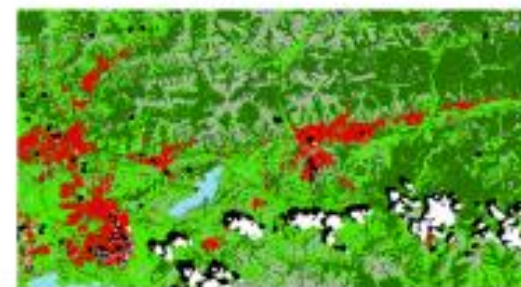
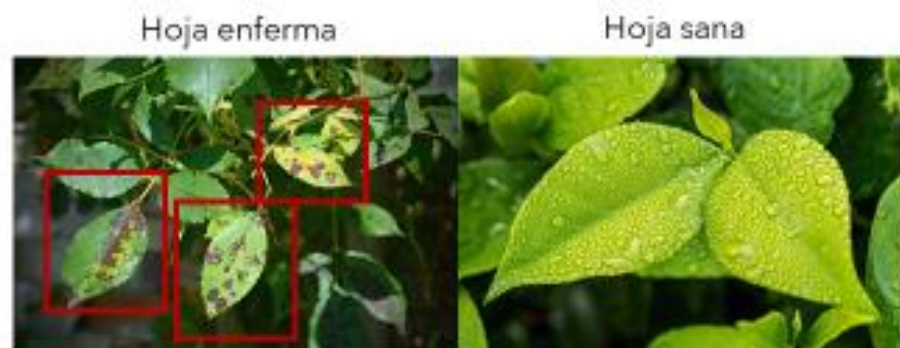


## Sesión 2. Clasificación. Regresión logística

- La tarea de clasificación.
- Discriminantes lineales.
- El algoritmo de regresión logística.

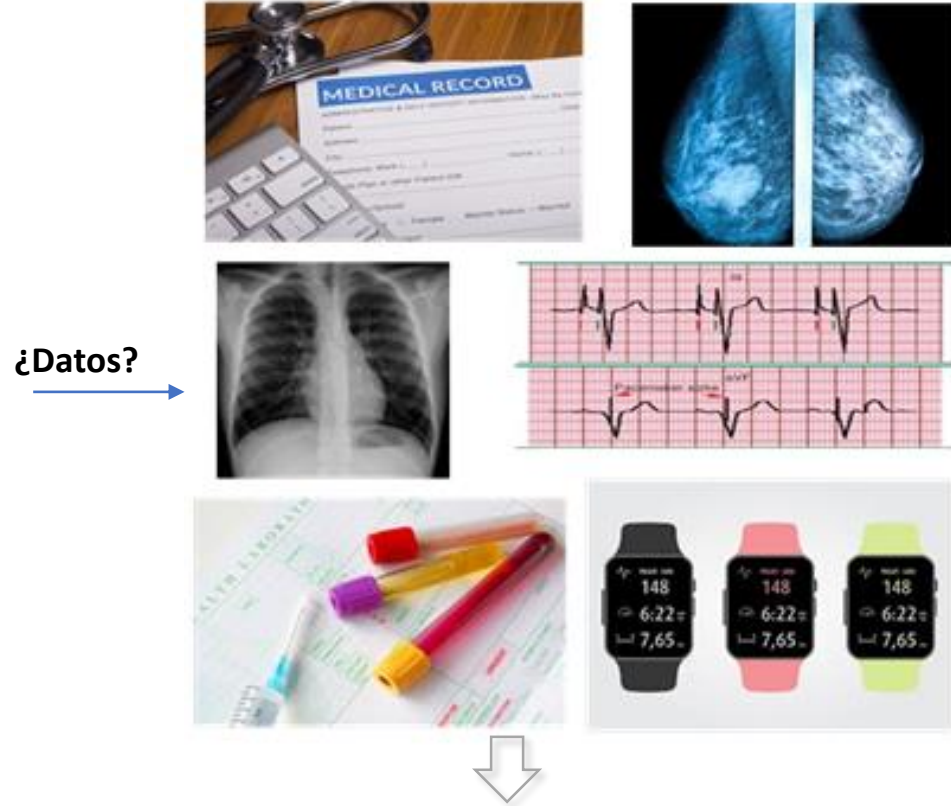
## ¿Cómo resolver una tarea de clasificación?



### Leyenda

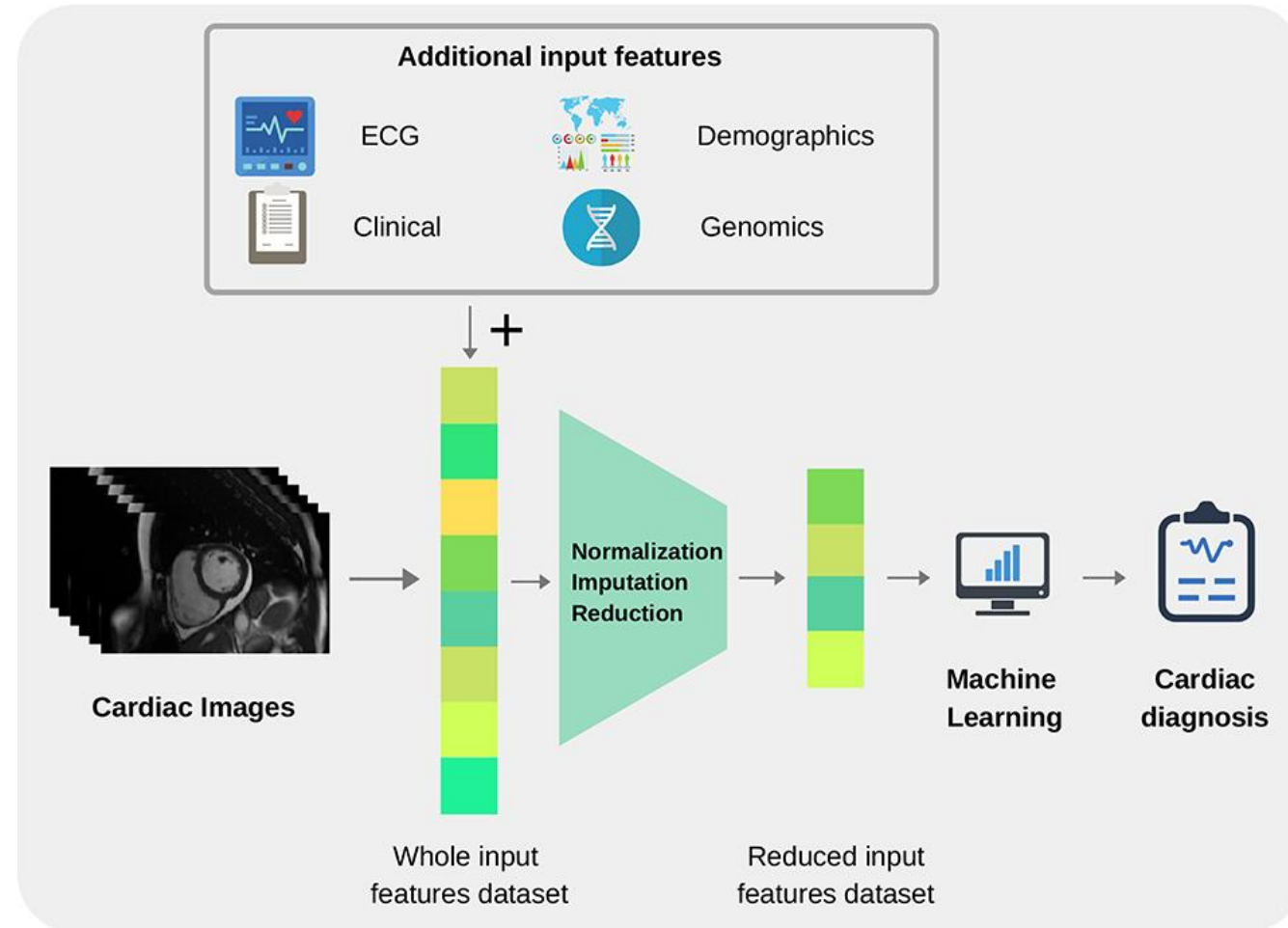
- Cuerpos de Agua
- Vegetación baja
- Vegetación alta
- Suelos desnudos
- Área urbana
- Nubes
- Sombra de nubes

## ML en salud



- Apoyo al diagnóstico de enfermedades.
- Interpretación de imágenes.
- Medicina personalizada.
- Planificación de radioterapia.
- Detección de anomalías en señales de equipos médicos.
- Gestión hospitalaria, ....

### Ejemplo: Image-Based Cardiac Diagnosis With Machine Learning: A Review



<https://www.frontiersin.org/articles/10.3389/fcvm.2020.00001/full>

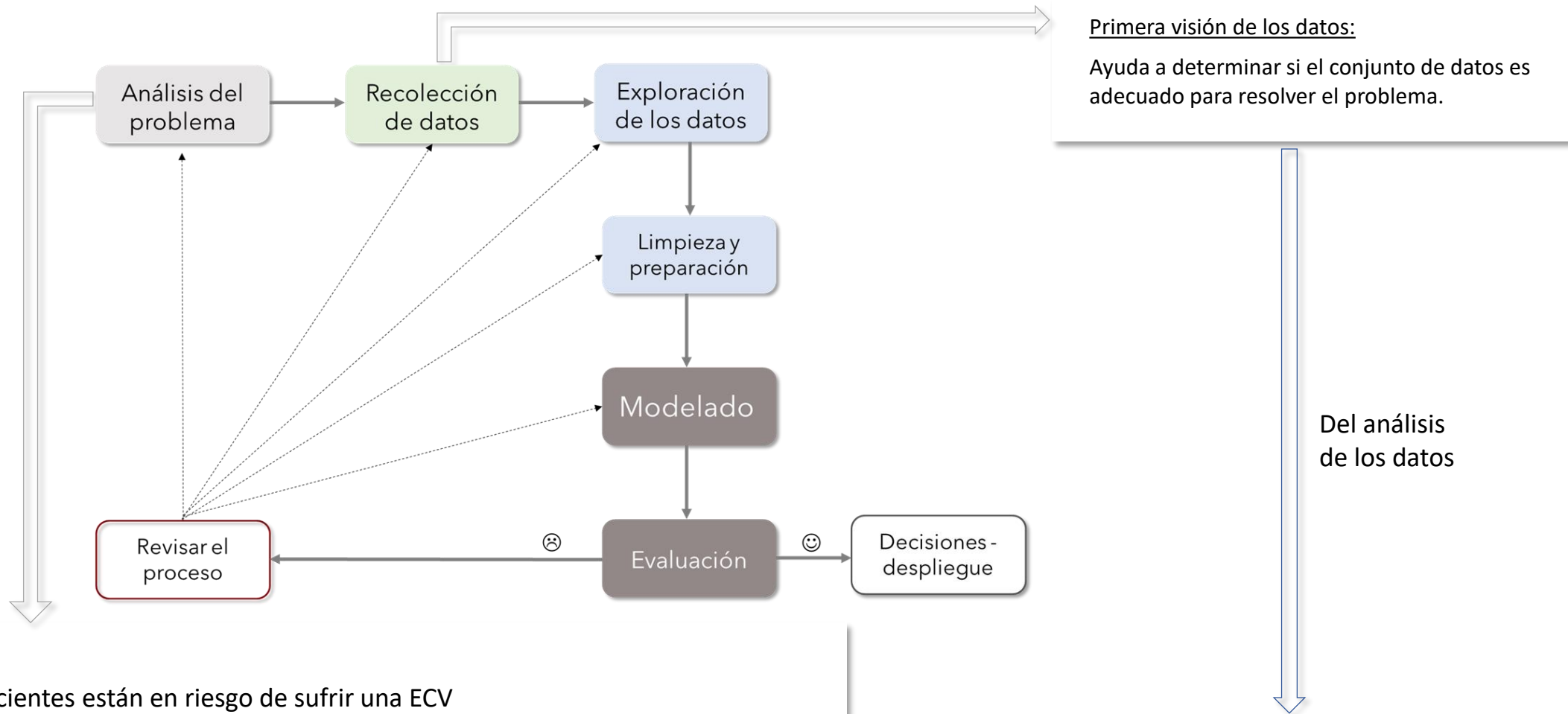
**Proyecto:** estudio sobre las enfermedades cardiovasculares con técnicas de machine learning.

Según la Organización Mundial de la Salud , las enfermedades cardiovasculares (ECV) son la principal causa de muerte en todo el mundo, y se estima que, para el 2030, 23,6 millones de personas morirán por alguna enfermedad de este tipo. Las ECV son un grupo de trastornos del corazón y de los vasos sanguíneos que incluyen, entre otras afecciones, a la enfermedad coronaria, enfermedad cerebrovascular y enfermedad cardíaca reumática. Cuatro de cada 5 muertes por ECV se deben a ataques cardíacos y accidentes cerebrovasculares, y un tercio de estas muertes ocurren prematuramente en personas menores de 70 años. Un pronóstico temprano podría ayudar a tomar decisiones sobre los cambios en el estilo de vida en pacientes de alto riesgo y, a su vez, reducir las complicaciones. En este sentido, se quiere llevar a cabo un estudio, tomando como base los datos de los pacientes registrados en una institución médica, que utilice técnicas de machine learning para la construcción de modelos que, no solo aporten en la identificación de los factores que más inciden en el padecimiento de enfermedad cardíaca, sino también puedan predecir qué pacientes están en riesgo de sufrir una ECV.

age	gender	height	weight	systolic blood pressure	diastolic blood pressure	cholesterol	glucose	smoke	Alcohol	active	ECV
18393	2	168	62.0	110	80	normal	normal	0	0	1	0
20228	1	156	85.0	140	90	well above normal	normal	0	0	1	1
18857	1	165	64.0	130	70	well above normal	normal	0	0	0	1
17623	2	169	82.0	150	100	normal	normal	0	0	1	1
17474	1	156	56.0	100	60	normal	normal	0	0	0	0
21914	1	151	67.0	120	80	above normal	above normal	0	0	0	0
22113	1	157	93.0	130	80	well above normal	normal	0	0	1	0
22584	2	178	95.0	130	90	well above normal	well above normal	0	0	1	1
17668	1	158	71.0	110	70	normal	normal	0	0	1	0
19834	1	164	68.0	110	60	normal	normal	0	0	0	0
22530	1	169	80.0	120	80	normal	normal	0	0	1	0

Número de datos: 69.997  
Número de variables: 12

## Ciclo de ML



### Objetivos:

- Predecir qué pacientes están en riesgo de sufrir una ECV
- Identificar los factores que más inciden en el padecimiento de enfermedad cardíaca.

### Tareas de aprendizaje:

Para alcanzar el primer objetivo se puede realizar una tarea de clasificación que utilice como variable target “ECV”, la cual indica la presencia o ausencia de enfermedad cardiovascular. Dependiendo del algoritmo a utilizar se podría hacer un análisis de las variables más significativas, lo cual podría permitir alcanzar el segundo objetivo. O utilizar alguna librería que permita realizar una interpretación del modelo.

Del análisis del problema

## Construcción de un modelo de clasificación

Solución propuesta: modelo de clasificación que permita predecir si un paciente tiene posibilidad de ser diagnosticado con ECV.





Variables de entrada =  $x$

Clase =  $y$

age	gender	height	weight	systolic blood pressure	diastolic blood pressure	cholesterol	glucose	smoke	Alcohol	active	ECV
18393	2	168	62.0	110	80	normal	normal	0	0	1	0
20228	1	156	85.0	140	90	well above normal	normal	0	0	1	1
18857	1	165	64.0	130	70	well above normal	normal	0	0	0	1
17623	2	169	82.0	150	100	normal	normal	0	0	1	1
17474	1	156	56.0	100	60	normal	normal	0	0	0	0
21914	1	151	67.0	120	80	above normal	above normal	0	0	0	0
22113	1	157	93.0	130	80	well above normal	normal	0	0	1	0
22584	2	178	95.0	130	90	well above normal	well above normal	0	0	1	1
17668	1	158	71.0	110	70	normal	normal	0	0	1	0
19834	1	164	68.0	110	60	normal	normal	0	0	0	0
22530	1	169	80.0	120	80	normal	normal	0	0	1	0
18815	2	173	60.0	120	80	normal	normal	0	0	1	0
14791	2	165	60.0	120	80	normal	normal	0	0	0	0
19809	1	158	78.0	110	70	normal	normal	0	0	1	0

...

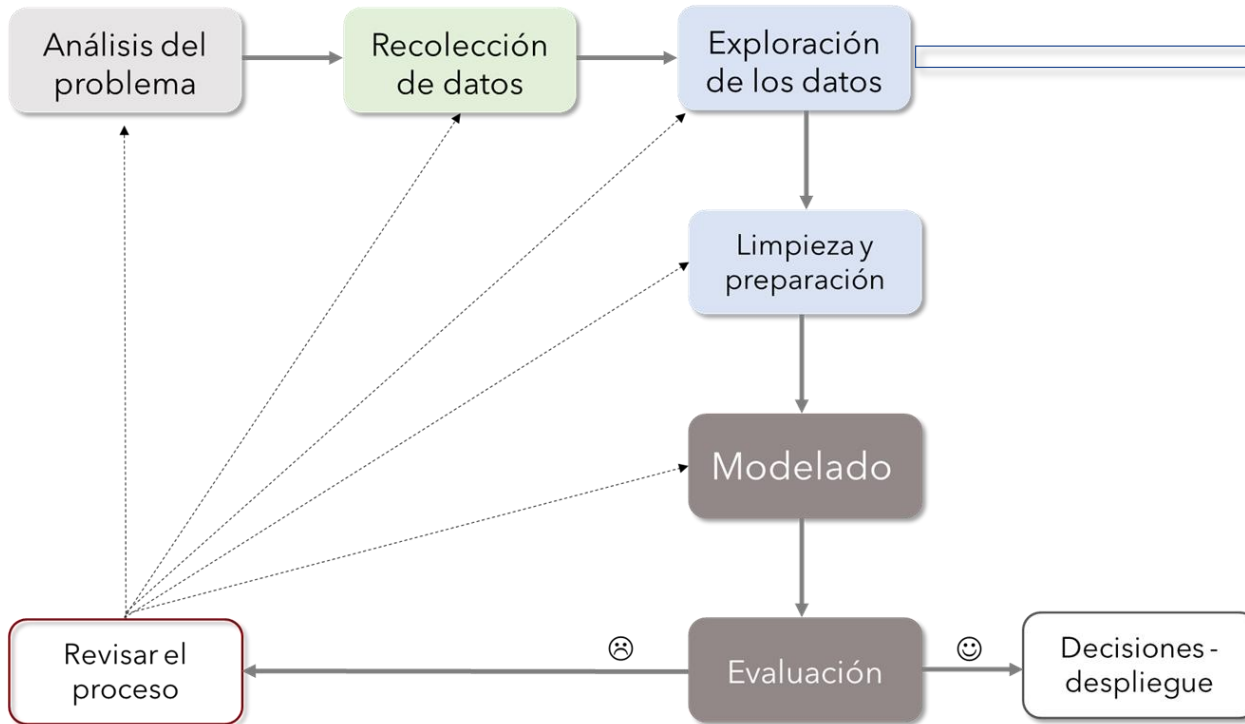
Algoritmo de clasificación

Modelo de clasificación

ECV = SI  
ECV = NO

Nuevo dato (paciente)

Edad	Género	Altura	Peso	Presión sistólica	Presión diastólica	Colesterol	Glucosa	Fumador	Alcohol	Activo
50	F	171	76.0	90	60	Normal	Elevado	NO	NO	SI



## Importante: Diccionario de los datos

There are 3 types of input features:

Objective: factual information.

Examination: results of medical examination.

Subjective: information given by the patient.

Features:

Age | Objective Feature | age | int (days)

Height | Objective Feature | height | int (cm) |

Weight | Objective Feature | weight | float (kg) |

Gender | Objective Feature | gender | categorical code |

Systolic blood pressure | Examination Feature | ap\_hi | int positive |

Diastolic blood pressure | Examination Feature | ap\_lo | int positive |

Cholesterol | Examination Feature | cholesterol | normal, above normal, well above normal |

Glucose | Examination Feature | gluc | normal, above normal, well above normal |

Smoking | Subjective Feature | smoke | binary |

Alcohol intake | Subjective Feature | alco | binary |

Physical activity | Subjective Feature | active | binary |

Presence or absence of cardiovascular disease | Target Variable | cardio | binary |

All the dataset values were collected at the moment of medical examination.

¿Limpieza de datos?

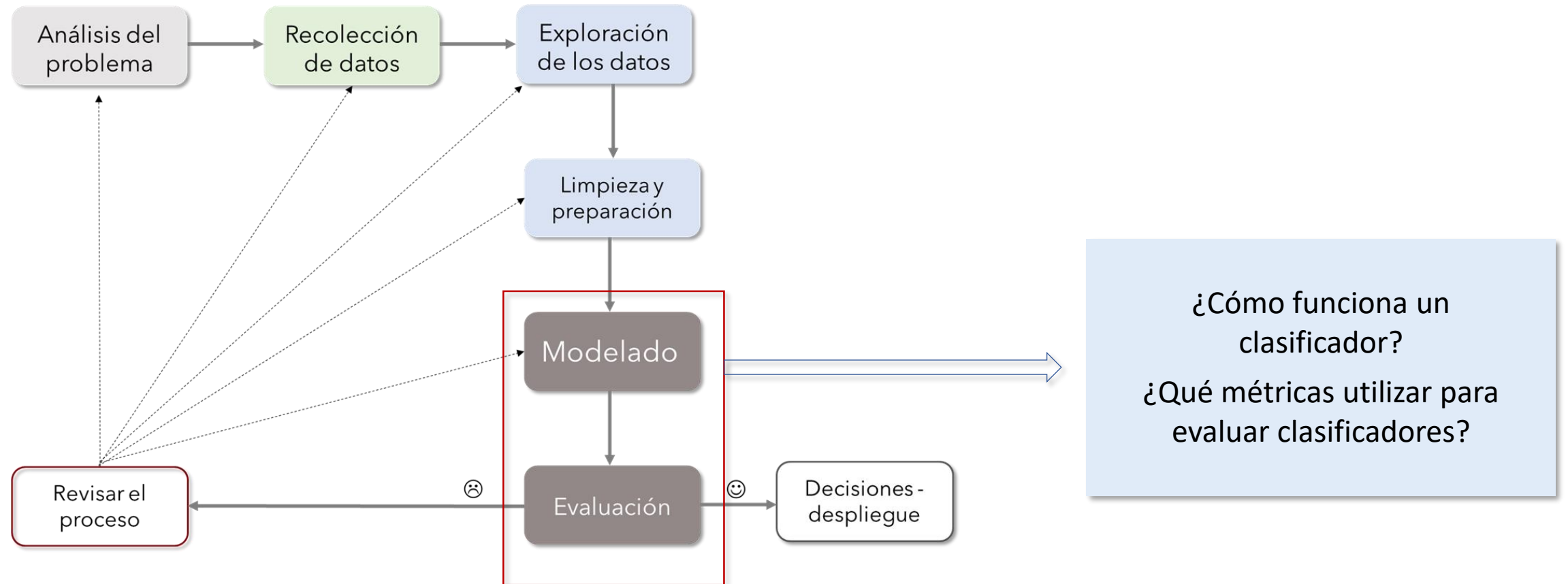
¿Cómo construir una buena representación de los datos para el algoritmo seleccionado?

¿Qué transformaciones aplicar?

Descripción de los datos  
(perfilamiento)

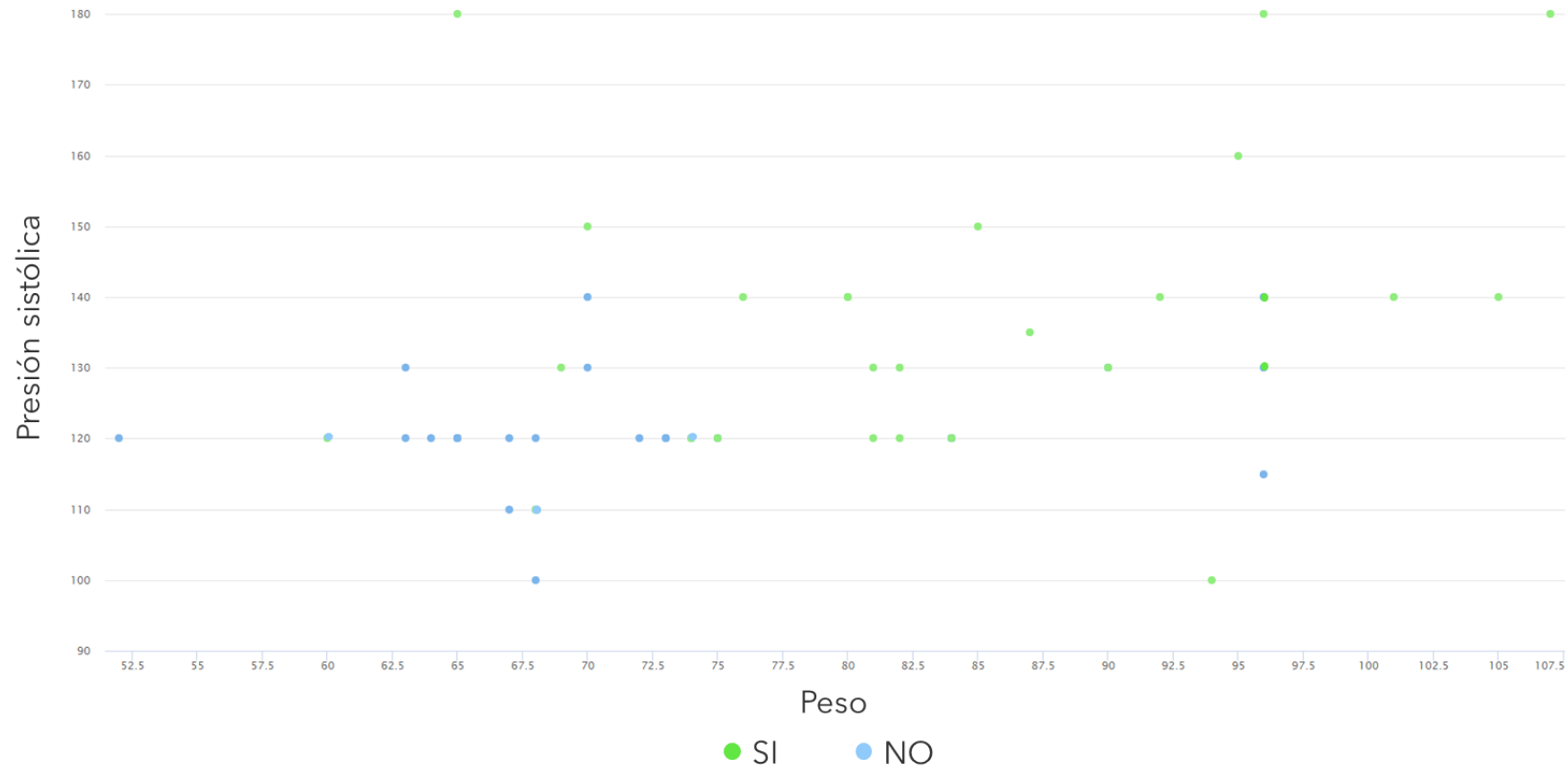
### Propuesta para la preparación de los datos:

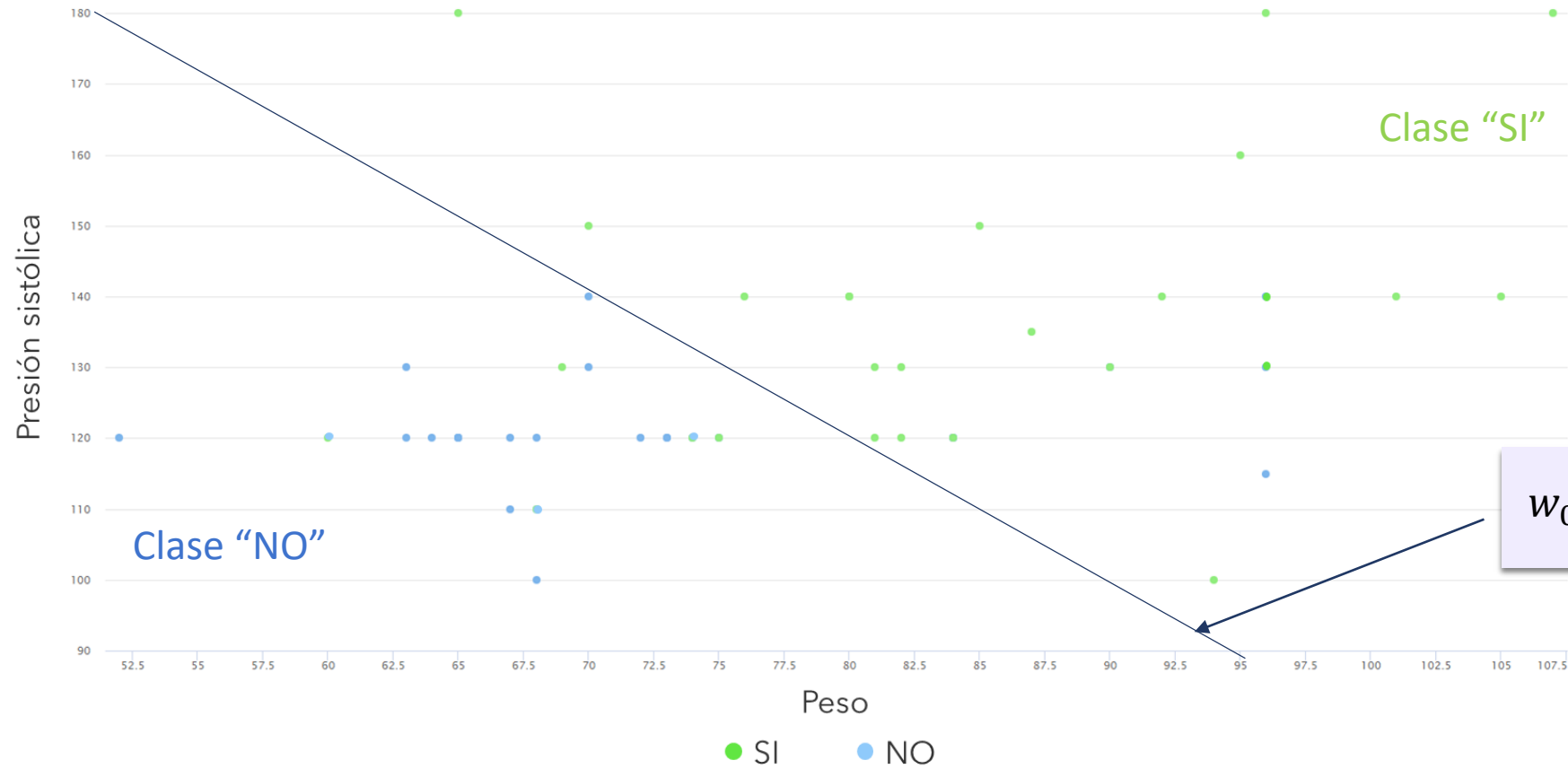
- Las variables “cholesterol” y “gluc” son categóricas. Como solo asumen tres categorías, se puede aplicar una numerización 1-de-n (one-hot).
- Con base en el algoritmo a utilizar, será conveniente normalizar o estandarizar los datos.





## ¿Qué hace un algoritmo de aprendizaje para resolver este problema de clasificación?

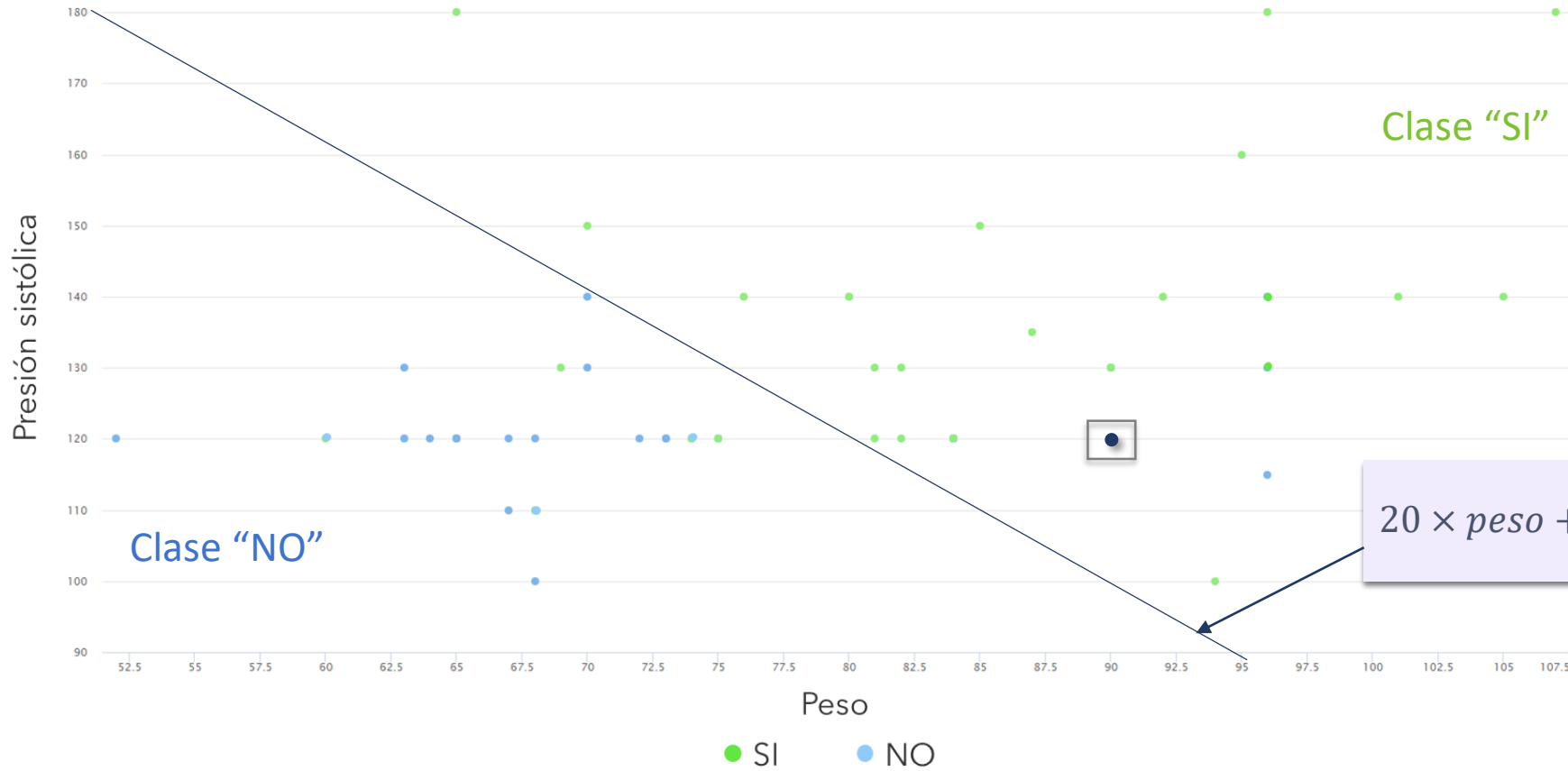




$$\text{Clase estimada} = f_{\text{umbral}}(w_0 + w_1x_1 + w_2x_2)$$

$$f_{\text{umbral}}(h) = \begin{cases} \text{clase "SI"} & \text{si } h \geq 0 \\ \text{clase "NO"} & \text{si } h < 0 \end{cases} \quad \leftarrow \text{umbral}$$

## ¿Cómo clasificar?



clase SI  $\Rightarrow 20 \times \text{peso} + 10 \times \text{presión sistólica} - 2800 \geq 0$

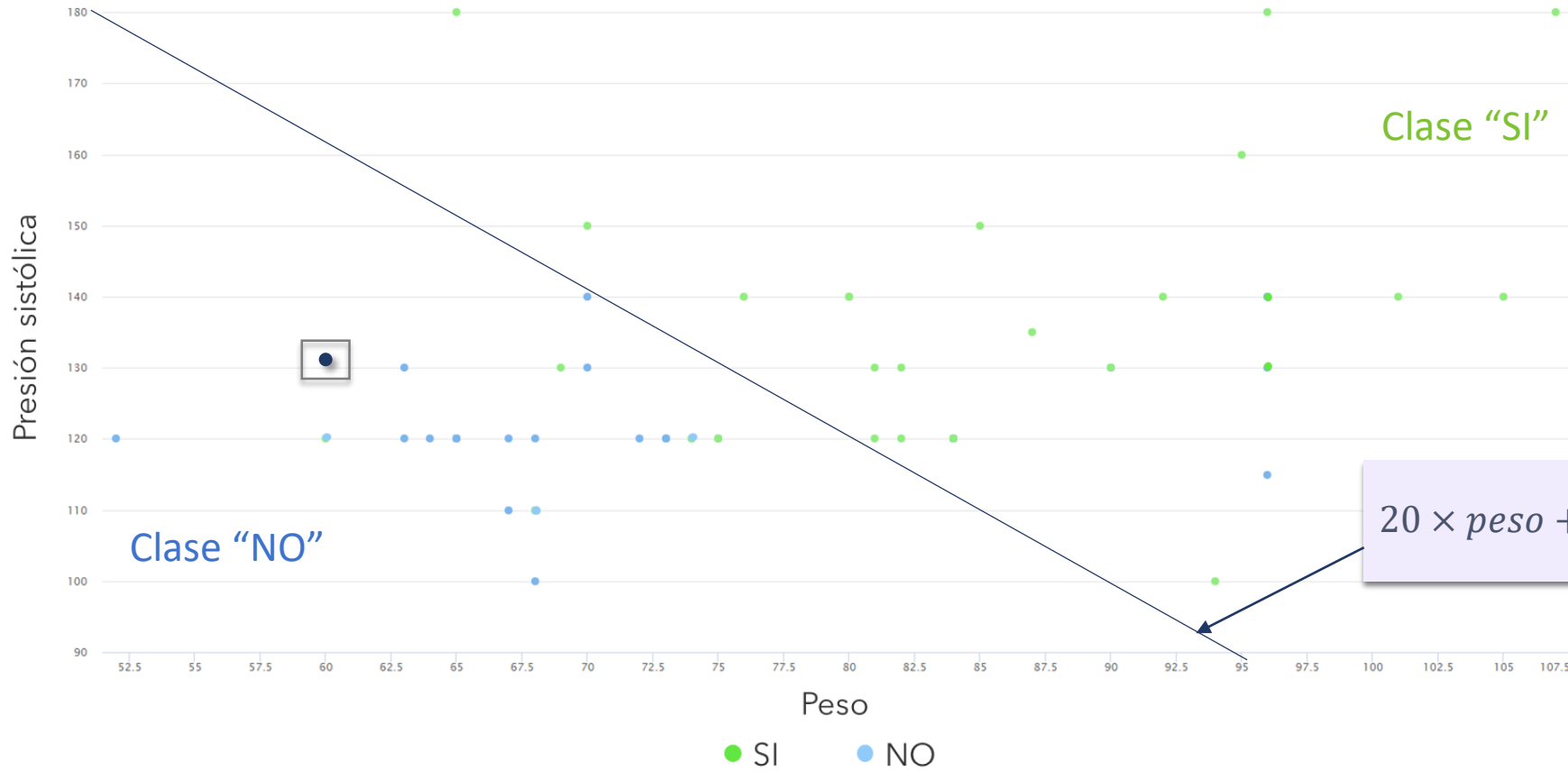
clase NO  $\Rightarrow 20 \times \text{peso} + 10 \times \text{presión sistólica} - 2800 < 0$

Si (90,120):

$$20 \times 90 + 10 \times 120 - 2800 = 200$$

$\Rightarrow$  Clase  
SI

## ¿Cómo clasificar?



clase SI si:  $20 \times \text{peso} + 10 \times \text{presión sistólica} - 2800 \geq 0$

clase NO si:  $20 \times \text{peso} + 10 \times \text{presión sistólica} - 2800 < 0$

Si (60,130):

$$20 \times 60 + 10 \times 130 - 2800 = -300$$

$\Rightarrow$  Clase  
NO

## Algunos algoritmos:

### **Clasificadores binarios**

- ☐ Regresión logística
- ☐ Máquinas de vectores de soporte

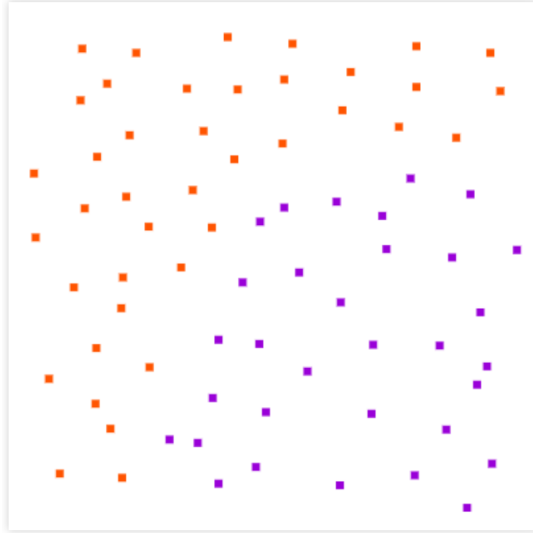
### **Clasificadores multiclase**

- ☐ Árboles de decisión
- ☐ Redes neuronales
- ☐ Basados en probabilidades (Naïve Bayes)
- ☐ K-vecinos más cercanos
- ☐ Ensembles (metal algoritmos). Ej. Random Forest
- ☐ ...

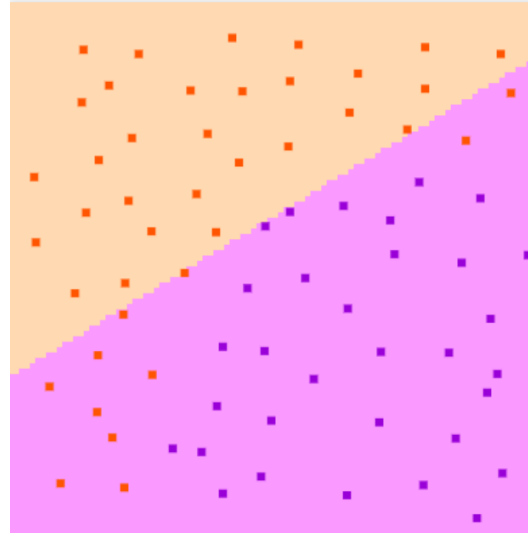


*En general...*

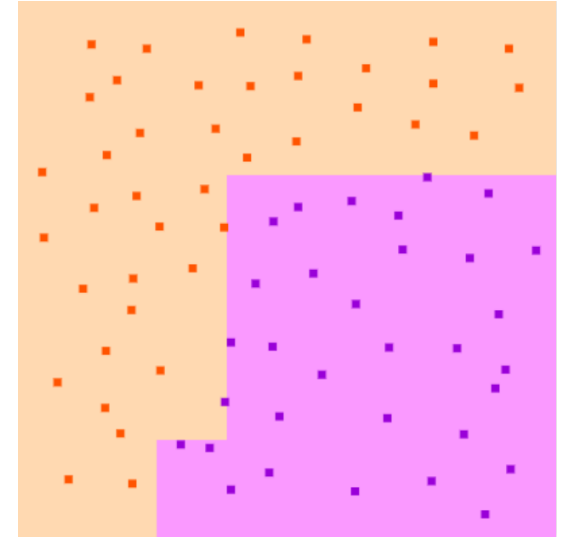
Datos



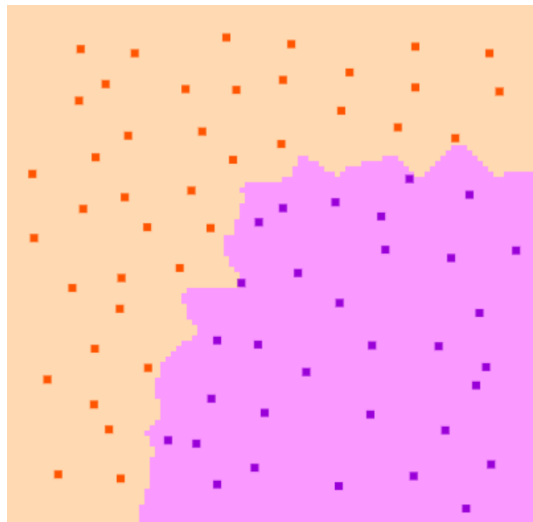
Discriminante lineal



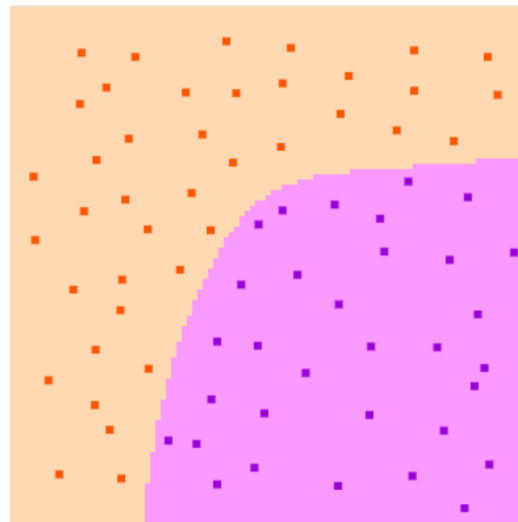
Árbol de decisión



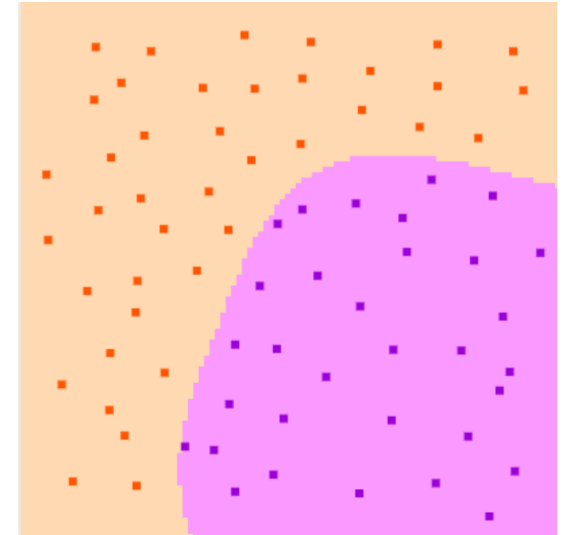
K-vecinos más cercanos



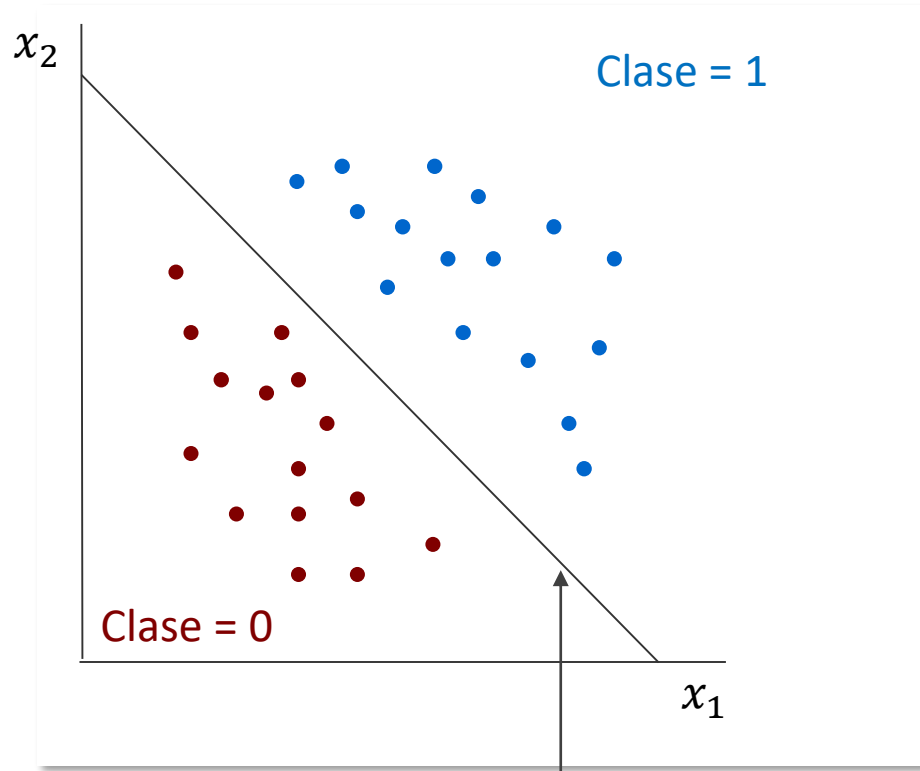
Red neuronal



Máquina de vectores de soporte



Problema de clasificación de dos clases:



Límite de decisión lineal

$$w_0 + w_1x_1 + w_2x_2 = 0$$

En vez de una función umbral:

$$f_{umbral}(w_0 + w_1x_1 + w_2x_2)$$

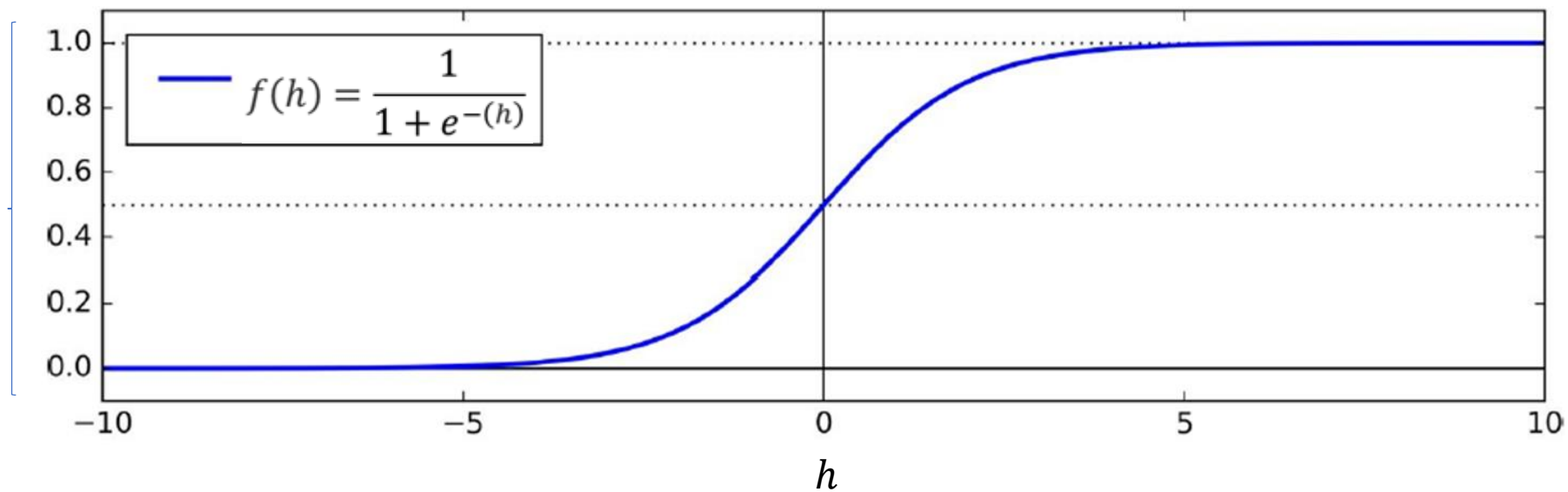
$$f_{umbral}(h) = \begin{cases} \text{clase "1"} & \text{si } h \geq 0 \\ \text{clase "0"} & \text{si } h < 0 \end{cases}$$

En regresión logística se utiliza la función logística:

$$f_{logística}(w_0 + w_1x_1 + w_2x_2)$$

$$f_{logística}(h) = \frac{1}{1 + e^{-(h)}}$$

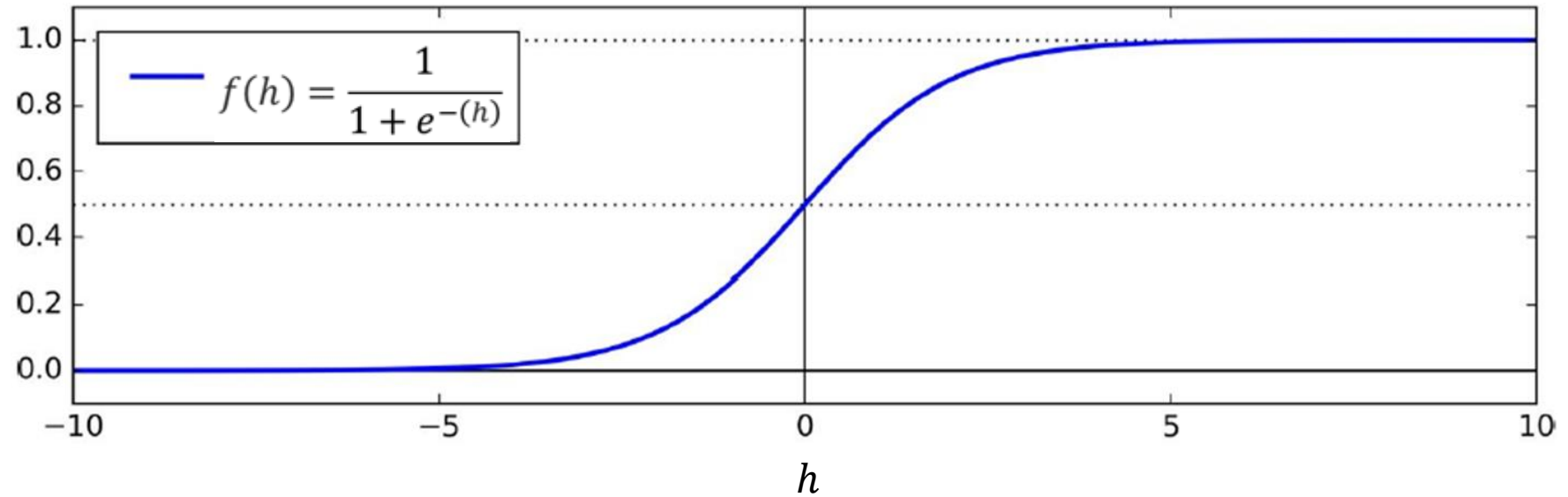
Rango =  $[0,1]$



- Puede ser interpretado como la probabilidad de un evento binario.
- Representa la probabilidad de la clase dada la instancia  $\mathbf{x}$  como:

$$P(y|\mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w} \cdot \mathbf{x} + b)}} = \text{probabilidad estimada de } y = 1, \text{ dado que se presenta } \mathbf{x}$$

$$\begin{cases} \text{Clase 1: } P(y = 1|\mathbf{x}) \\ \text{Clase 0: } P(y = 0|\mathbf{x}) = 1 - P(y = 1|\mathbf{x}) \end{cases}$$



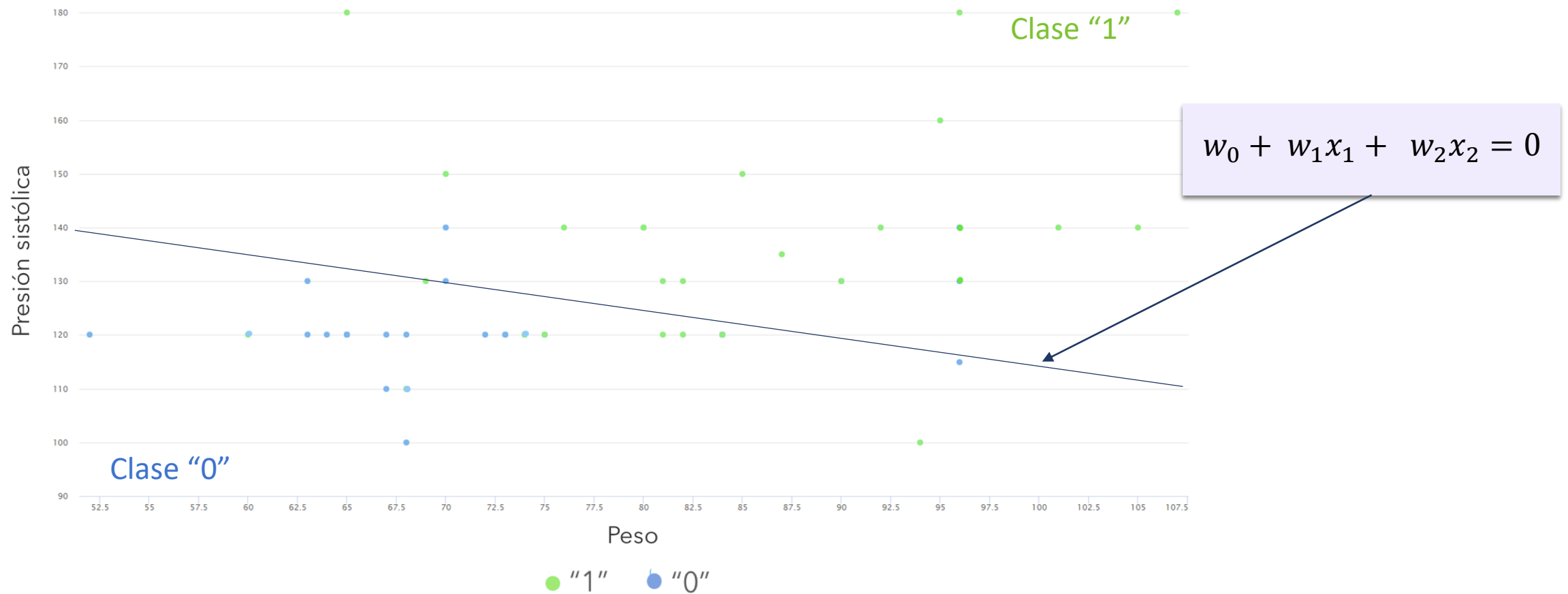
Límite de decisión  $\Rightarrow \mathbf{w} \cdot \mathbf{x} = 0$

$$\rightarrow \theta(\mathbf{w} \cdot \mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w} \cdot \mathbf{x} + b)}} = \frac{1}{1 + e^0} = 0.5$$

- Para clasificar se utiliza un umbral:

*predice clase 1 si  $f_{logística}(w_0 + w_1x_1 + w_2x_2) \geq 0.5$*

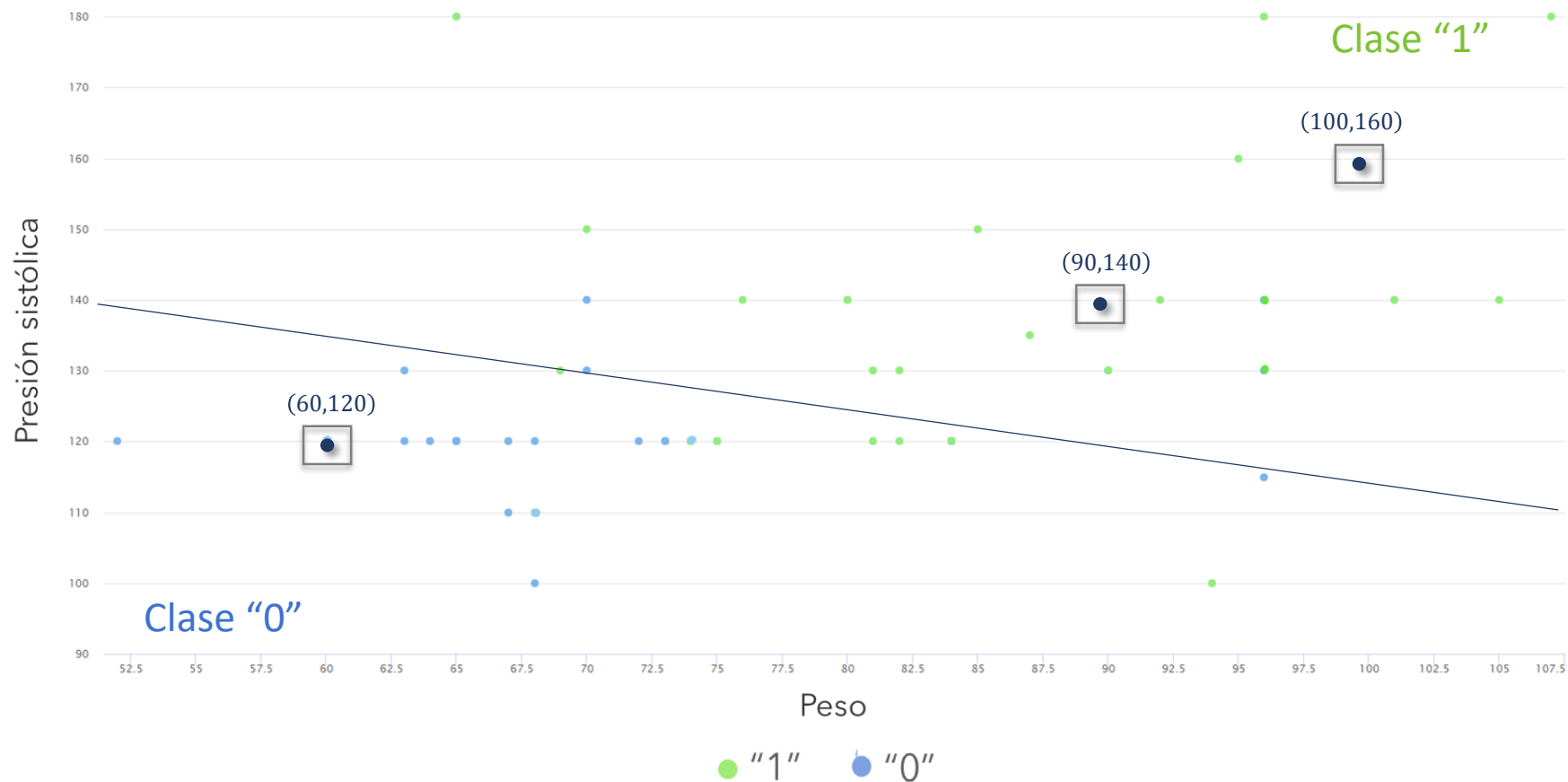
*predice clase 0 si  $f_{logística}(w_0 + w_1x_1 + w_2x_2) < 0.5$*



$$\text{Límite de decisión} = 0.017 \times \text{peso} + 0.066 \times \text{presión sistólica} - 9.15 = 0$$

$$P(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-(0.017 \times \text{peso} + 0.066 \times \text{presión sistólica} - 9.15)}}$$





clase 1:  $f_{logística}(h) \geq 0.5$

clase 0:  $f_{logística}(h) < 0.5$

$$P(y = 1 | x = (90, 140)) = \frac{1}{1 + e^{-(0.017 \times 90 + 0.066 \times 140 - 9.15)}} = 0.71$$

→ Clase "1"

$$P(y = 1 | x = (60, 120)) = \frac{1}{1 + e^{-(0.017 \times 60 + 0.066 \times 120 - 9.15)}} = 0.40$$

→ Clase "0"

$$P(y = 1 | x = (100, 160)) = \frac{1}{1 + e^{-(0.017 \times 100 + 0.066 \times 160 - 9.15)}} = 0.96$$

→ Clase "1"

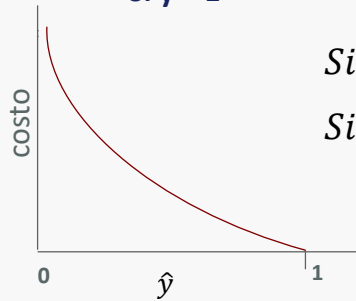
¿Por qué la diferencia en las probabilidades?

## ¿Cómo determinar los parámetros?

Recordar: el  $\text{costo}(\hat{y}, y)$  mide la discrepancia entre lo que predice el modelo y el valor real.

$$\text{Costo para una instancia} = \begin{cases} -\log(\hat{y}) & \text{si } y = 1 \\ -\log(1 - \hat{y}) & \text{si } y = 0 \end{cases}$$

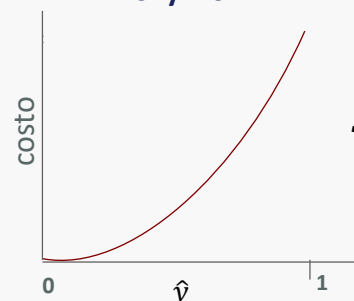
Si  $y = 1$



Si  $\hat{y} = 1 \Rightarrow \text{costo} = 0$

Si  $\hat{y} = 0 \Rightarrow \text{costo} \rightarrow \infty$

Si  $y = 0$



Si  $\hat{y} = 0 \Rightarrow \text{costo} = 0$

Si  $\hat{y} = 1 \Rightarrow \text{costo} \rightarrow \infty$

Otra representación:

$$\text{Costo}(\hat{y}, y) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

Entonces, función de costo de la regresión logística

$$J(\mathbf{W}) = -\frac{1}{N} \sum_{i=1}^N y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)$$

¿Cómo funciona el algoritmo a utilizar?  
¿Qué requerimientos tiene?  
¿Cómo ajustar sus parámetros?

Solo atributos numéricos

	$x_1$	$x_2$		$x_d$	$y$
Instancia 1					
Instancia 2					
Instancia N					

Parámetros del algoritmo:  
 $C$ , hiperparámetro de regularización

Costo log (cross entropy binaria)

$$J(W) = -\frac{1}{N} \sum_{i=1}^N y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i) + C \sum_{j=1}^d (w_j)^2$$

Aprendizaje supervisado  
Tarea de clasificación

Regresión logística

Procedimiento de  
Optimización

Función de costo

Algoritmo de descenso por el gradiente

Entrada = Conjunto de datos  $D = \{(x_i, y_i)\}_{i=1..N}$

Comienzo\_procedimiento

Inicializar los pesos  $W(0)$  y la constante de aprendizaje  $\alpha$

Para  $k = 0$  hasta  $\text{num\_iteraciones}$

Calcular el gradiente  $g_k = \nabla J(w(k))$

Actualizar los parámetros:

$$W(k+1) = W(k) - \alpha g_k$$

Fin-procedimiento

Salida =  $\{W\}$

Modelo estimado

$$\hat{y} = \text{función}_{\text{logística}}(w_0 + w_1 \times x_1 + w_2 \times x_2 + \dots + w_d \times x_d)$$

parámetros del modelo =  $\{w_0, w_1, \dots, w_d\}$

$d = \text{No. de variables}$

## Próxima sesión

---

- ¿Cómo evaluar los modelos?
- Capacidad de generalización, complejidad de modelos y sobreajuste.
- Dilema bias-varianza.
- Regularización. Hiperparámetros.
- Selección de modelos y técnicas de validación.