta de etiquetas

No siempre se dispone de datos para todos los fallos posibles



métodos  
semisupervisados

# Clasificador Naïve Bayes

## teorema de Bayes

### Posibilidades para estimar los componentes

(i no las únicas !)

$$p(C_i) \approx \frac{\text{num datos de fallo}}{\text{num datos totales}} \quad (\text{o si los hay, estadísticos históricos de frecuencia del fallo})$$

$$p(\mathbf{x} | C_i) \approx \text{normpdf}(\mathbf{x}, \mu_i, \sigma_i)$$

asumir distribución normal  
calculando la media  $\mu$  y la desviación típica  $\sigma$   
muestrales de los datos de la clase  $C_i$

$$p(\mathbf{x}) = \sum_{k=1}^n p(\mathbf{x} | C_k)P(C_k) \quad (\text{teorema de probabilidad total})$$

“prior”  
probabilidad  
de pertenecer a la clase  $C_i$   
sin otra consideración

densidad de probabilidad  
condicional de  $\mathbf{x}$   
sabiendo que pertenece a  $C_i$

“posterior”  
probabilidad de pertenecer  
a la clase  $C_i$  conocida  $\mathbf{x}$

$$P(C_i | \mathbf{x}) = \frac{p(\mathbf{x} | C_i)P(C_i)}{p(\mathbf{x})} = \frac{p(\mathbf{x} | C_i)P(C_i)}{\sum_{k=1}^n p(\mathbf{x} | C_k)P(C_k)}$$

densidad de probabilidad  
de  $\mathbf{x}$ , sea de la clase que sea

# Clasificador Naïve Bayes

## ejemplo en Matlab

```
% TRAINING DATA
n1 = 100;
n2 = 200;
p1 = randn(n1,1)*2+1;
p2 = randn(n2,1)*0.5-1;
p=[p1;p2];

% PRIOR PROBABILITIES FOR CLASS C1, C2
PC1 = n1/(n1+n2);
PC2 = n2/(n1+n2);

% TEST DATA
x=linspace(-10,10,1000);

% ESTIMATE THE PARAMETERS FOR THE DISTRIBUTIONS (suppose normally distributed)
mu1e = mean(p1);
mu2e = mean(p2);
sigma1e = std(p1);
sigma2e = std(p2);

% OBTAIN THE CONDITIONAL PROBABILITY DENSITIES FOR EACH CLASS P(X|Ci)
pXC1=normpdf(x,mu1e,sigma1e);
pXC2=normpdf(x,mu2e,sigma2e);

% TOTAL PROBABILITY DENSITY FUNCTION
px=pXC1*PC1+pXC2*PC2;

% APPLY BAYES RULE TO OBTAIN POSTERIOR PDF p(Ci|x)
PC1x=pXC1*PC1./px;
PC2x=pXC2*PC2./px;
```

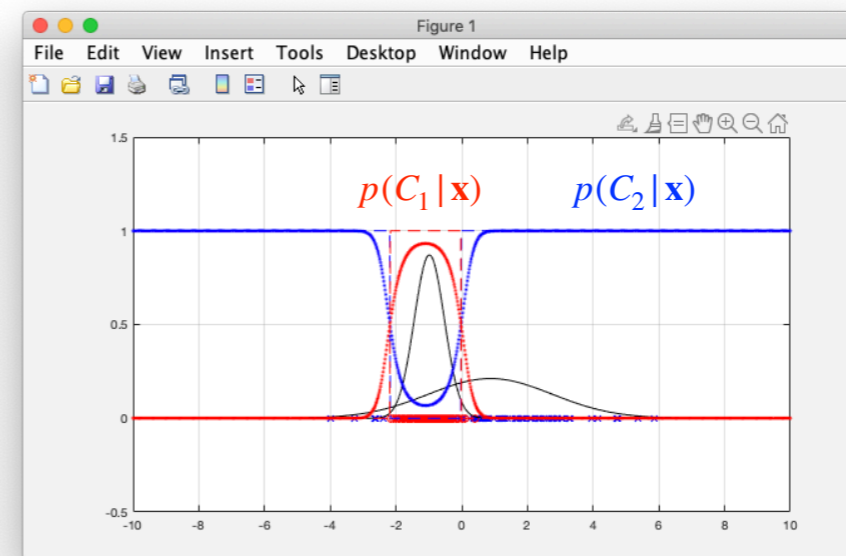
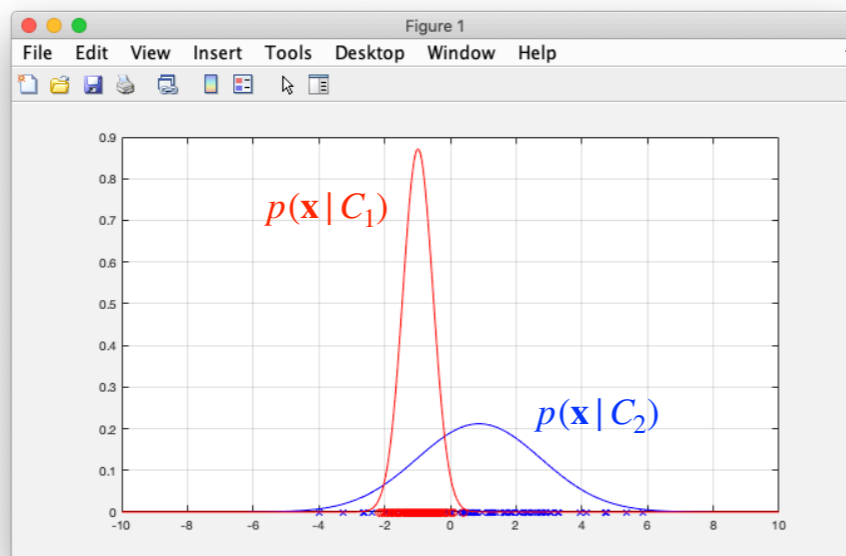
```
% PLOT RESULTS
figure;

% training data
plot(p1,zeros(n1,1),'xb',p2,zeros(n2,1),'or')
title('training data');
axis([-10 10 -0.5 1.5]);
hold on;

% conditioned pdf of each class p(x|Ci)
plot(x,pXC1,'b');plot(x,pXC2,'r');
title('conditioned pdf of each class p(x|Ci)');

% posterior pdf's p(Ci|x)
plot(x,PC1x,'.b');plot(x,PC2x,'.r');
title('posterior pdfs p(Ci|x)');

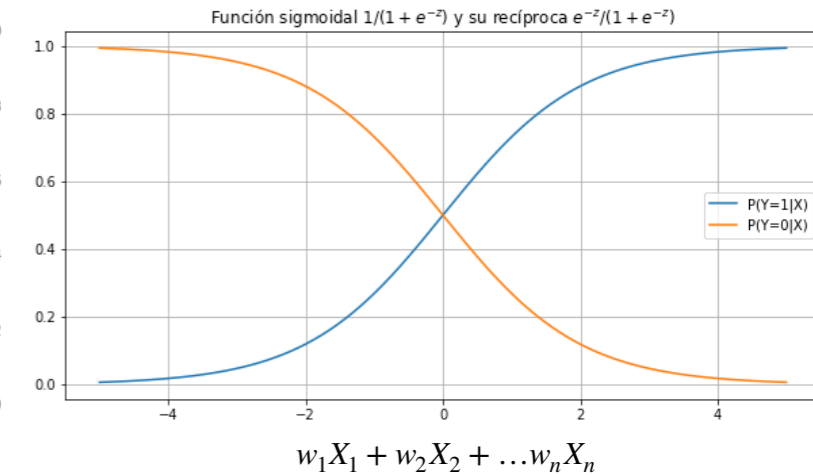
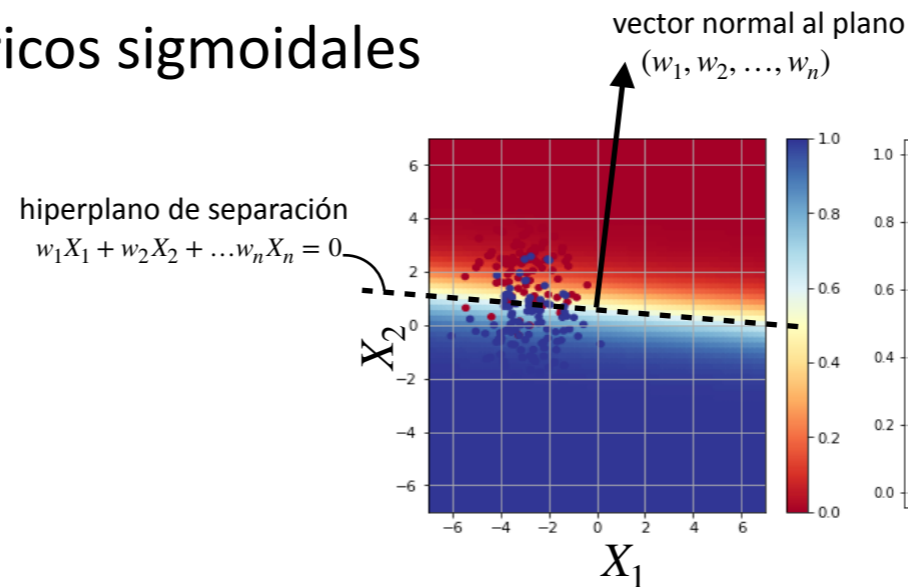
% classification: assign the class with largest posterior probability
plot(x,PC1x>PC2x,'--b');plot(x,PC2x>PC1x,'--r');
title('classification: assign the class with largest posterior probability');
```



Suponemos modelos paramétricos sigmoidales

$$P(Y = 0 | X) = \frac{e^{-\sum_i w_i X_i}}{1 + e^{-\sum_i w_i X_i}}$$

$$P(Y = 1 | X) = \frac{1}{1 + e^{-\sum_i w_i X_i}}$$

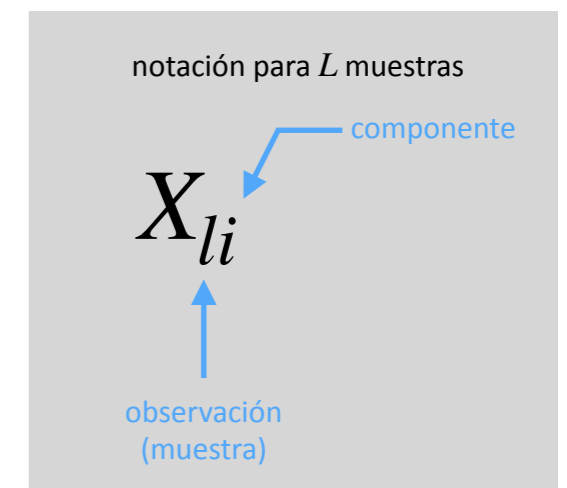


La regla de clasificación sería

$$1 < \frac{P(Y = 0 | X)}{P(Y = 1 | X)} \quad \Leftrightarrow \quad 0 < - \sum_i w_i X_i$$

Si tenemos  $L$  observaciones,  $X_1, \dots, X_L$ ,  
la *verosimilitud* de acuerdo con el modelo  
paramétrico  $P(Y | X, w)$  sería

$$\prod_l P(Y_l | X_l, w)$$



# Clasificación

## Regresión Logística

Buscamos los pesos  $w$  del modelo que maximicen la verosimilitud

$$w \leftarrow \arg \max_w \prod_l P(Y_l | X_l, w)$$

Tomando logaritmos, se convierten en sumas

$$w \leftarrow \arg \max_w \sum_l \ln P(Y_l | X_l, w) \quad \text{o bien} \quad w \leftarrow \arg \max_w L(w)$$

*log-likelihood*

$Y_l$  solo toma dos valores  $\{0, 1\}$ , podemos separar sumandos en dos:  $Y_l = 1$ ,  $Y_l = 0$

$$L(w) = \sum_l \left\{ Y_l \ln P(Y_l = 1 | X_l, w) + (1 - Y_l) \ln P(Y_l = 0 | X_l, w) \right\}$$

$$L(w) = \sum_l \left\{ Y_l \ln \frac{P(Y_l = 1 | X_l, w)}{P(Y_l = 0 | X_l, w)} + \ln P(Y_l = 0 | X_l, w) \right\}$$

$$L(w) = \sum_l \left\{ Y_l \left( \sum_i w_i X_{li} \right) + \ln \frac{e^{-\sum_i w_i X_{li}}}{1 + e^{-\sum_i w_i X_{li}}} \right\}$$



derivando respecto a los pesos  $w_i$  obtenemos el gradiente

$$\frac{\partial L(w)}{\partial w_i} = \sum_l X_{li}(Y_l - P(Y_l = 1 | X_l, w))$$

... y del gradiente la regla de aprendizaje

$$w_i \leftarrow w_i + \nu \sum_l X_{li}(Y_l - P(Y_l = 1 | X_l, w))$$

ejemplo programándolo (¡muy sencillo!)

```
# Clase para iteración de regresión logística
# ¡¡La solución es muy simple!!

class LogisticRegression():
    def fit(self, X, y, n_iter=4000, lr=0.01):
        self.w = np.random.rand(X.shape[1])
        for _ in range(n_iter):
            self.w -= lr * (self.predict(X) - y).dot(X)
    def predict(self, X):
        return sigmoid(X.dot(self.w))
```

```
# Instanciamos nuestra clase LogisticRegression
a = LogisticRegression()

# ejecutamos el método fit para entrenar
a.fit(X,y,lr=0.0001)

# Definimos datos de test (Xt)

# ejecutamos el método predict()
yt = a.predict(Xt)
```

ejemplo con *scikit-learn*

```
from sklearn.datasets import load_iris
from sklearn.linear_model import LogisticRegression

X, y = load_iris(return_X_y=True)

clf = LogisticRegression(random_state=0).fit(X, y)

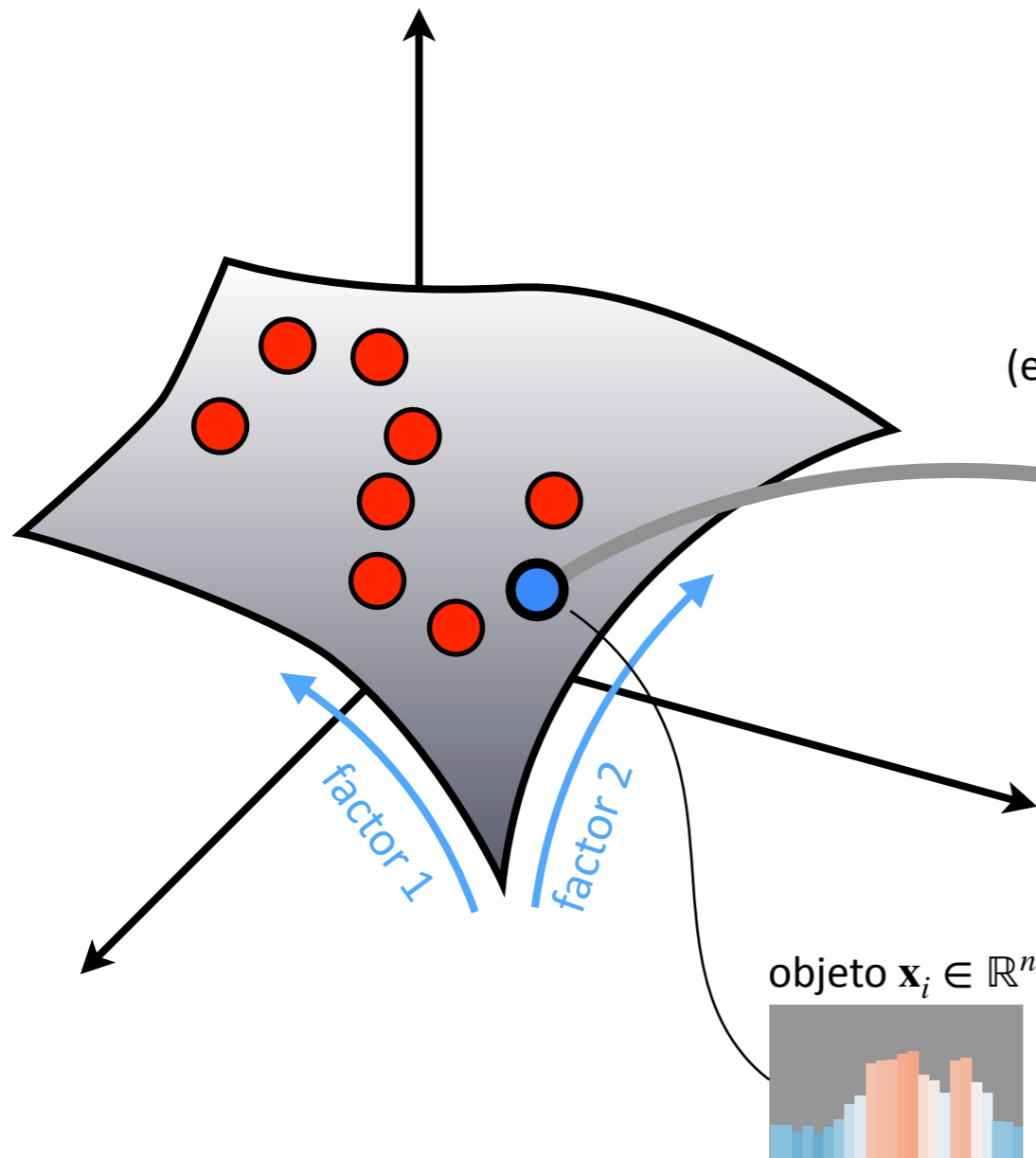
clf.predict(X[:2, :])
clf.predict_proba(X[:2, :])
clf.score(X, y)
```

# Reducción de la Dimensionalidad

# Reducción de la dimensionalidad

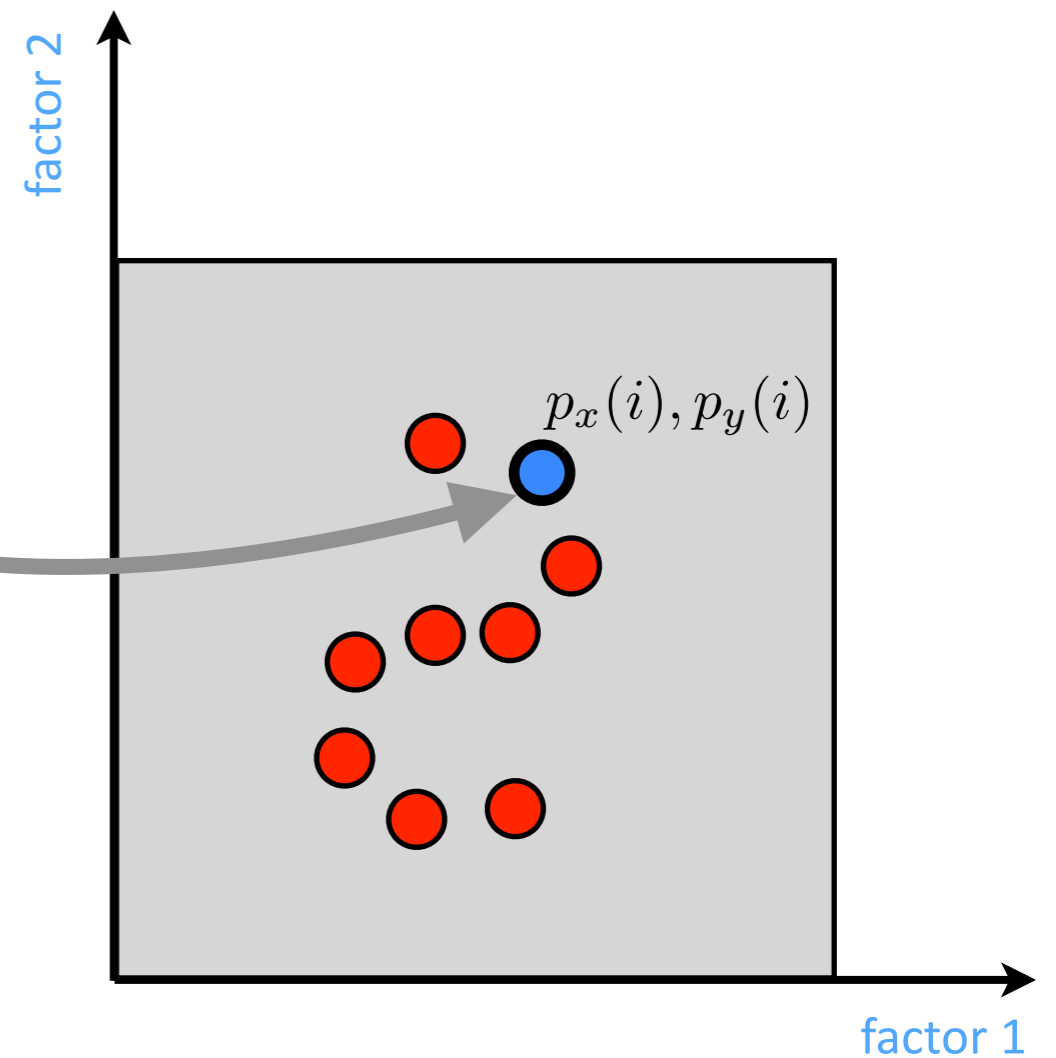
## preservación de similitudes

espacio de datos  $\mathbb{R}^n$



mapeo DR  
(ej. *t*-SNE, UMAP,  
AE, VAE, ...)

espacio latente  $\mathbb{R}^d$   
(o de visualización,  $d = 2$ )



Principio de espacialización:  
“parecido”  $\leftrightarrow$  “cercano”

# Reducción de la dimensionalidad

## ejemplo en scikit-learn

```
import matplotlib.pyplot as plt

from sklearn import datasets
from sklearn.decomposition import PCA
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis

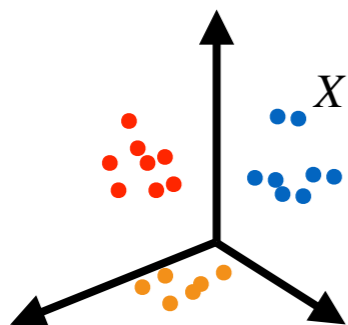
iris = datasets.load_iris()

X = iris.data
y = iris.target
target_names = iris.target_names
```

$$X \in \mathbb{R}^{n,4}$$

|     | sepal length (cm) | sepal width (cm) | petal length (cm) | petal width (cm) |
|-----|-------------------|------------------|-------------------|------------------|
| 0   | 5.1               | 3.5              | 1.4               | 0.2              |
| 1   | 4.9               | 3.0              | 1.4               | 0.2              |
| 2   | 4.7               | 3.2              | 1.3               | 0.2              |
| 3   | 4.6               | 3.1              | 1.5               | 0.2              |
| 4   | 5.0               | 3.6              | 1.4               | 0.2              |
| ... | ...               | ...              | ...               | ...              |
| 145 | 6.7               | 3.0              | 5.2               | 2.3              |
| 146 | 6.3               | 2.5              | 5.0               | 1.9              |
| 147 | 6.5               | 3.0              | 5.2               | 2.0              |
| 148 | 6.2               | 3.4              | 5.4               | 2.3              |
| 149 | 5.9               | 3.0              | 5.1               | 1.8              |

[150 rows x 4 columns]



### proyecciones lineales

#### PCA (Principal Component Analysis)

```
pca = PCA(n_components=2)
X_r = pca.fit(X).transform(X)
```

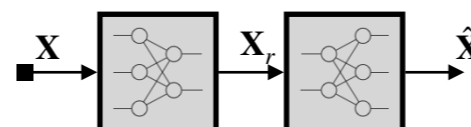
### manifold learning

UMAP, t-SNE, ...

```
from sklearn.manifold import TSNE
tsne = TSNE(n_components=2, perplexity=20)
X_r = tsne.fit_transform(X)
```

### autoencoders

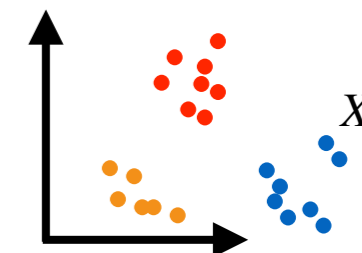
AE, deep AE, ...



$$X_r \in \mathbb{R}^{n,2}$$

|     | componente 1 | componente 2 |
|-----|--------------|--------------|
| 0   | -2.684126    | 0.319397     |
| 1   | -2.714142    | -0.177001    |
| 2   | -2.888991    | -0.144949    |
| 3   | -2.745343    | -0.318299    |
| 4   | -2.728717    | 0.326755     |
| ... | ...          | ...          |
| 145 | 1.944110     | 0.187532     |
| 146 | 1.527167     | -0.375317    |
| 147 | 1.764346     | 0.078859     |
| 148 | 1.900942     | 0.116628     |
| 149 | 1.390189     | -0.282661    |

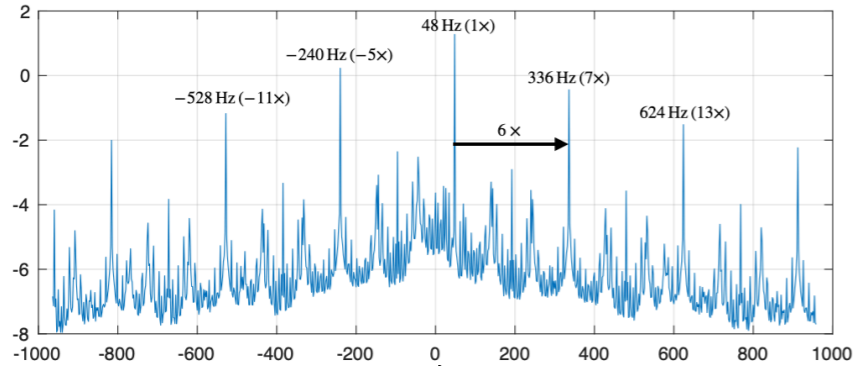
[150 rows x 2 columns]



# Reducción de la dimensionalidad

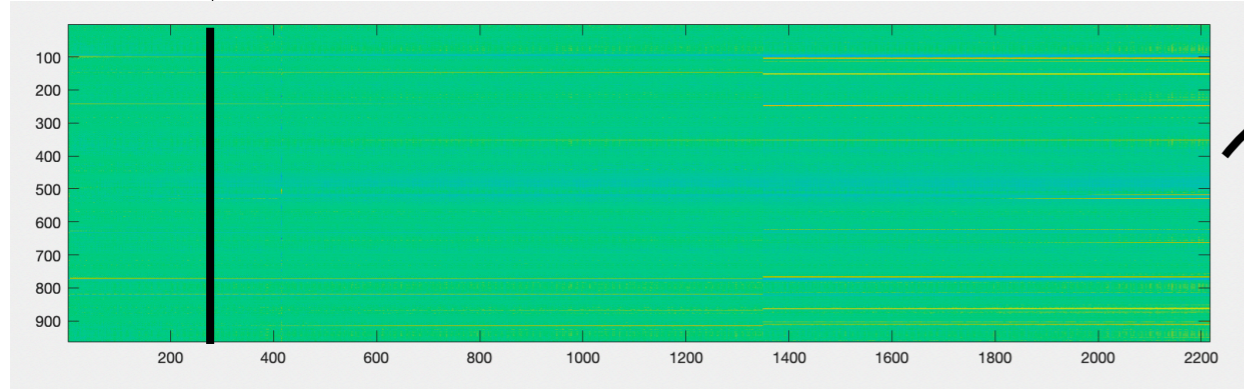
## ejemplo estados motor de inducción

espectro de corriente motor



x (1001 armonicos)

espectrograma de corrientes



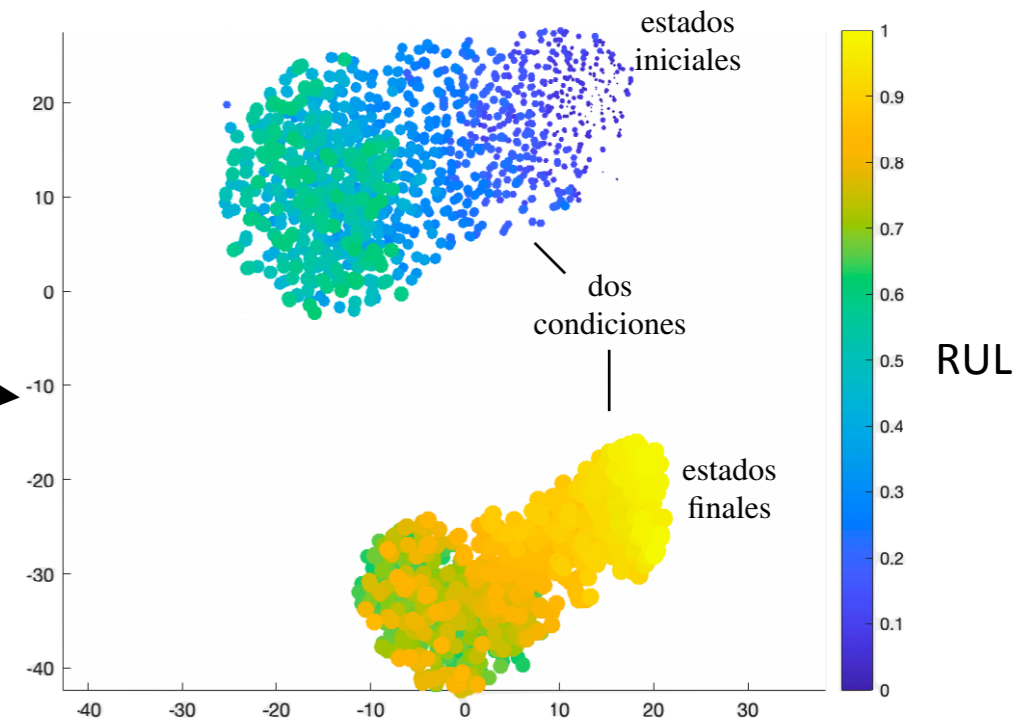
$x_i$

tiempo

tsne

$\mathbb{R}^{1001} \rightarrow \mathbb{R}^2$

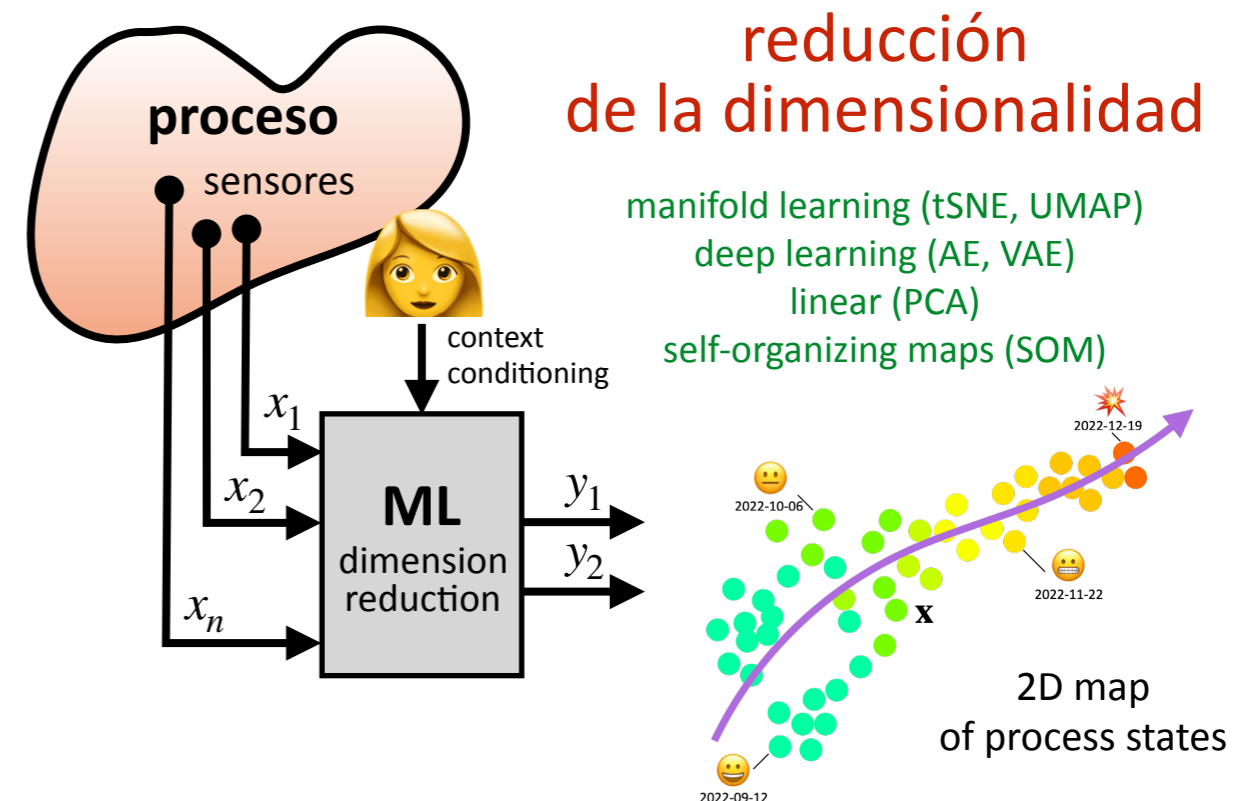
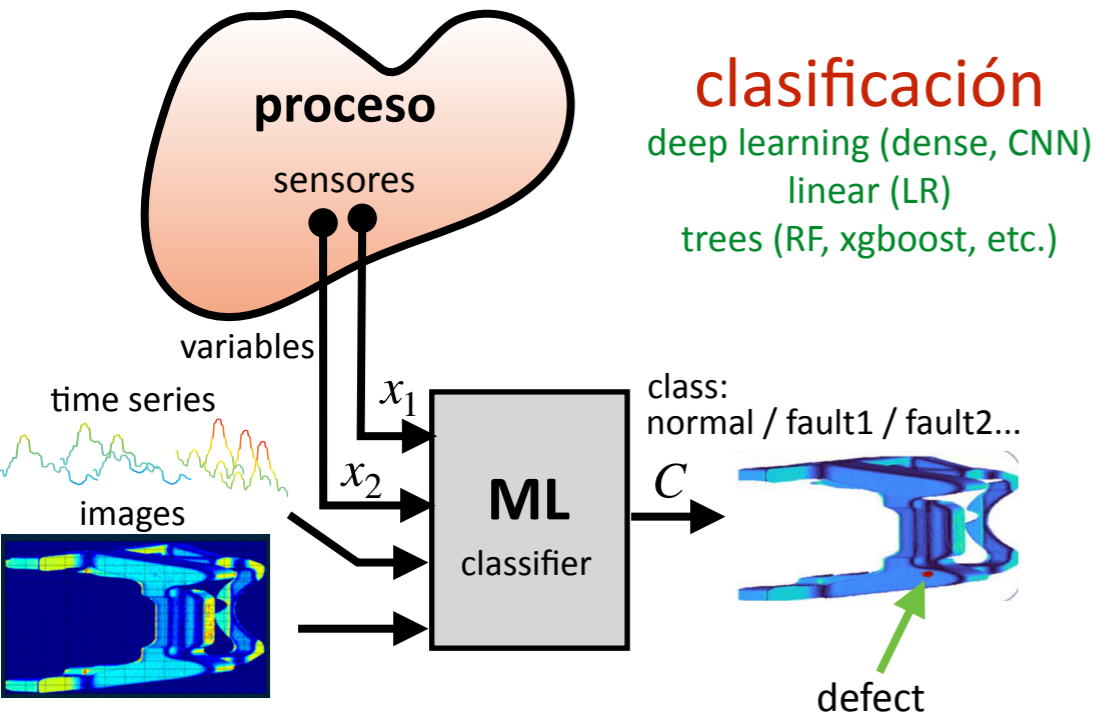
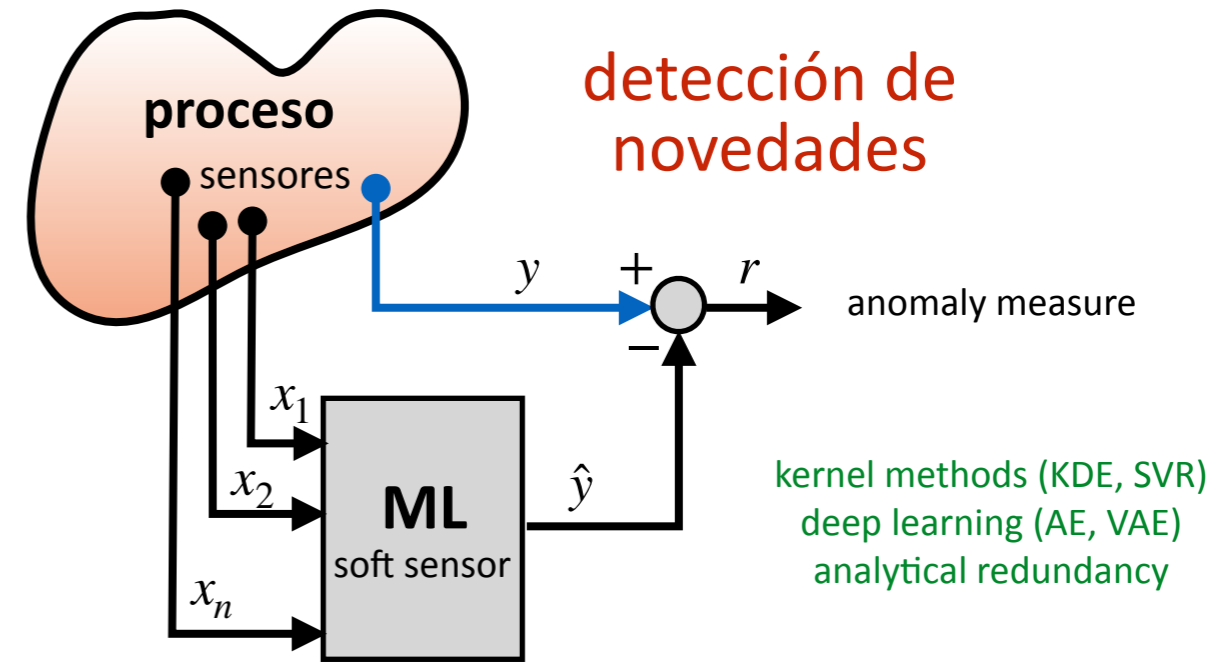
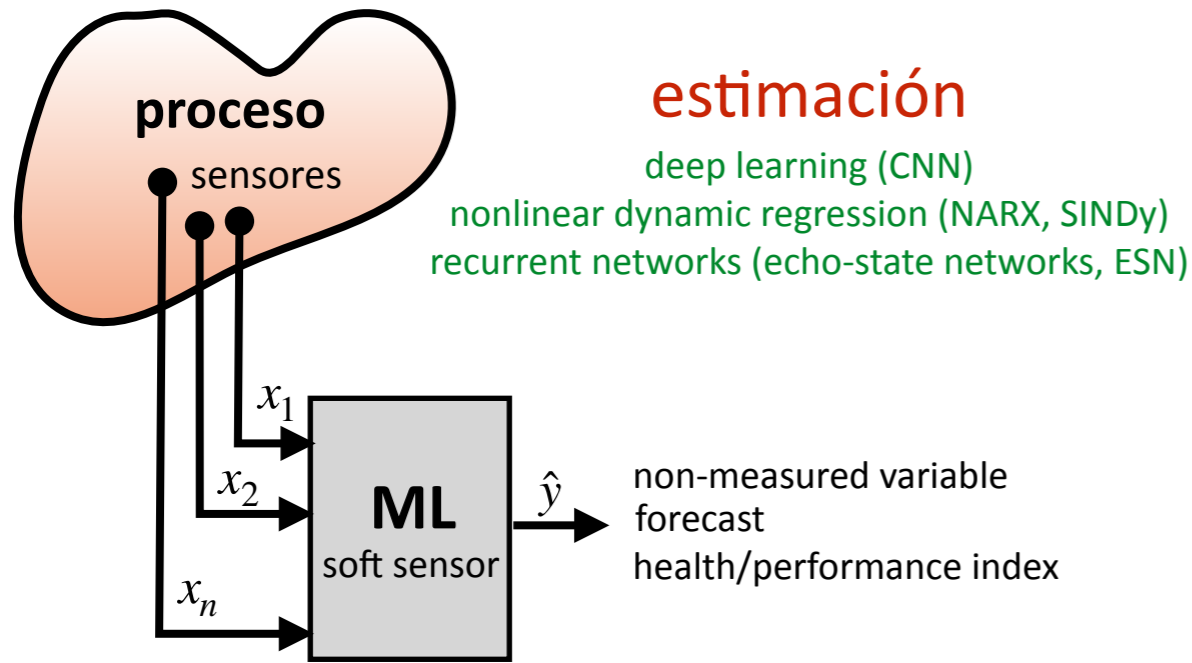
mapa visual de estados



# Aplicaciones

# Machine Learning

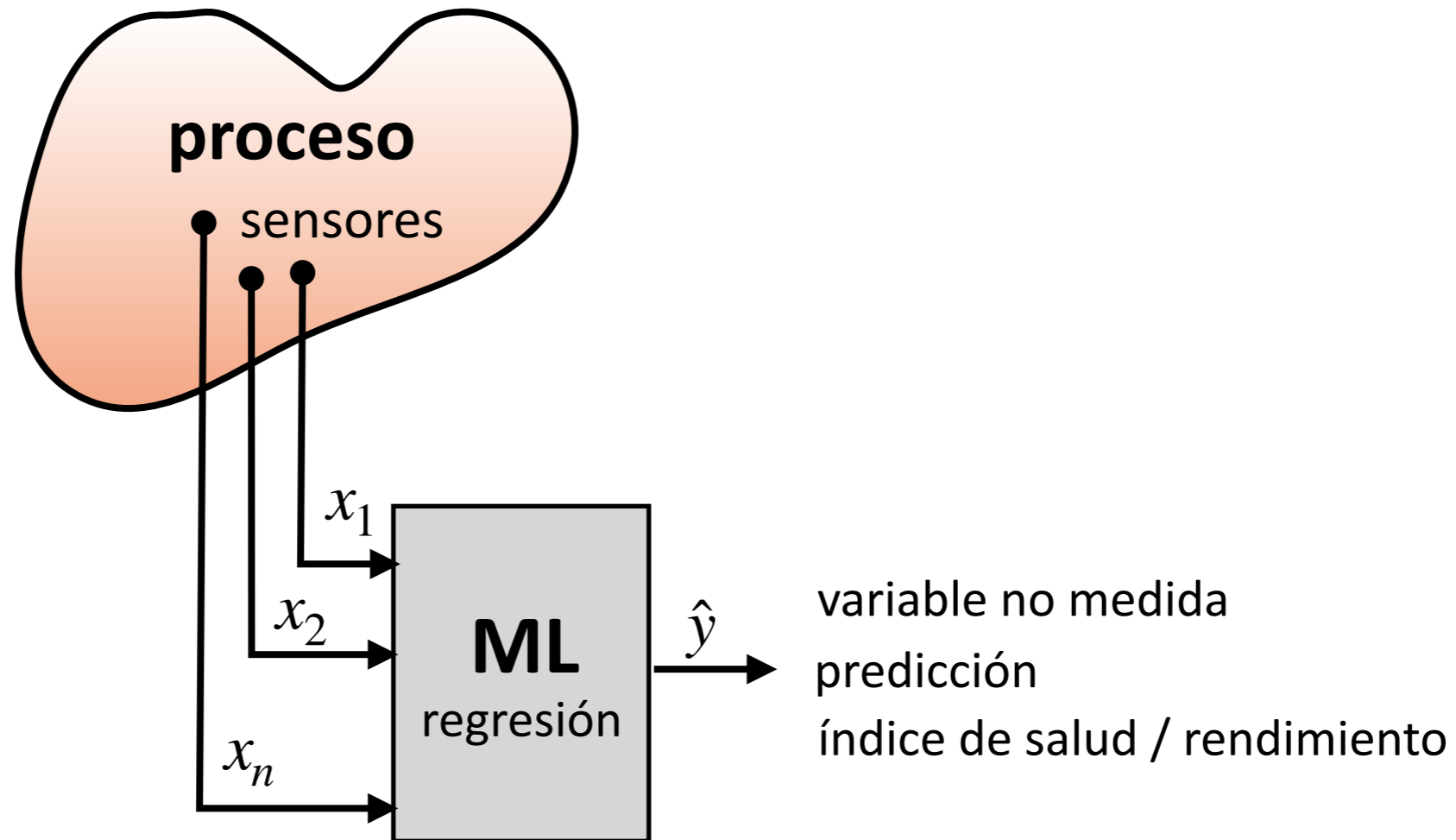
## tareas más frecuentes





# Sensores virtuales (soft sensors)

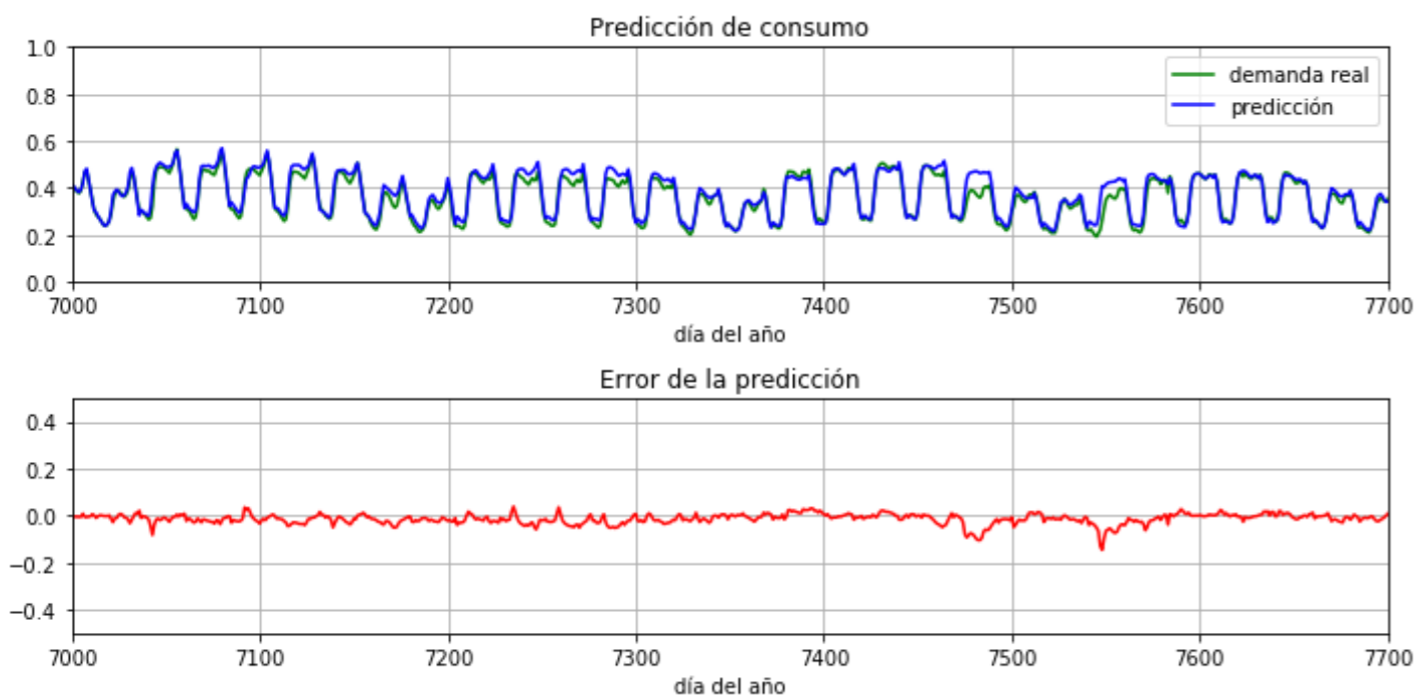
## Estimación - Predicción - Índices de rendimiento



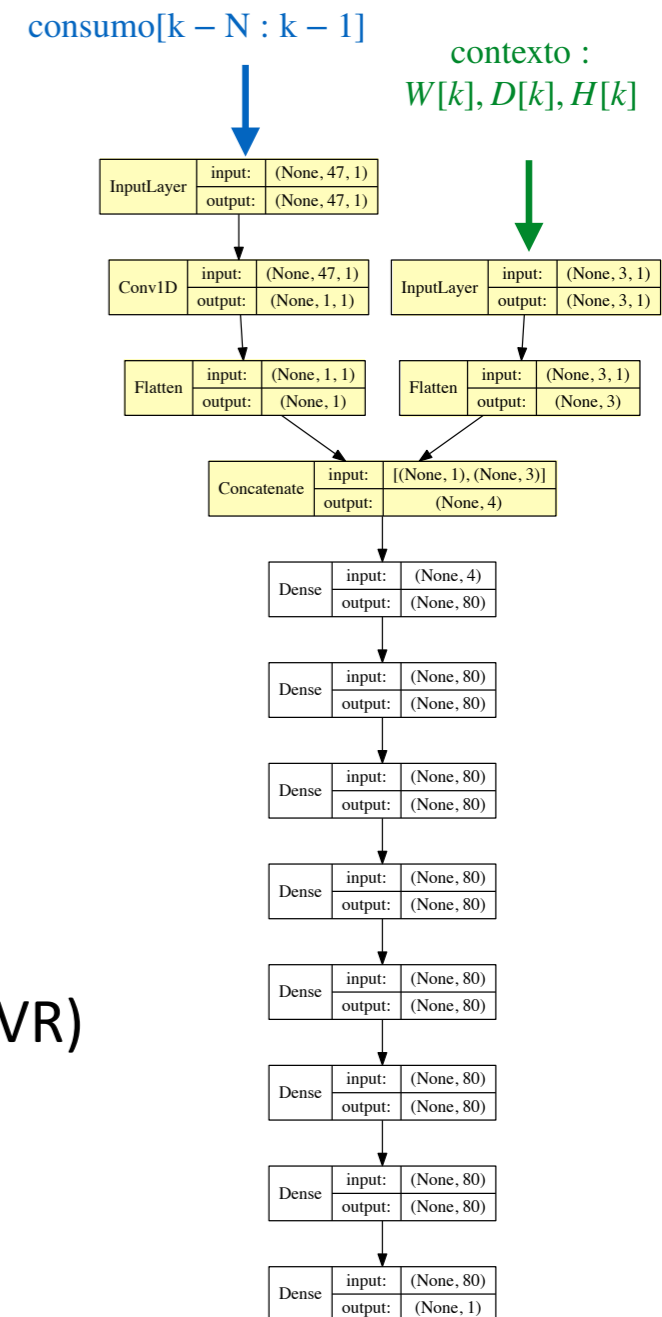


## Predicción

### Predicción de demanda eléctrica



### 3. Red Convolutiva (CNN)



(González, 2017)

### Extracción de características

```
# CREAMOS REGRESOR DINÁMICO (TAP-DELAY) DE DATOS DE ENTRENAMIENTO
X = []
Y = []

# entrenamos con los datos de un año (8760 registros horarios)
N = 8760

# estimamos los datos de varias horas después (ej. ahead = 24 horas)
ahead = 24
for k in range(169, N):
    X.append(np.hstack((
        C[k-24:k-1], # vector de consumo diario (24h)
        W[k], # semana del año
        D[k], # día de la semana
        H[k], # hora del día
        1 # término independiente
    )))
    Y.append(C[k+ahead].tolist())

# CREAMOS REGRESOR DINÁMICO (TAP-DELAY) DE DATOS DE TEST
Xtst = []
Ytst = []
Ntst = 10000

for k in range(169, Ntst):
    Xtst.append(np.hstack((
        C[k-24:k-1], # vector de consumo diario (24h)
        W[k], # semana del año
        D[k], # día de la semana
        H[k], # hora del día
        1 # término independiente
    )))
    Ytst.append(C[k+ahead].tolist())
```

### 1. Regresión Lineal (Ridge)

```
# PREDICCIÓN LINEAL
from sklearn import linear_model
clf = linear_model.Ridge(alpha = 0.1)
clf.fit(X_norm, y_norm)
yest_norm = clf.predict(Xtst_norm)
```

### 2. Regresión Vectores de Soporte (SVR)

```
# PREDICCIÓN Support Vector Regression (SVR)
from sklearn.svm import SVR
import numpy as np

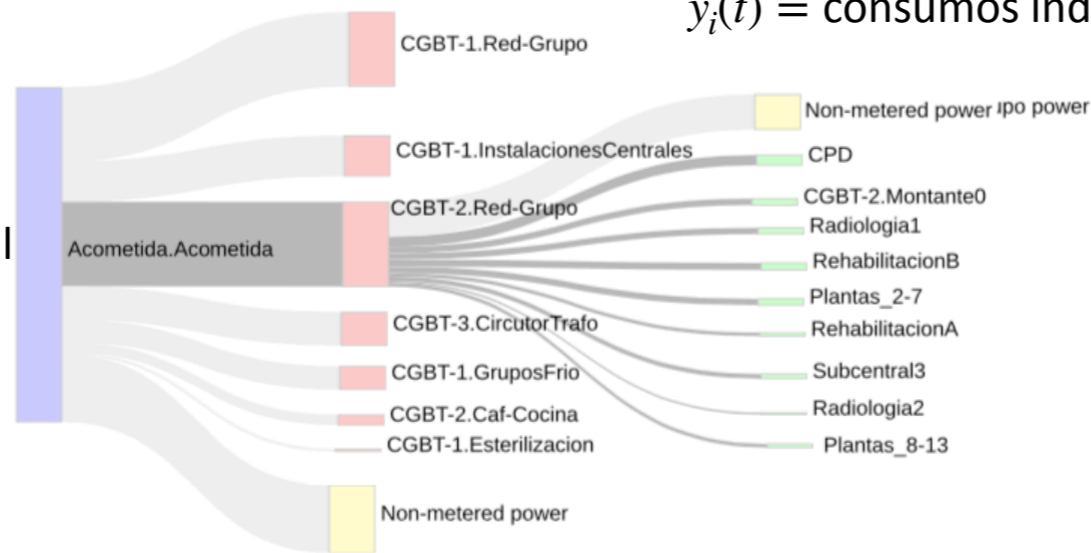
# este método permite calcular los coeficientes de una
aproximación rbf
# mediante una optimización "robusta". No está limitado a rbf. Es
extensible a otros tipos de aproximaciones.
clf = SVR(kernel='rbf', gamma=.5, C=10, epsilon=.02)

clf.fit(X_norm, y_norm)
yest_norm = clf.predict(Xtst_norm)
```

# Sensores virtuales (soft sensors)

Estimación → NILM

(García et al, 2021)

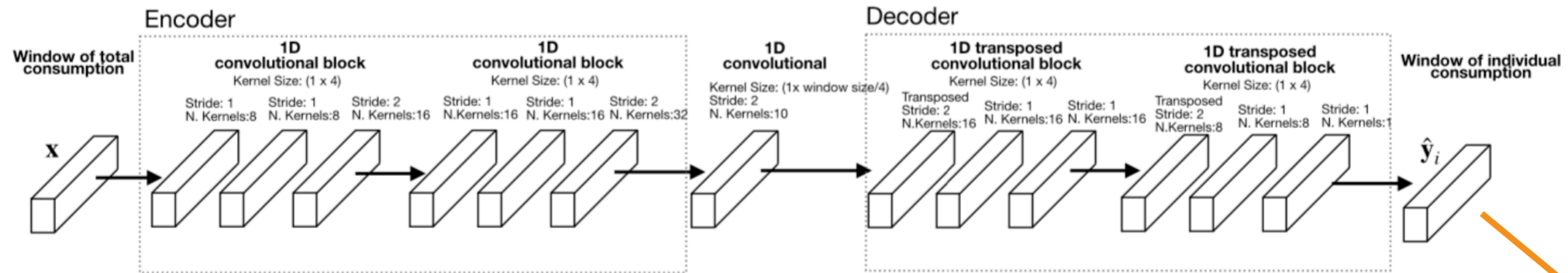


$y_i(t) = \text{consumos individuales}$

$$x(t) = y_1(t) + y_2(t) + \dots + y_n(t) + e(t)$$

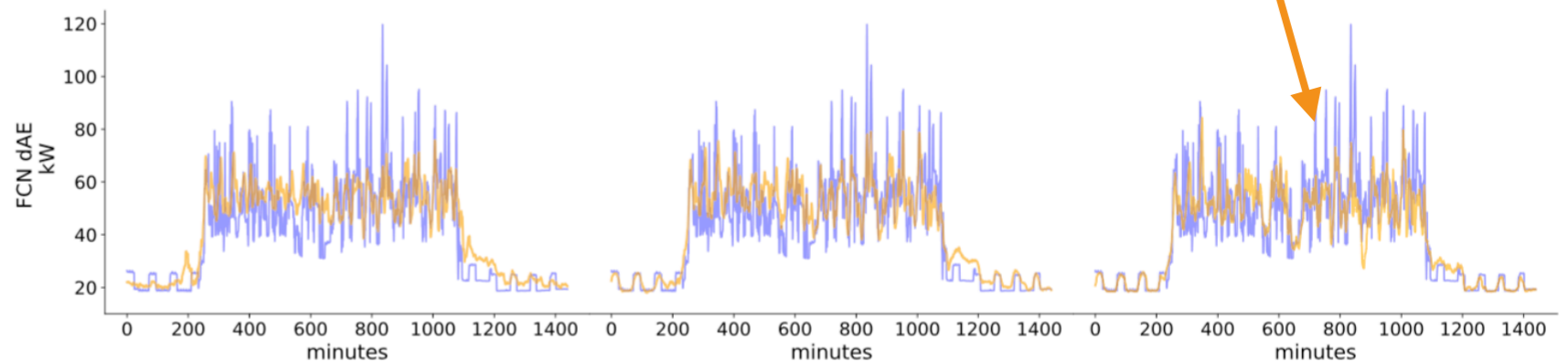
$$f_i : \mathbf{x} \rightarrow \mathbf{y}_i$$

Estimación seq2seq



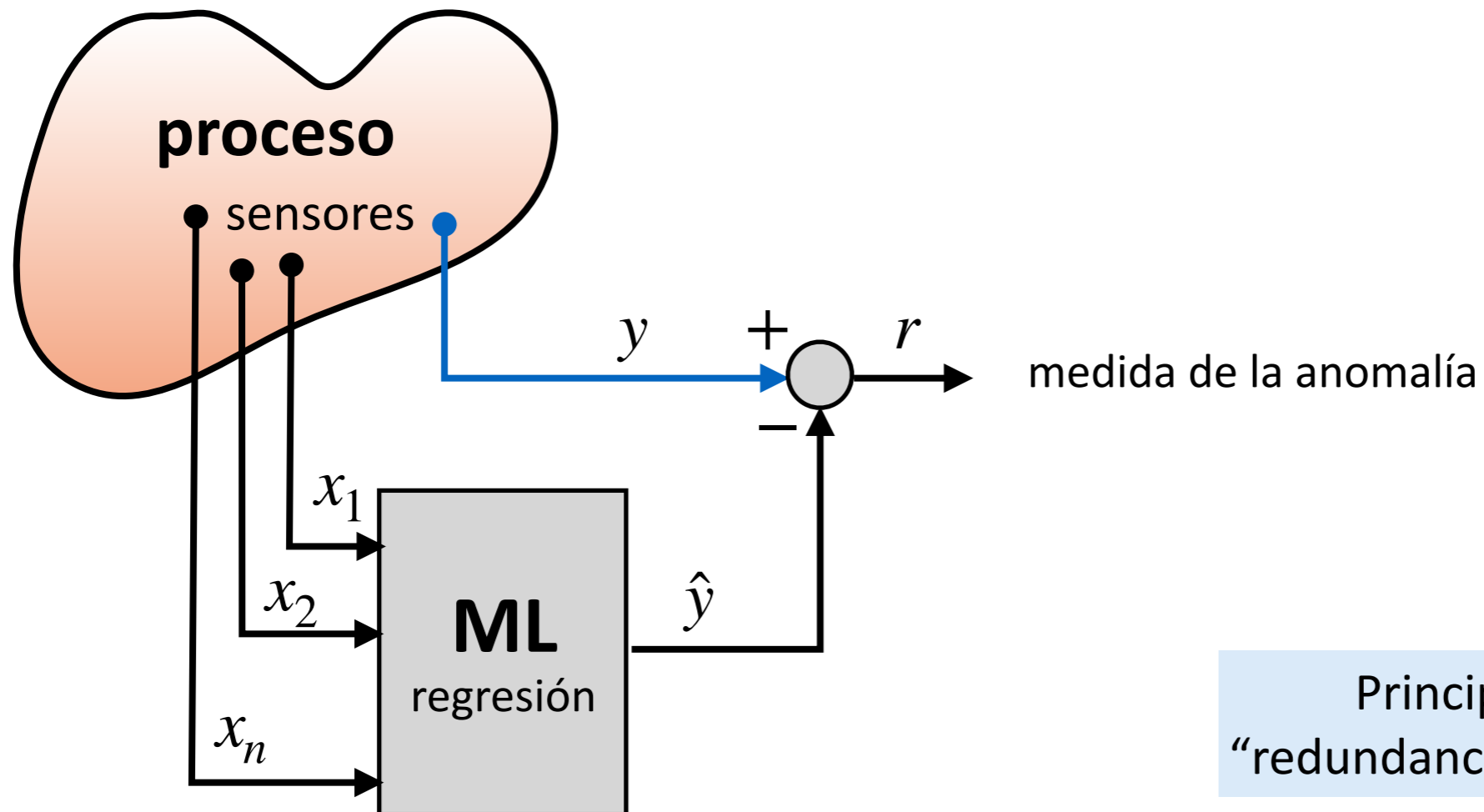
serie temporal de consumos globales

serie temporal de consumos individuales



# Detección de Novedades

## Detección de fallos / Anomalías

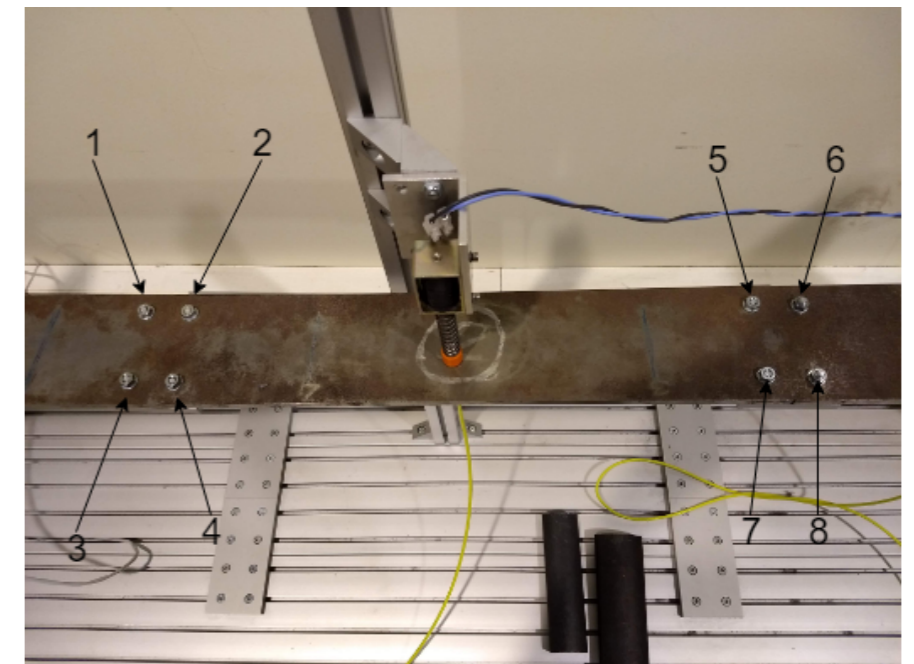
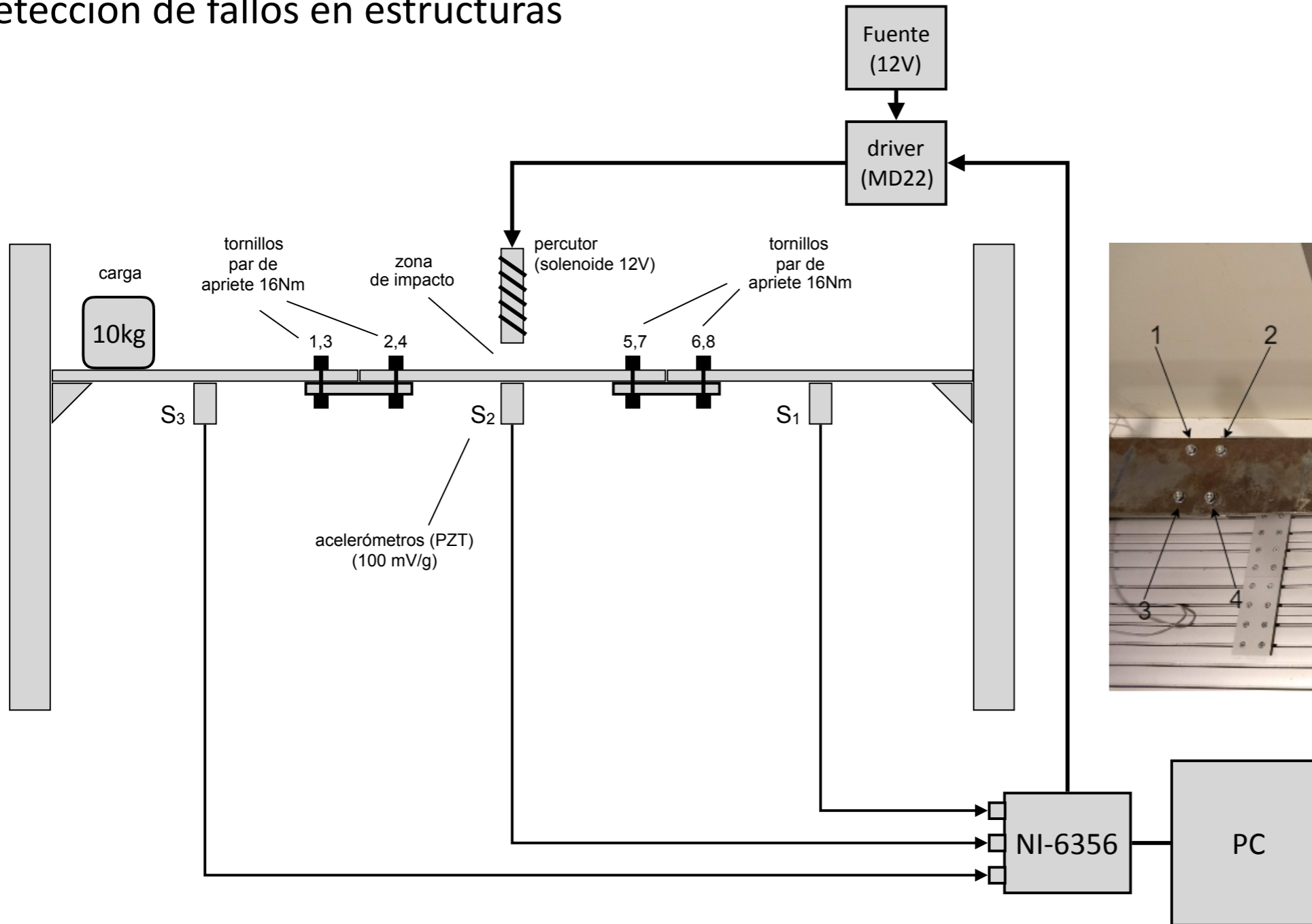


Principio de  
“redundancia analítica”

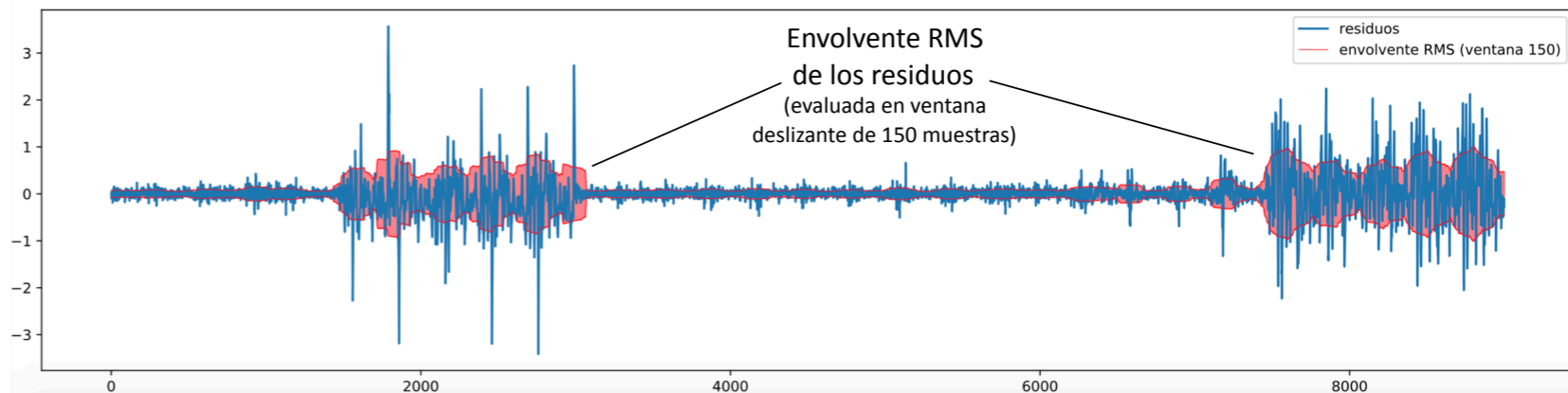
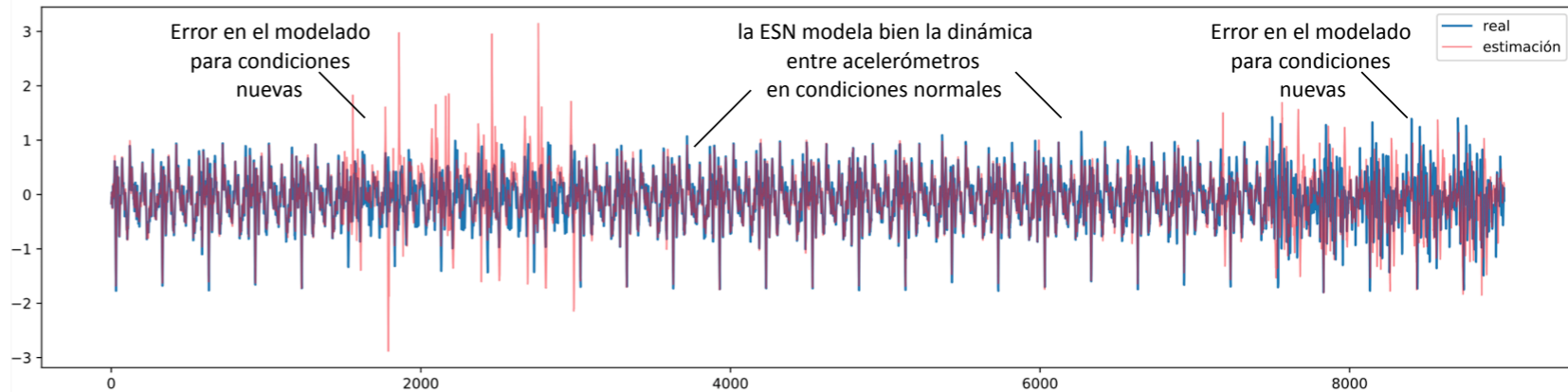
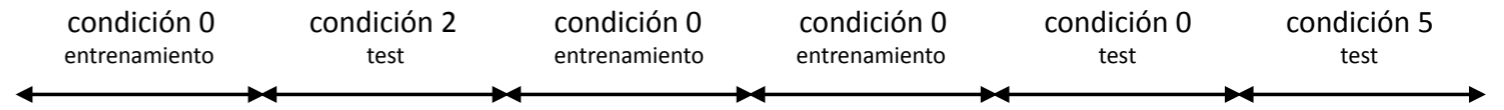
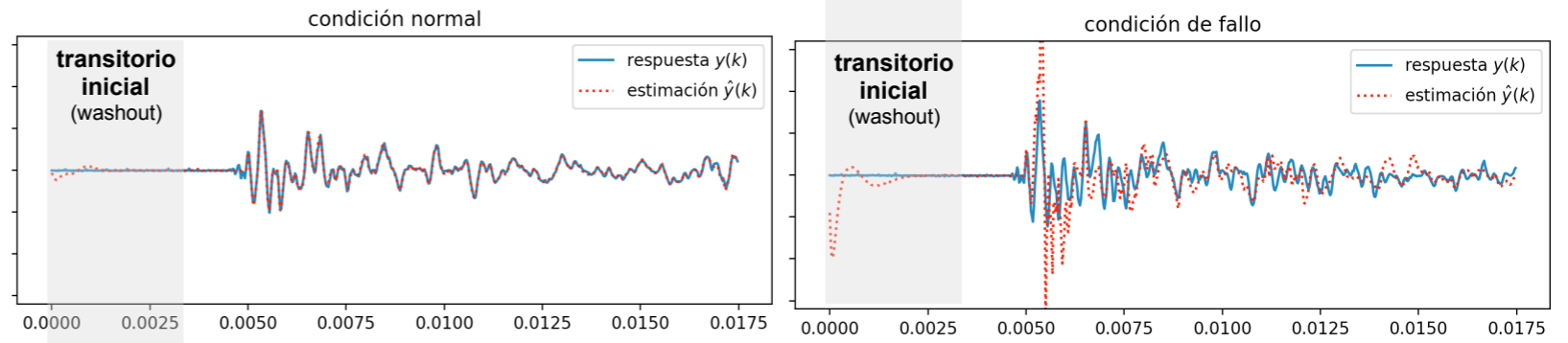
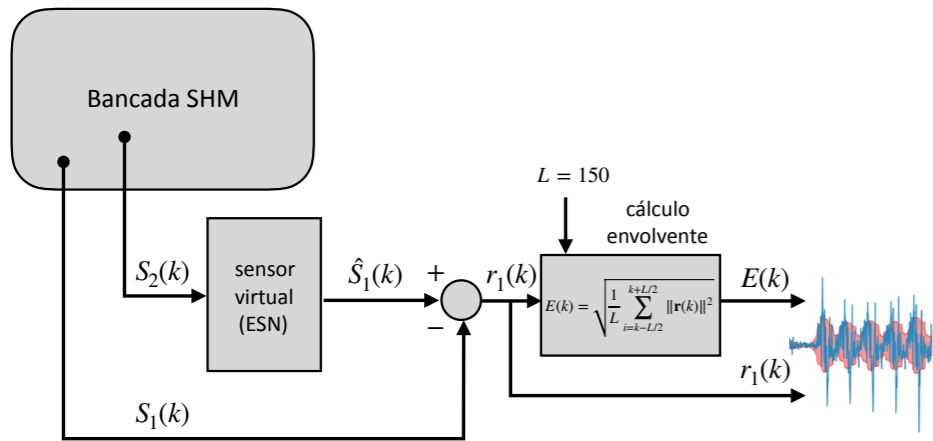
# Detección de Novedades

## Detección de fallos / Anomalías

### Detección de fallos en estructuras



### Detección de fallos en estructuras



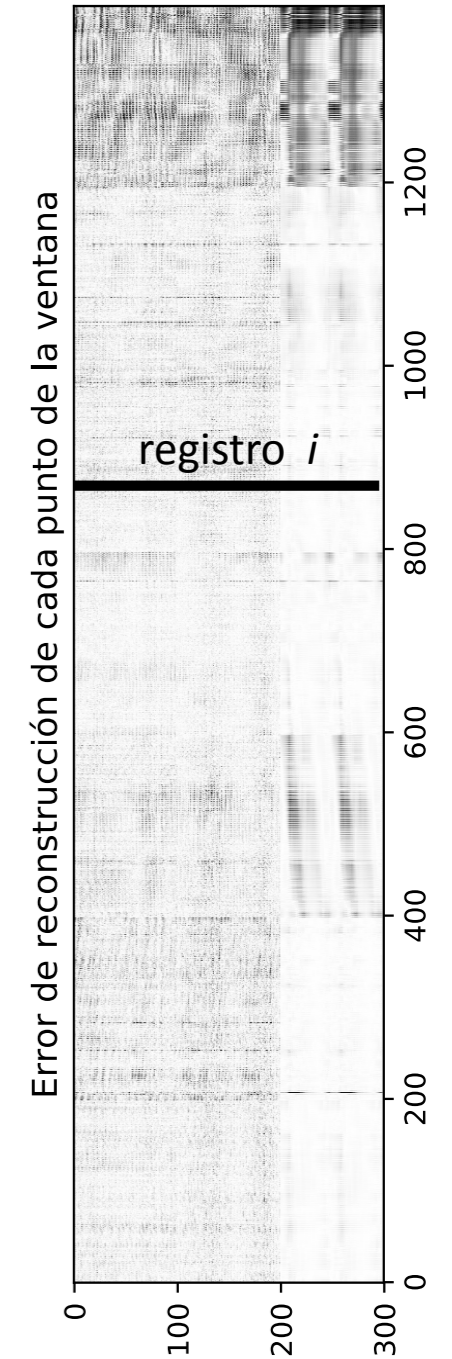
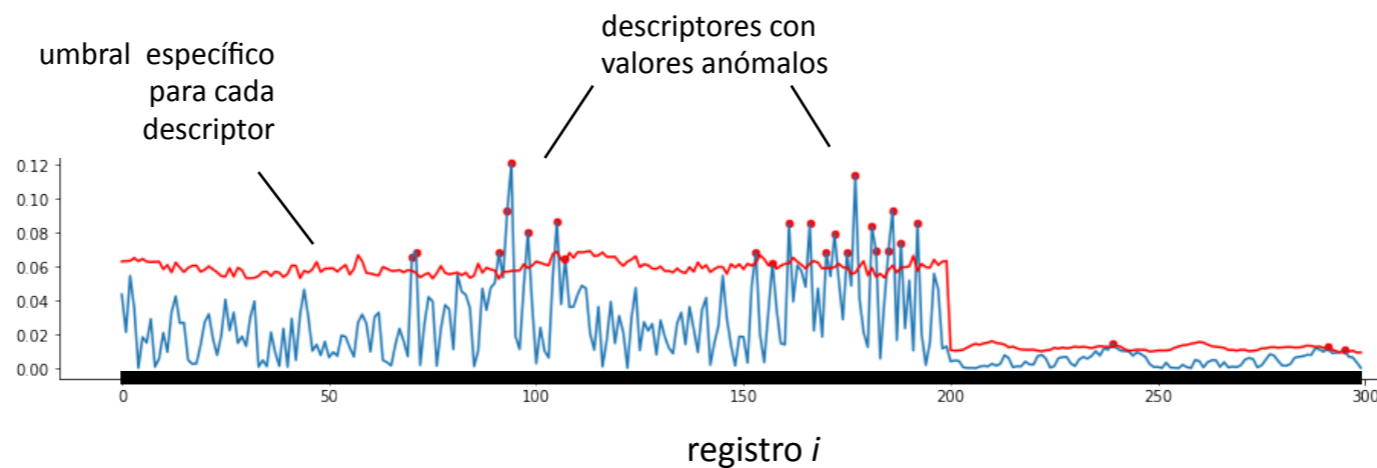
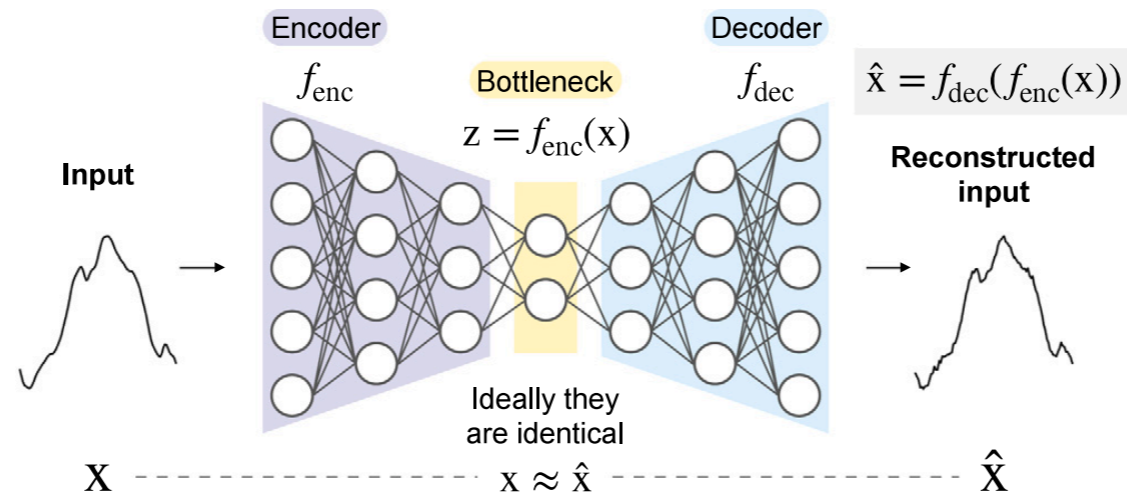


# Detección de Novedades

## Detección de fallos / Anomalías

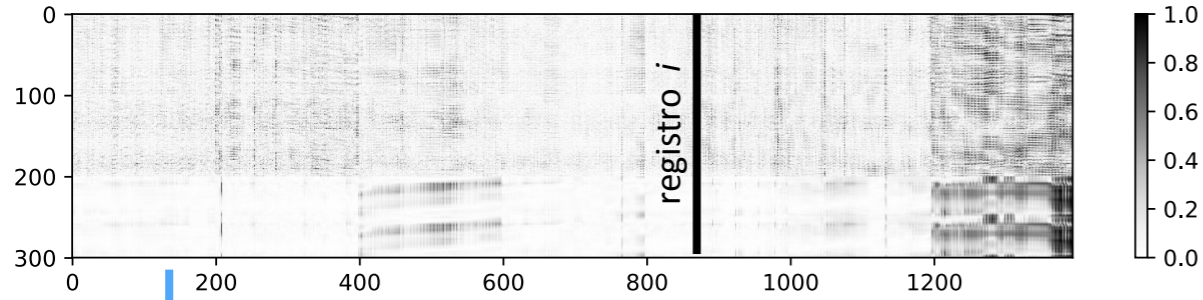
### Anomalías en máquinas eléctricas

fuelle: (González et al. 2022)

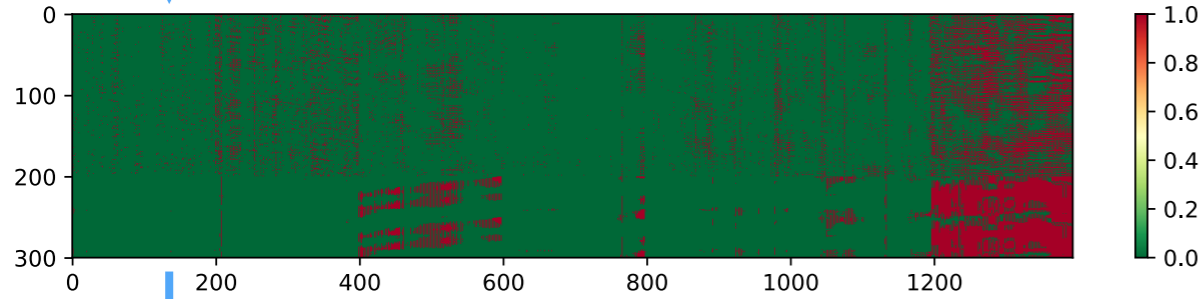


### 1600 registros con 300 descriptores

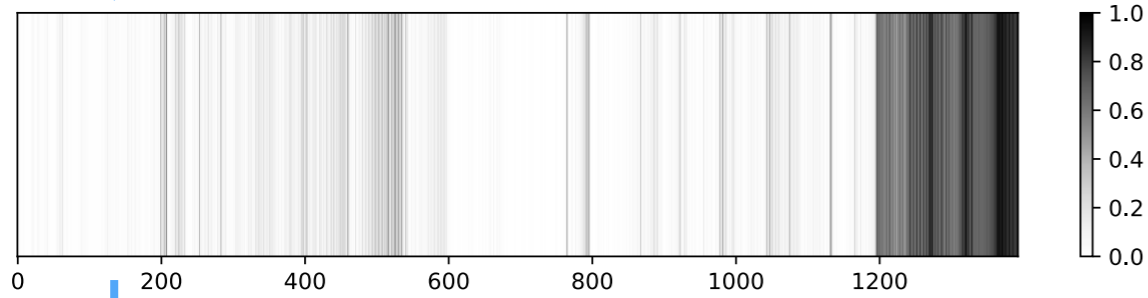
Error de reconstrucción para cada registro de descriptores



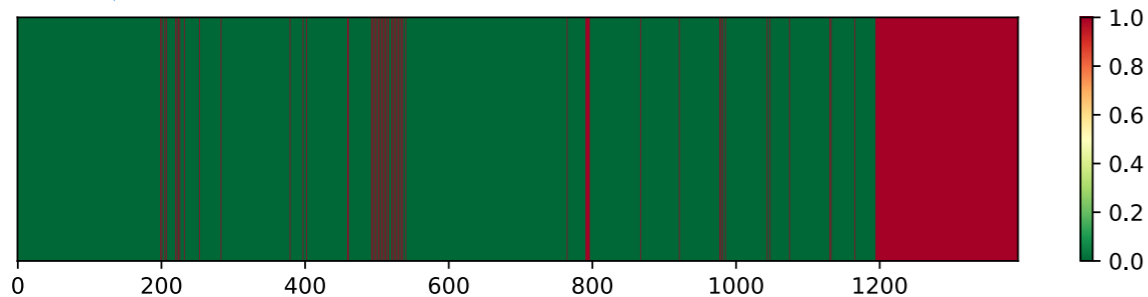
Clasificación de los descriptores de cada registro (percentil 95)



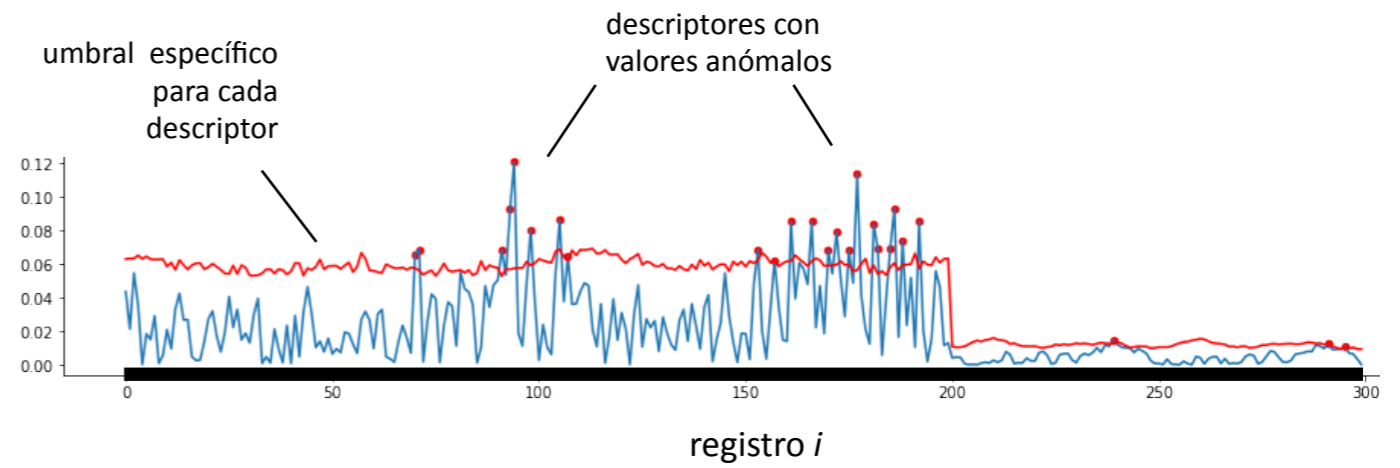
Número de descriptores anómalos en cada registro



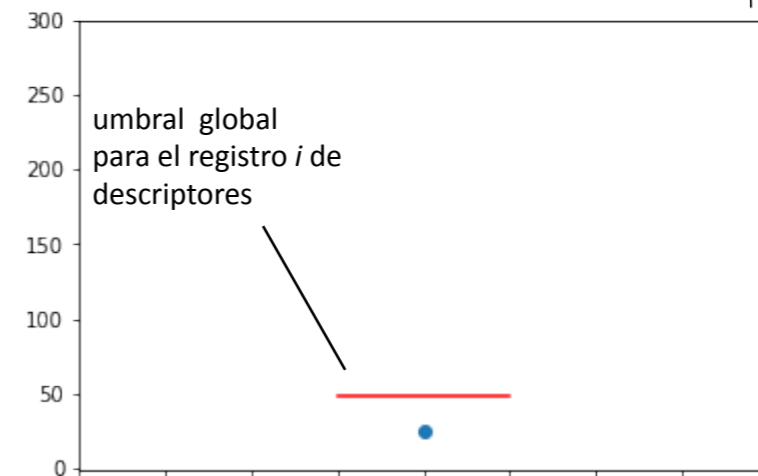
Clasificación de cada registro (perc. 95)



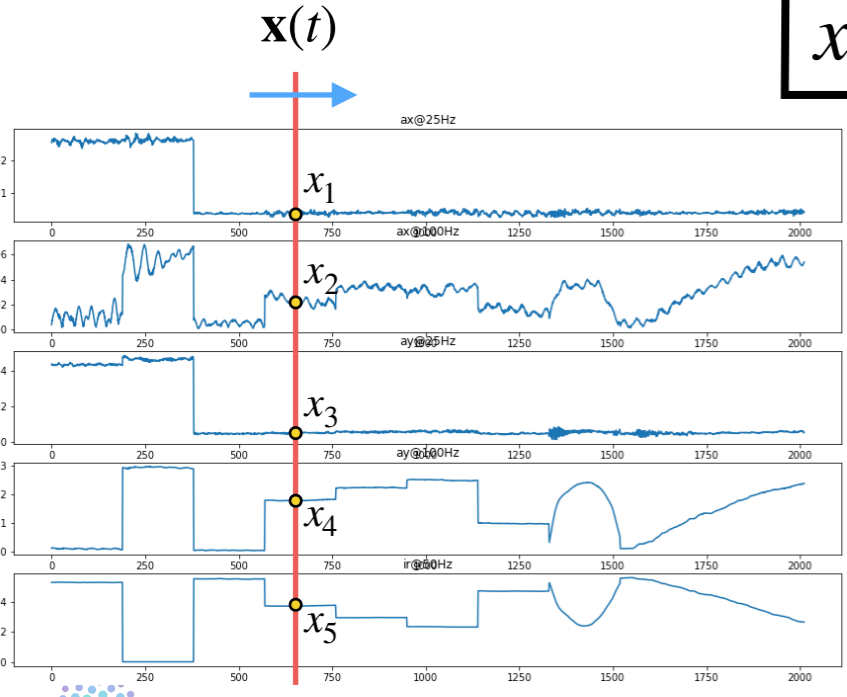
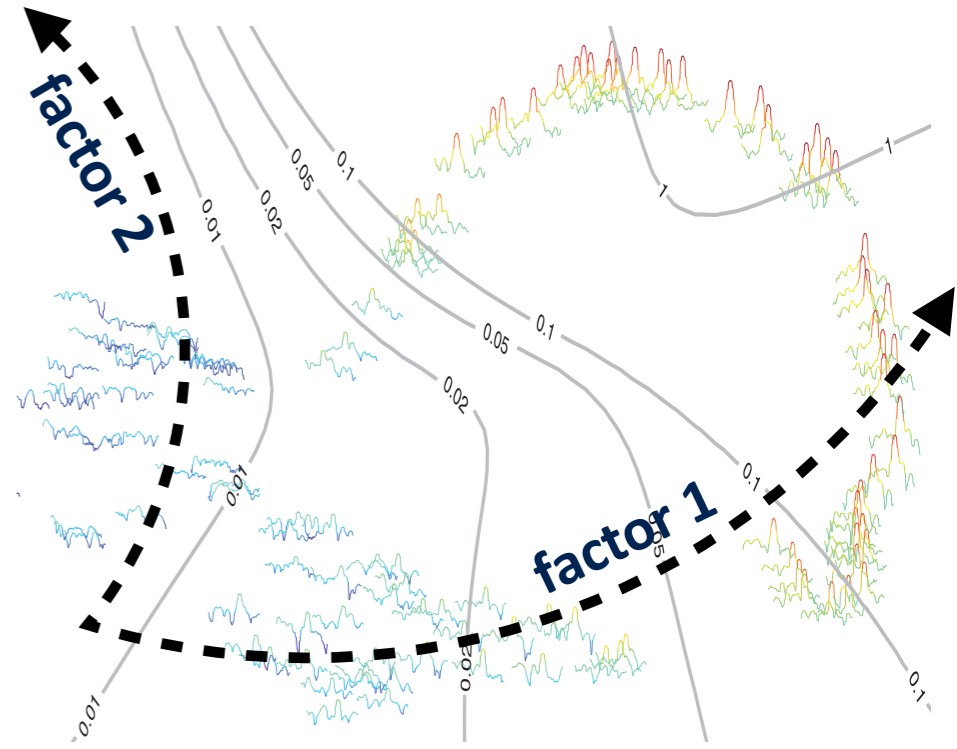
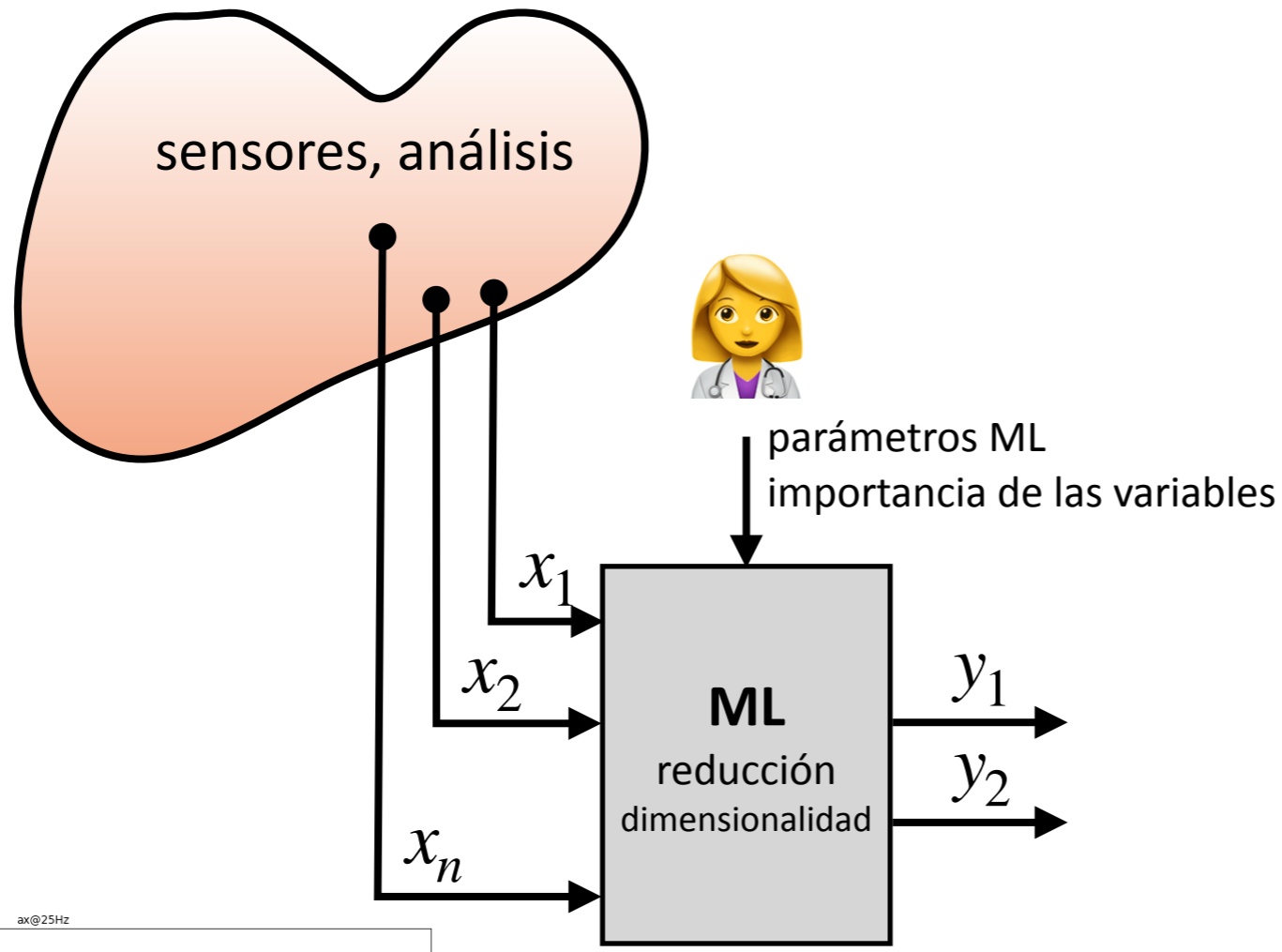
### Análisis de anomalías del registro



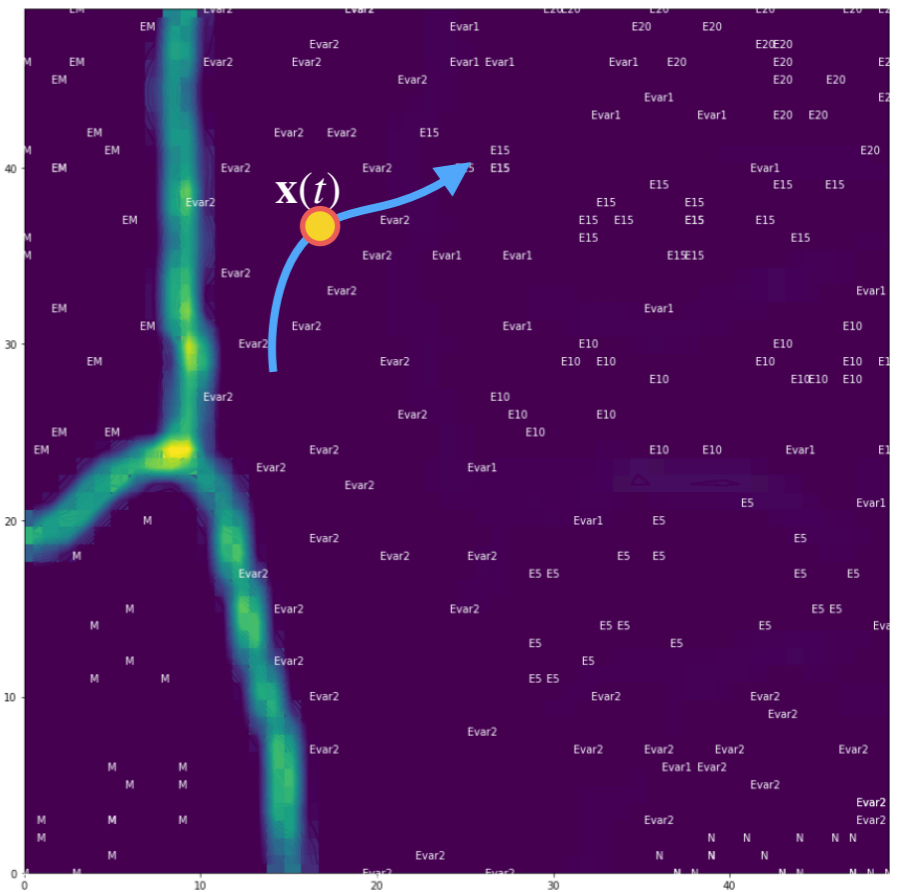
Nº de descriptores anómalos + Clasificación de la ventana



# Machine learning visualización



**Visualización**  
representación 2D/3D  
interacción, exploración,  
animación

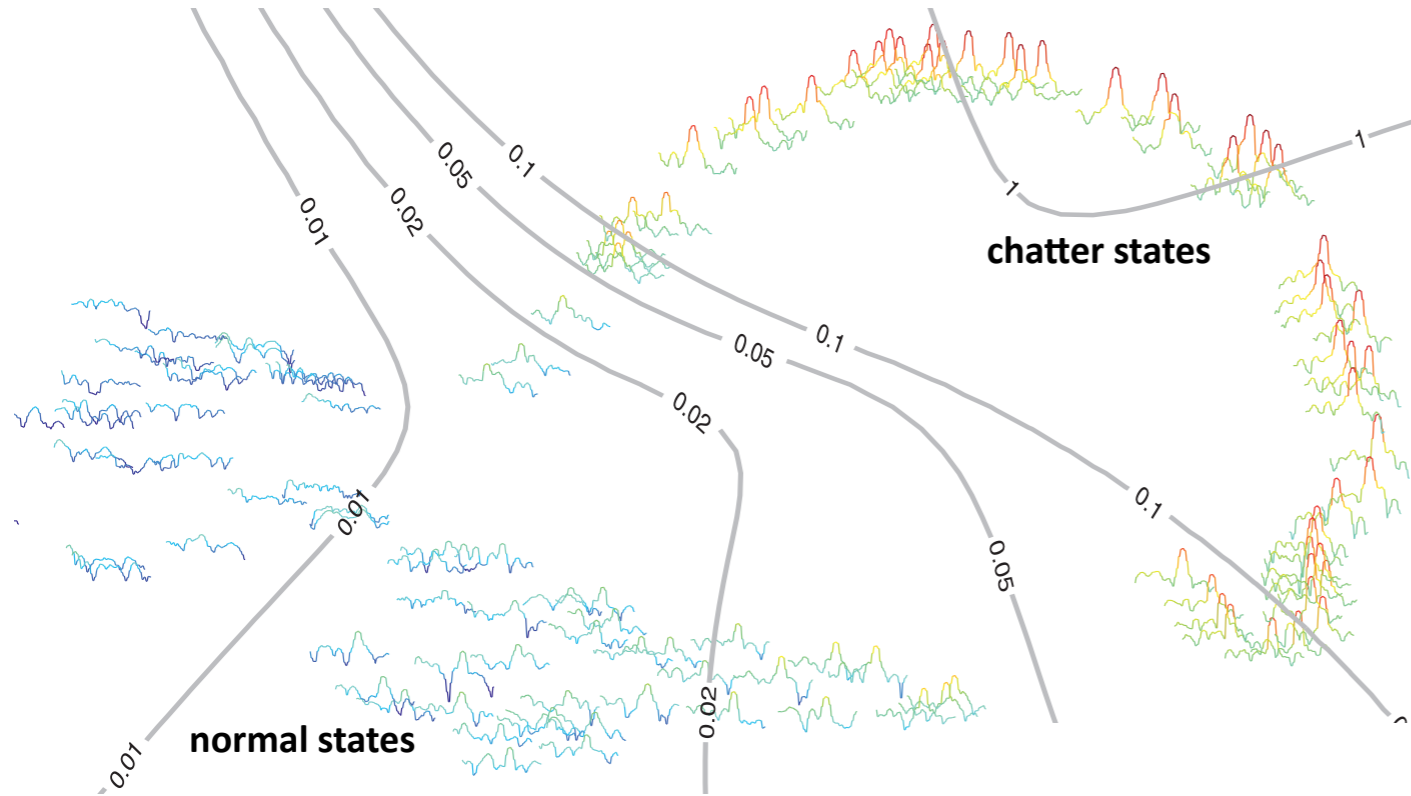




# Reducción de la dimensionalidad

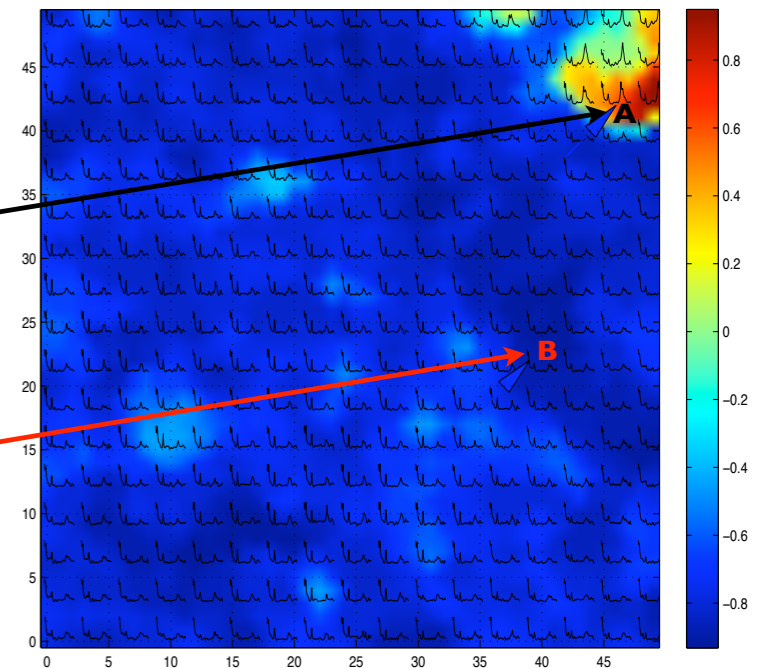
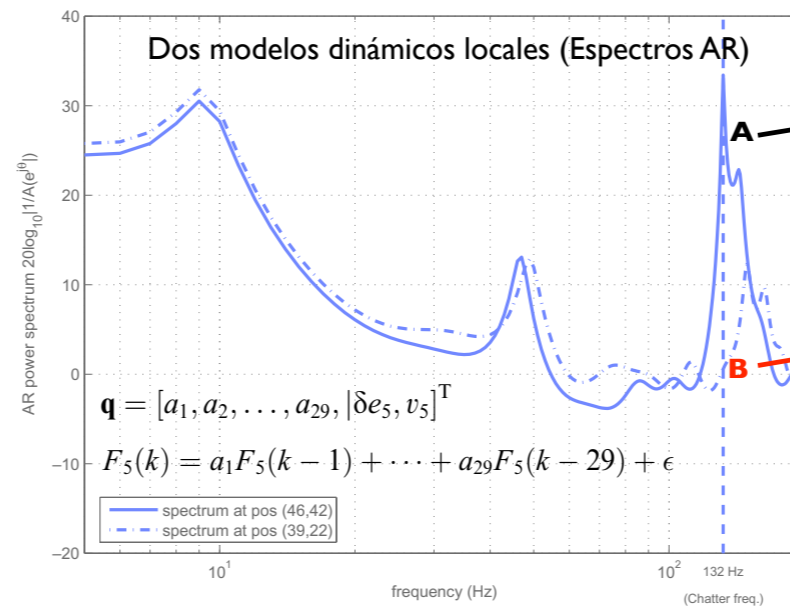
## mapas de estados del proceso

### Mapas de vibración



### Mapas de modelos: fuerzas de laminación

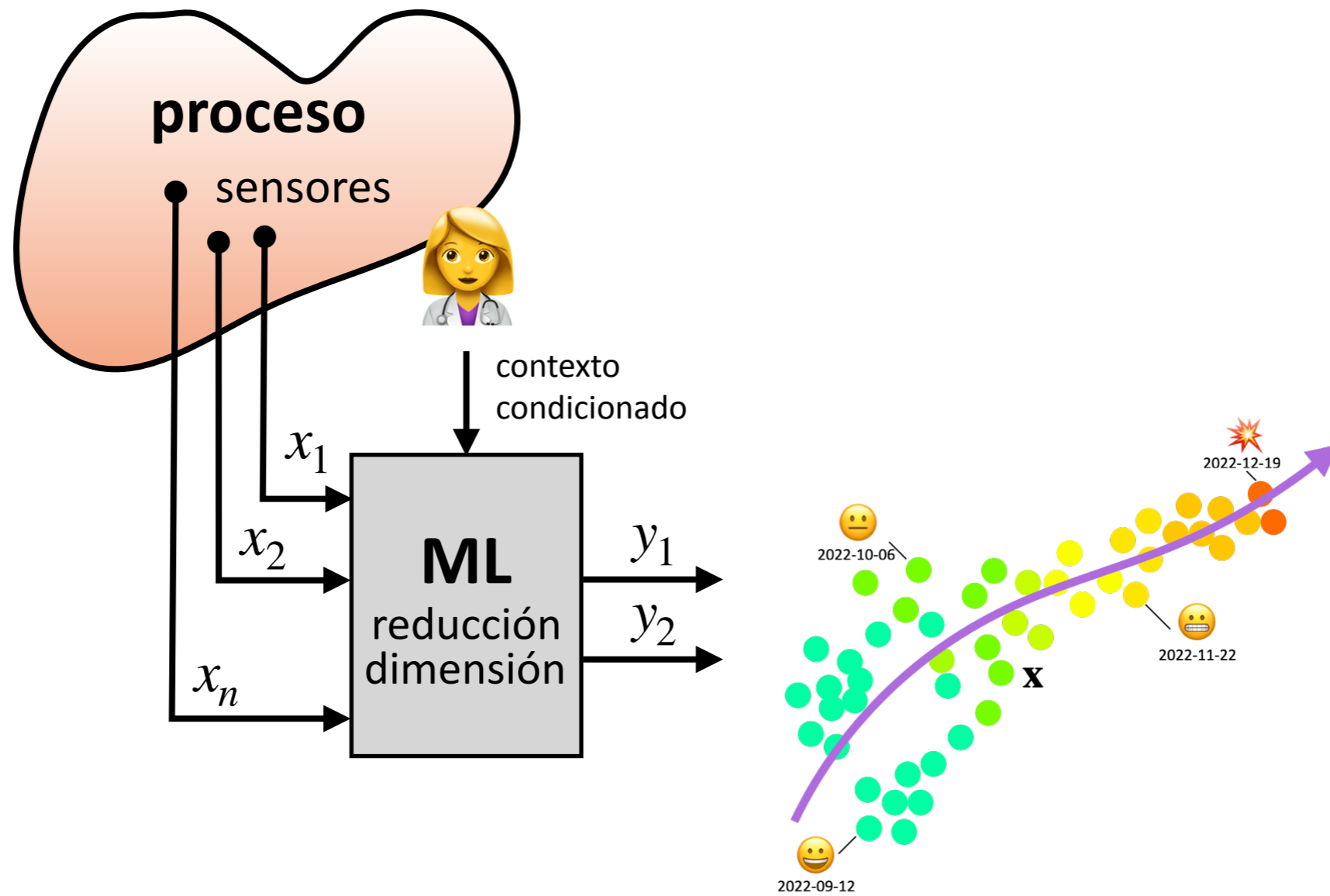
$\delta e_5$



# Reducción de la dimensionalidad

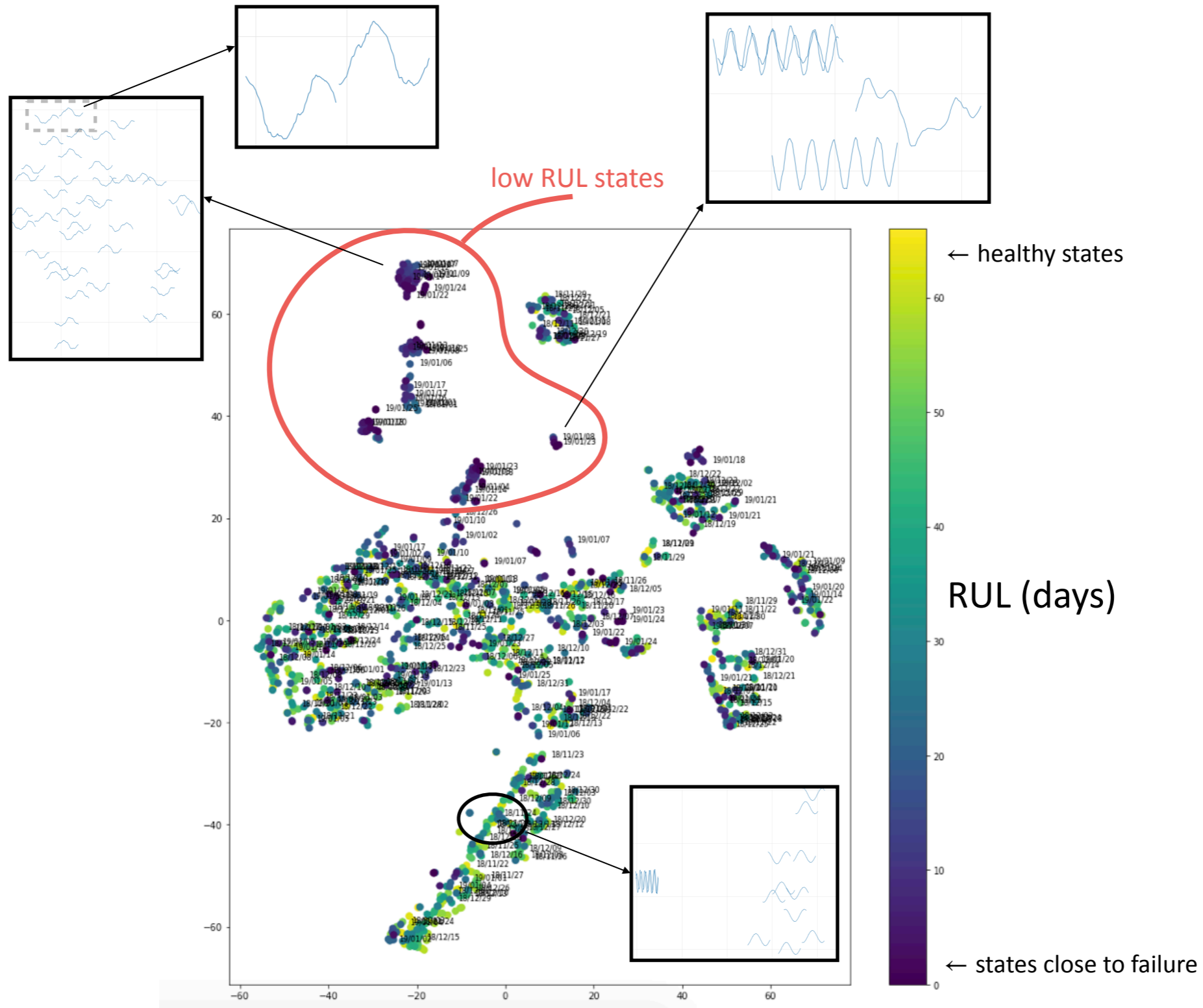
## Visualización de datos / Extracción de características

Análisis del RUL (*Remaining Useful Life*)

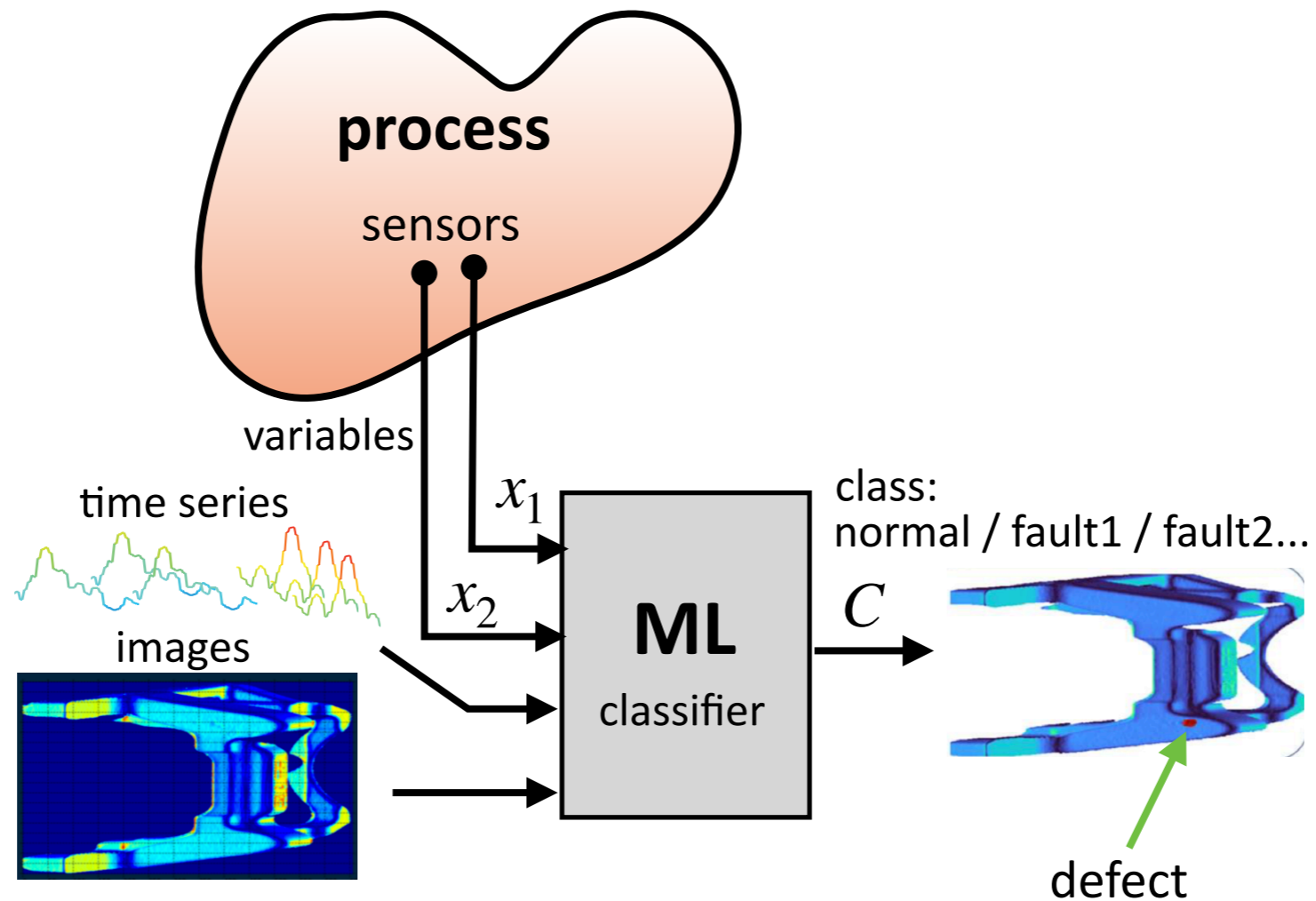


# Reducción de la dimensionalidad

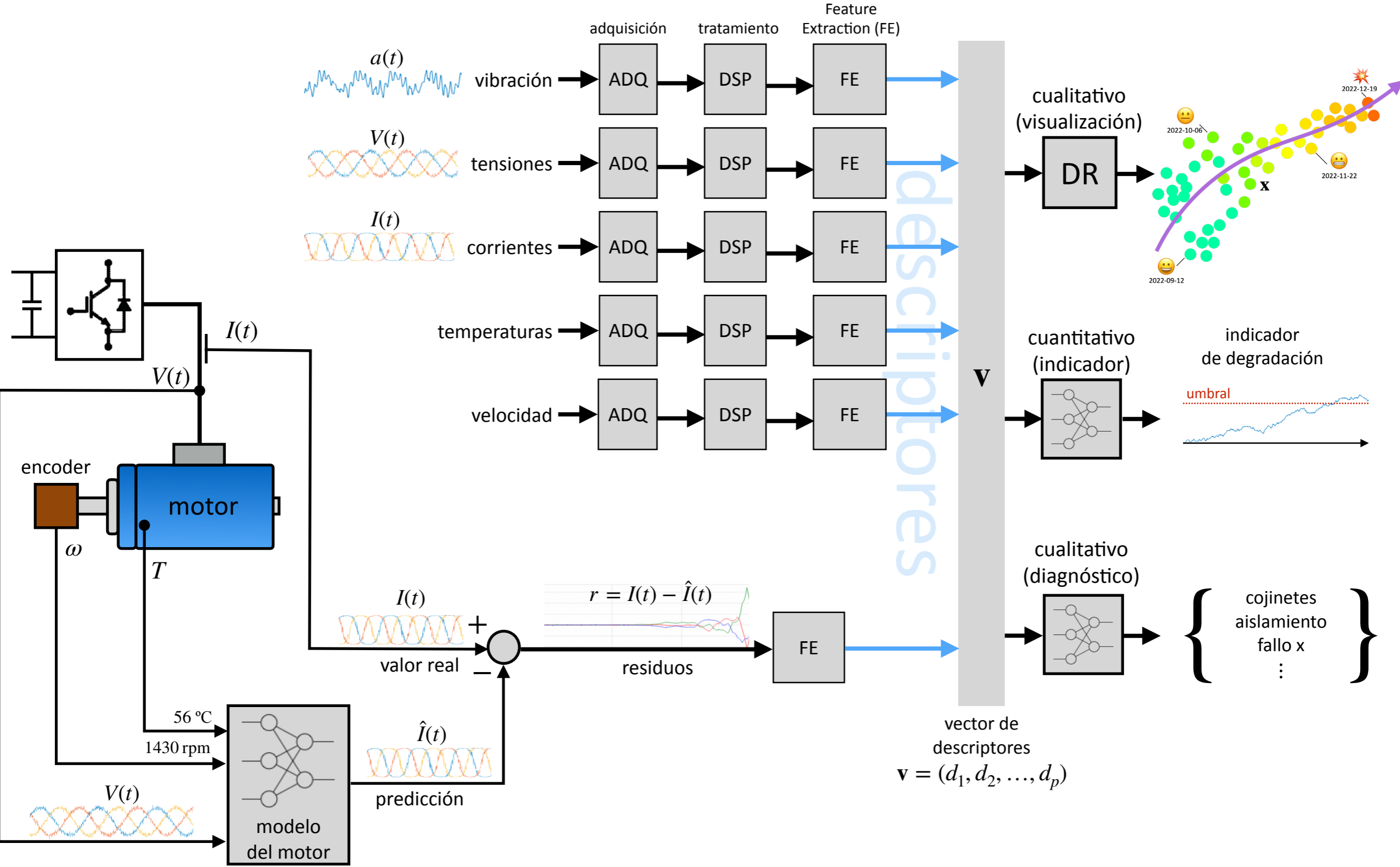
## mapas de estados del proceso



work financed by ArcelorMittal  
under project FUIO-20-227



# Ejemplo de aplicación analítica de datos en motores



redundancia analítica

# Conclusiones

- ML permite construir **modelos** a partir de datos
- Existen **librerías** potentes que simplifican mucho el entrenamiento e inferencia de modelos ML  
Python, Matlab, R, ...
- Permite muchos tipos de **tareas**:
  - detección y diagnóstico de fallos
  - sensores virtuales, gemelo digital
  - predicción, pronóstico
  - análisis exploratorio / visualización de datos
- **Aplicaciones** en problemas de ingeniería eléctrica, electrónica, mecánica...
- Un **futuro** impactante
  - gemelo digital
  - modelos subrogados
  - edge / cloud computing
  - modelos de language (LLM)