



Montreal, July 4–8

# DLM12022

## Introduction to machine learning

By

Pierre-Marc Jodoin

(Hugo Larochelle)



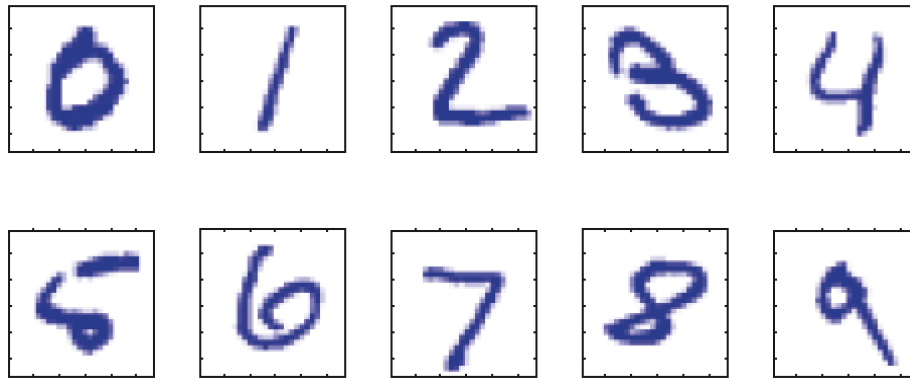
UNIVERSITÉ DE  
SHERBROOKE

Killer question...

What is machine learning?



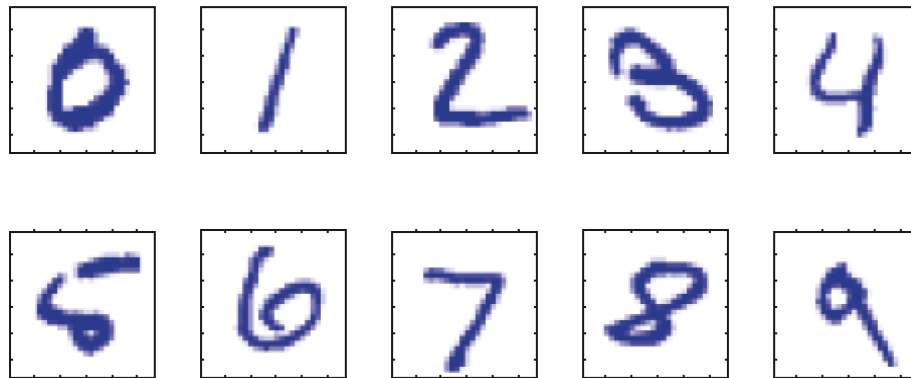
**Question :** how can one recognize hand written digits?



**Answer :** Design your own rules?

- A series of aligned pixels => '1'
- A circle of pixels => '0'
- Etc.

**Question :** how can one recognize hand written digits?

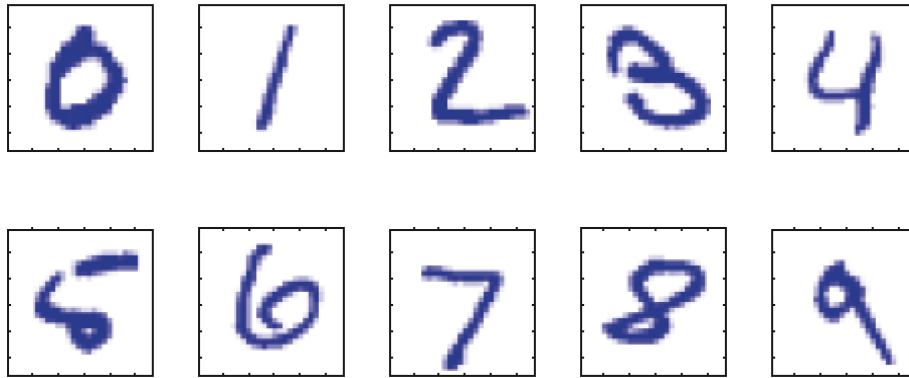


**Answer :** ~~Design your own rules?~~ **Wrong**

➤ Bad generalization



**Question :** how can one recognize hand written digits?

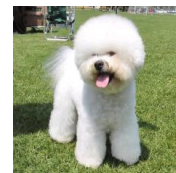


**Answer :** ~~Design your own rules?~~ **Wrong**

➤ Bad generalization



➤ Often difficult



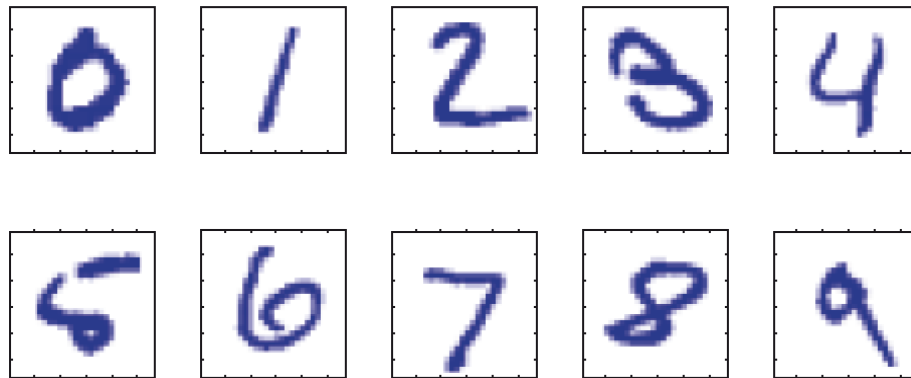
Vs



Dogs

Birds

**Question** : how can one recognize hand written digits?



**Answer** : Let the computer « **learn** » the rules

➤ *Main goal of machine learning*

# Two large families

Supervised learning

Unsupervised learning

# Two large families

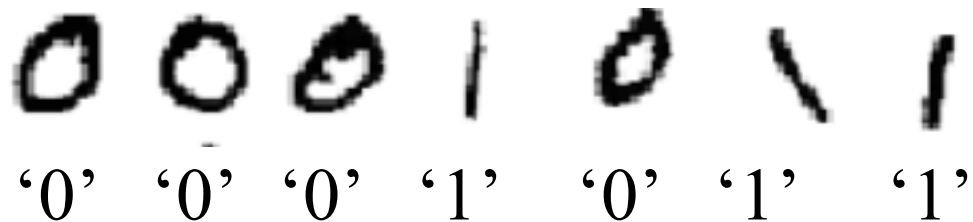
## **Supervised learning**

## Unsupervised learning

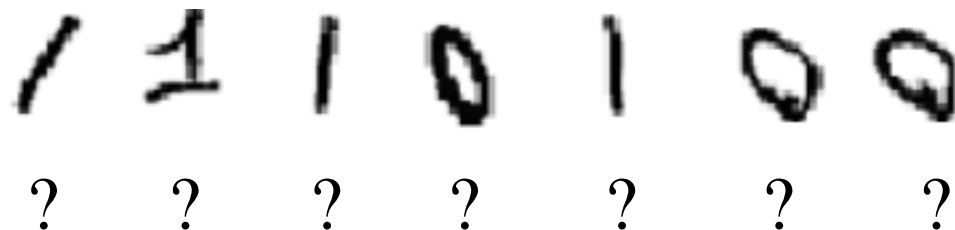


# Supervised learning

Provide the algorithm with **annotated training data**

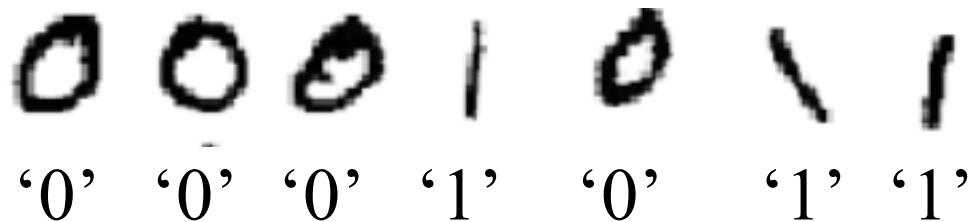


...and the algorithm returns a function capable of **generalizing** on new data



# Supervised learning

Provide the algorithm with **annotated training data**



The **training dataset**

$$D = \{(\vec{x}_1, t_1), (\vec{x}_2, t_2), \dots, (\vec{x}_N, t_N)\}$$

where  $\vec{x}_i \in \mathbb{R}^d$  is an **input** and  $t_i$  is a **target**

# Goal of a supervised machine learning method

From a **training dataset**:  $D = \{(\vec{x}_1, t_1), (\vec{x}_2, t_2), \dots, (\vec{x}_N, t_N)\}$

$\vec{x}_i \in \mathfrak{R}^d$  input data

$t_i$  target associated to  $\vec{x}_i$

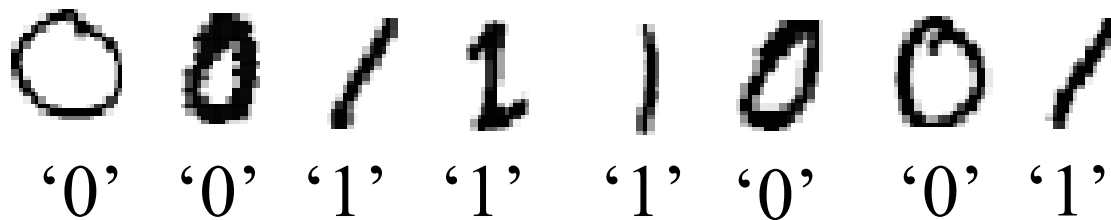
the goal is **to learn** a function that may predict  $t_i$  given  $\vec{x}_i$

$$y_{\vec{w}}(\vec{x}_i) \rightarrow t_i$$

where  $\vec{w}$  are the **parameters** of the model.

# Supervised learning

Once the model  $y_{\vec{w}}(\vec{x})$  is trained, we use a **test set**  $D_{test}$  to gauge the **generalization** capabilities of the model.



# Two large families

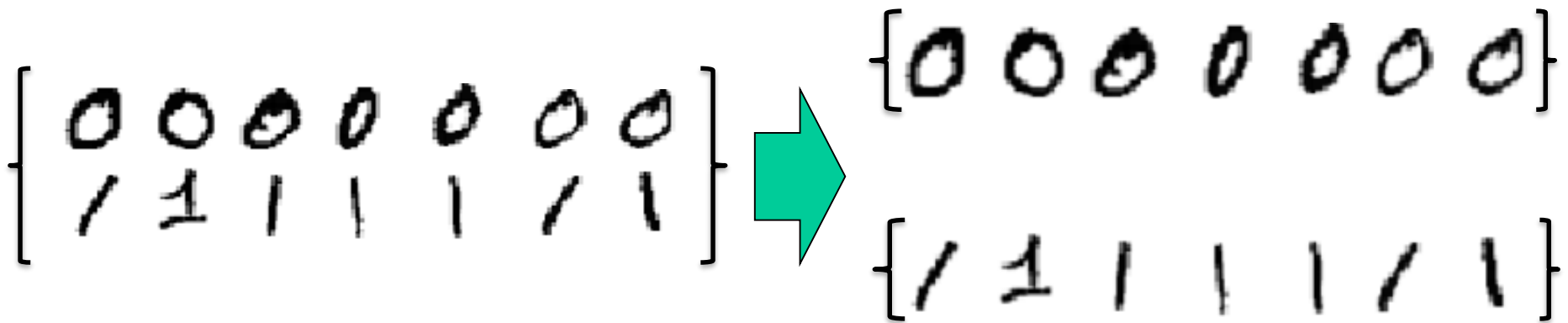
Supervised learning

**Unsupervised learning**

# Unsupervised learning

When no target is explicitly provided

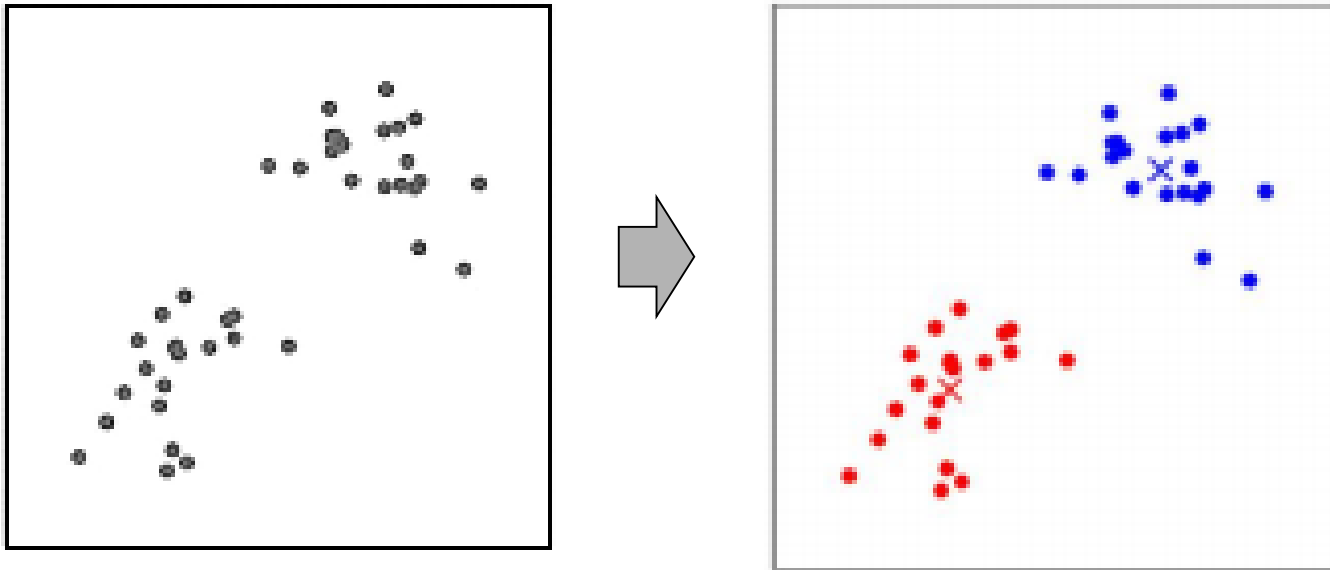
➤ E.g. data *clustering*



# Unsupervised learning

When no target is explicitly provided

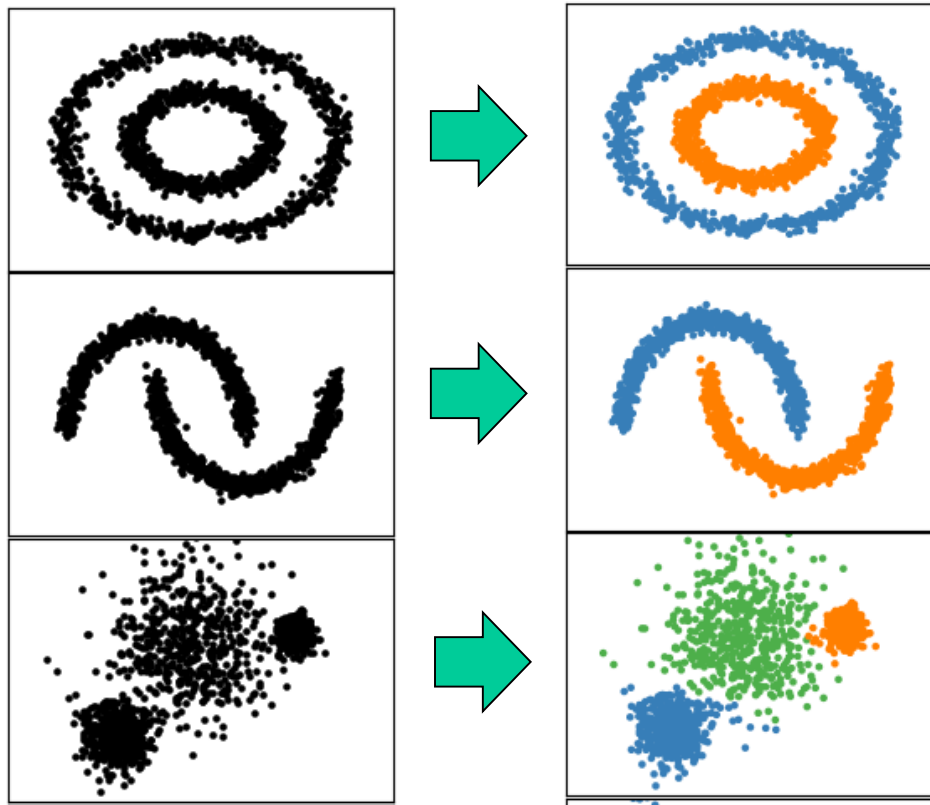
➤ E.g. data *clustering*



# Unsupervised learning

When no target is explicitly provided

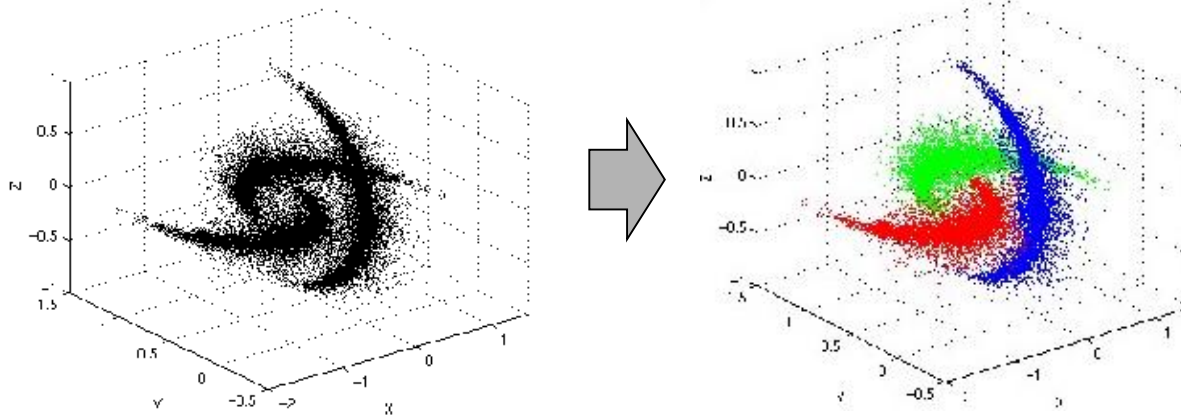
➤ E.g. data *clustering*





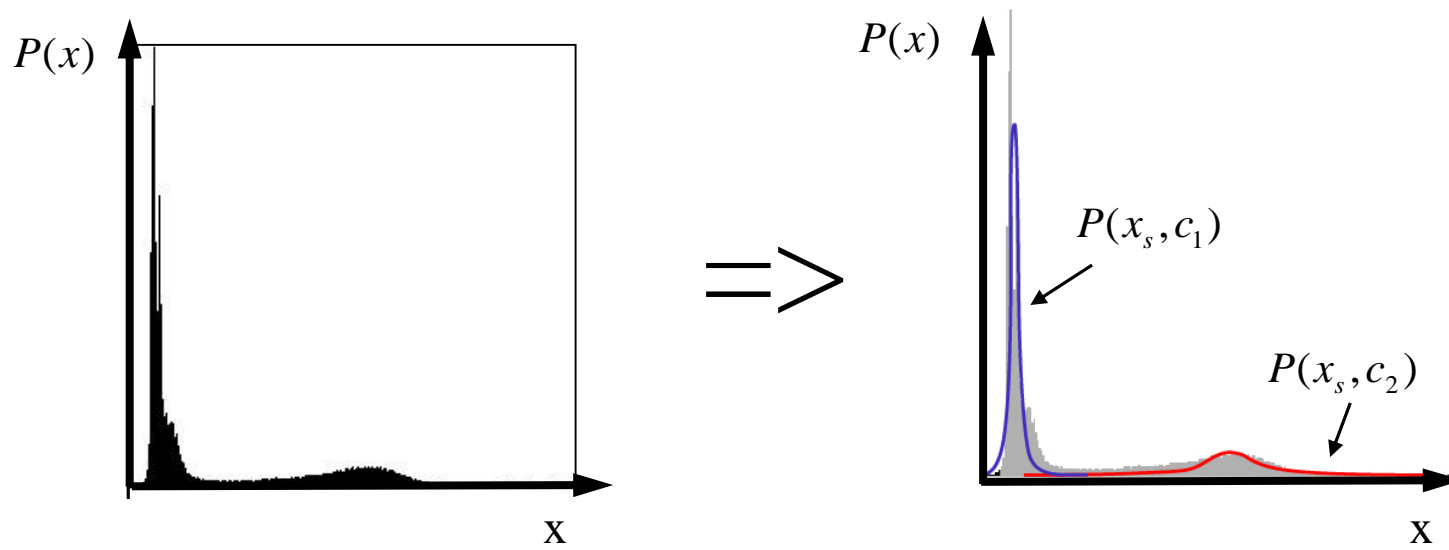
# Unsupervised learning

No limit to dimensionality. Could be 3D, 4D,...100kD



# Unsupervised learning

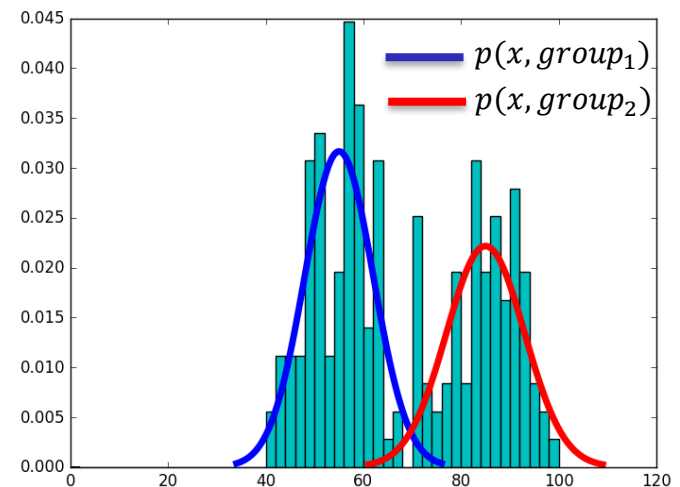
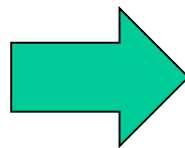
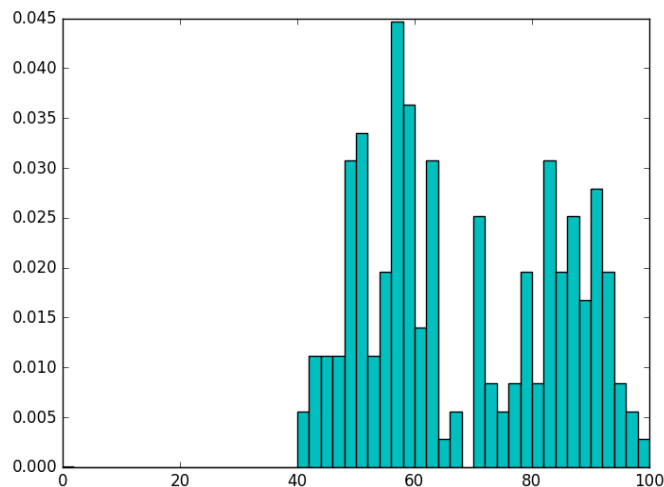
The goal is often to learn the distribution  $p(x)$  from which the data has been generated from



# Unsupervised learning

The goal is often to learn the distribution  $p(x)$  from which the data has been generated from

Example : find two groupes of patients following a memory test



# Unsupervised learning

The goal is often to learn the distribution  $p(x)$  from which the data has been generated from

## Other applications

- Data compression
- Data visualisation
- Population analysis
- Image segmentation
- etc.

# Supervised vs non-supervised



**Main topic of  
the school**

**Supervised learning** : there is a target

$$D = \{(\vec{x}_1, t_1), (\vec{x}_2, t_2), \dots, (\vec{x}_N, t_N)\}$$

**Unsupervised learning** : unknown target

$$D = \{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_N\}$$

# Supervised vs non-supervised

**Supervised learning** : there is a target

$$D = \{(\vec{x}_1, t_1), (\vec{x}_2, t_2), \dots\}$$

**Logistic regression**  
**Perceptron**  
**Multilayer perceptron**  
**Convolutional neural networks**  
**Recurrent neural networks**  
**Graph Neural Nets**  
**Etc.**

**Unsupervised learning** : unknown target

$$D = \{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_N\}$$

# Supervised vs non-supervised

**Supervised learning** : there is a target

$$D = \{(\vec{x}_1, t_1), (\vec{x}_2, t_2), \dots, (\vec{x}_N, t_N)\}$$

**Unsupervised learning** : unknown

**Autoencoders**  
**Variational autoencoders**  
**GANs**

$$D = \{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_N\}$$

# Back to supervised learning



# Supervised learning

## Two main applications

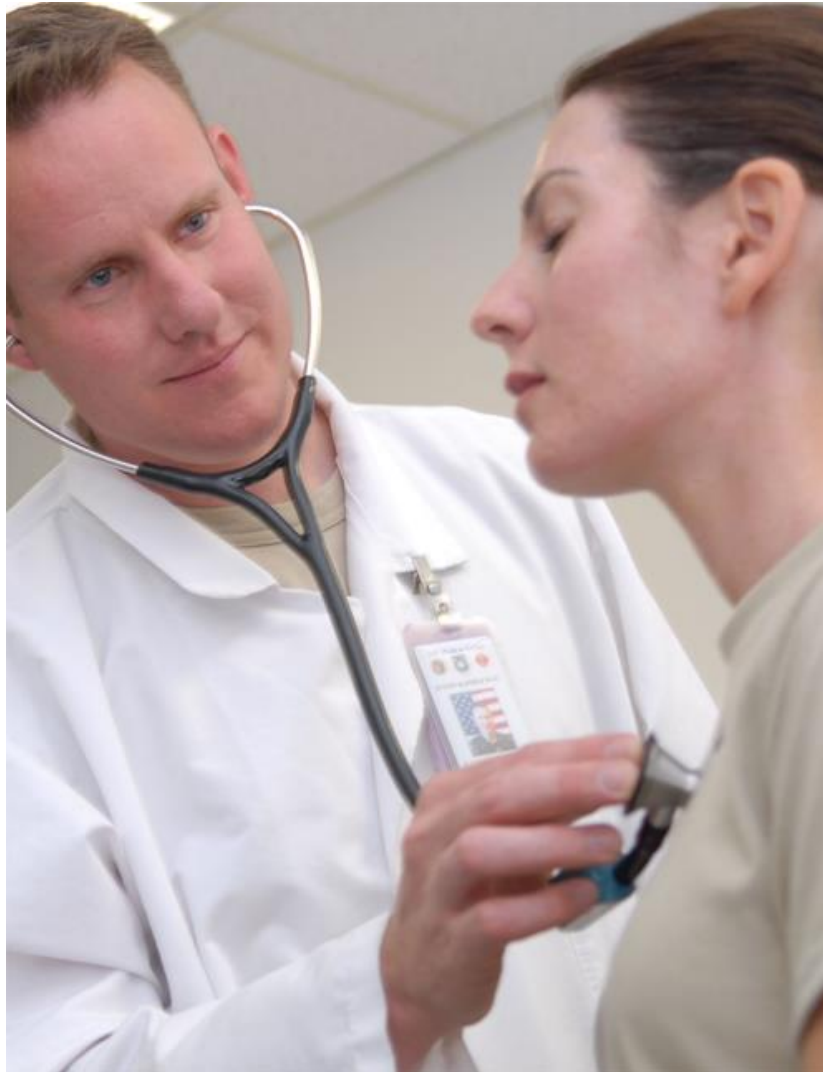
- **Classification** : the target is a class label  $t \in \{1, \dots, K\}$ 
  - Exemple : disease recognition
    - ✓  $\vec{x}$  : vector of medical measures, age, sex, etc.
    - ✓  $t$  : myocardial infarction, dilated cardiomyopathy, hypertrophic cardiomyopathy, normal
  
- **Régression** : the target is a real number  $t \in \mathbb{R}$ 
  - Exemple : prediction of life expectancy
    - ✓  $\vec{x}$  : vector of medical measures, age, sex, etc.
    - ✓  $t$  : number of months before death.

# Supervised learning

## Two main applications

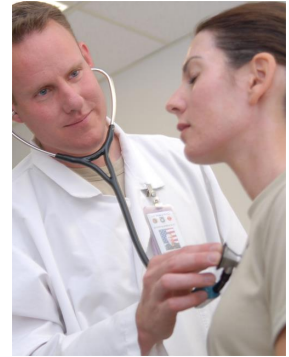
- **Classification** : the target is a class label  $t \in \{1, \dots, K\}$ 
  - Exemple : disease recognition
    - ✓  $\vec{x}$  : vector of medical measures, age, sex, etc.
    - ✓  $t$  : myocardial infarction, dilated cardiomyopathy, hypertrophic cardiomyopathy, normal
  
- **Régression** : the target is a real number  $t \in \mathbb{R}$ 
  - Exemple : prediction of life expectancy
    - ✓  $\vec{x}$  : vector of medical measures, age, sex, etc.
    - ✓  $t$  : number of months before death.

# Simple example of binary classification

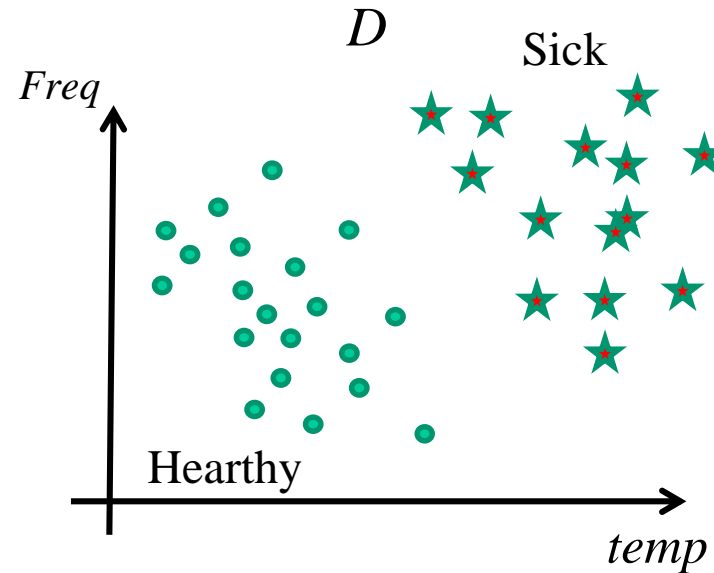


From Wikimedia Commons  
the free media repository

# Simple example of binary classification



$D$		
	( temp, freq)	Diagnostic
Patient 1	(37.5, 72)	heathy
Patient 2	(39.1, 103)	sick
Patient 3	(38.3, 100)	sick
	(...)	...
Patient N	(36.7, 88)	heathy
	$\bar{x}$	$t$

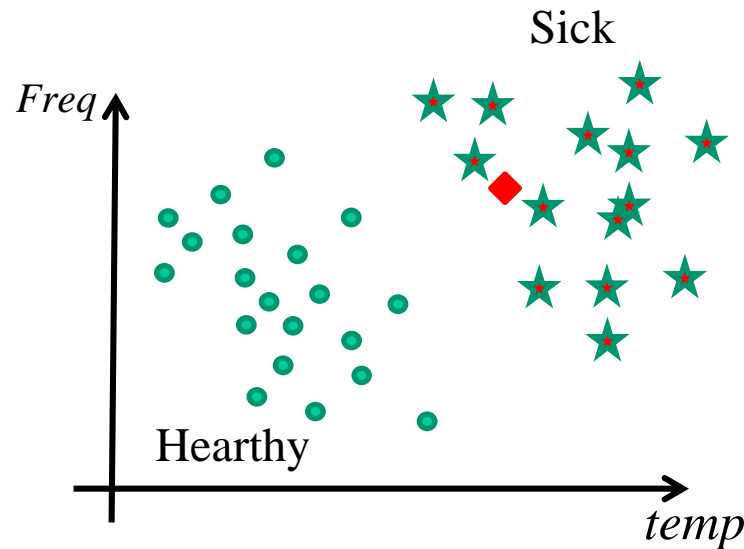


# Simple example of binary classification

A new patient shows up at the hospital  
**How can we predict its state?**



From Wikimedia Commons  
the free media repository

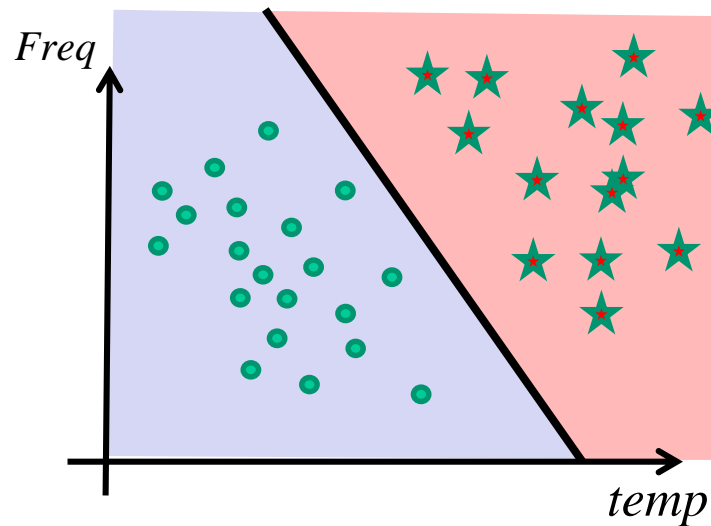


# Solution



From Wikimedia Commons  
the free media repository

Divide the feature space in two regions : **healthy** and **sick**

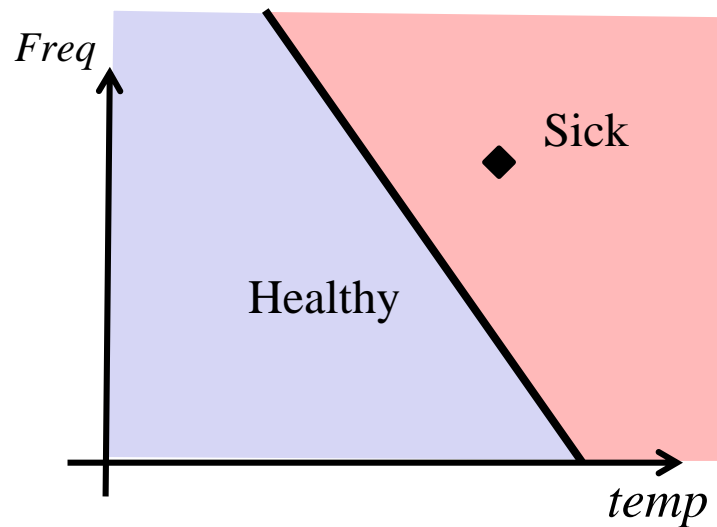


# Solution



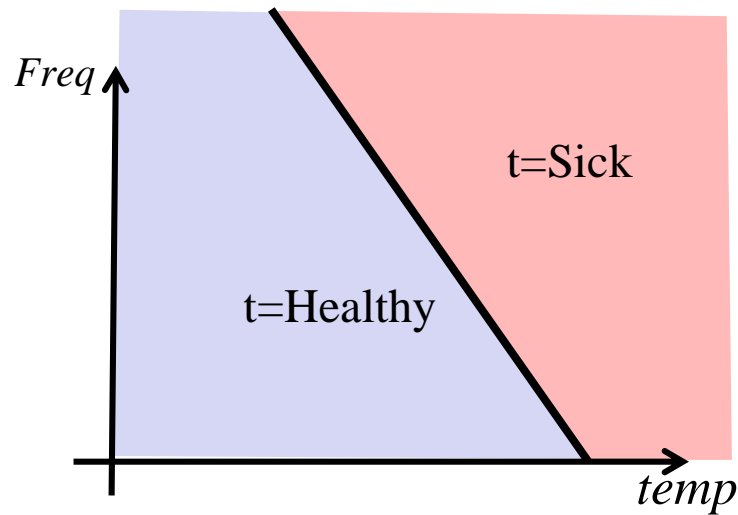
From Wikimedia Commons  
the free media repository

Divide the feature space in two regions : **healthy** and **sick**



# More formally

$$y_{\vec{w}}(\vec{x}) = \begin{cases} \textit{Healthy} & \text{if } \vec{x} \text{ is in the blue region} \\ \textit{Sick} & \text{otherwise} \end{cases}$$

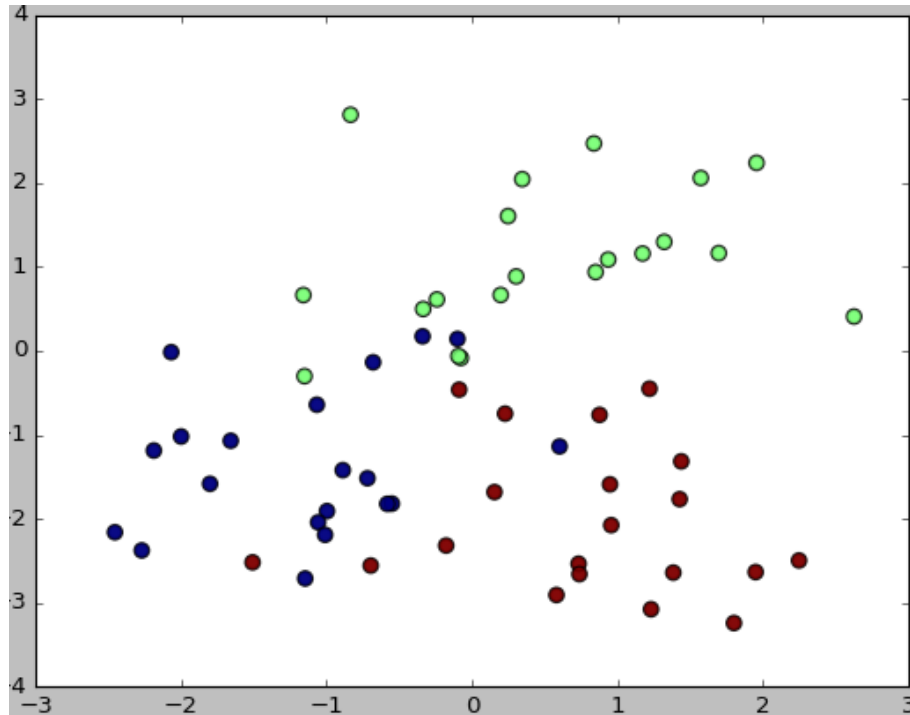


From Wikimedia Commons  
the free media repository

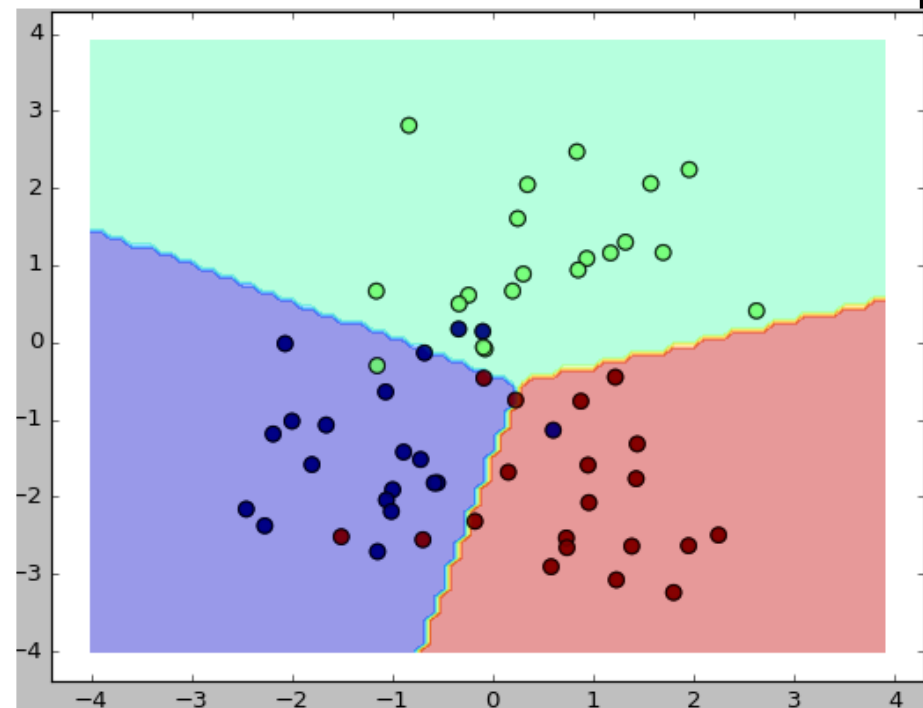


# Classification

## 3-class case



3 classes ●, ●, ● in a 2D feature space



Once training is over

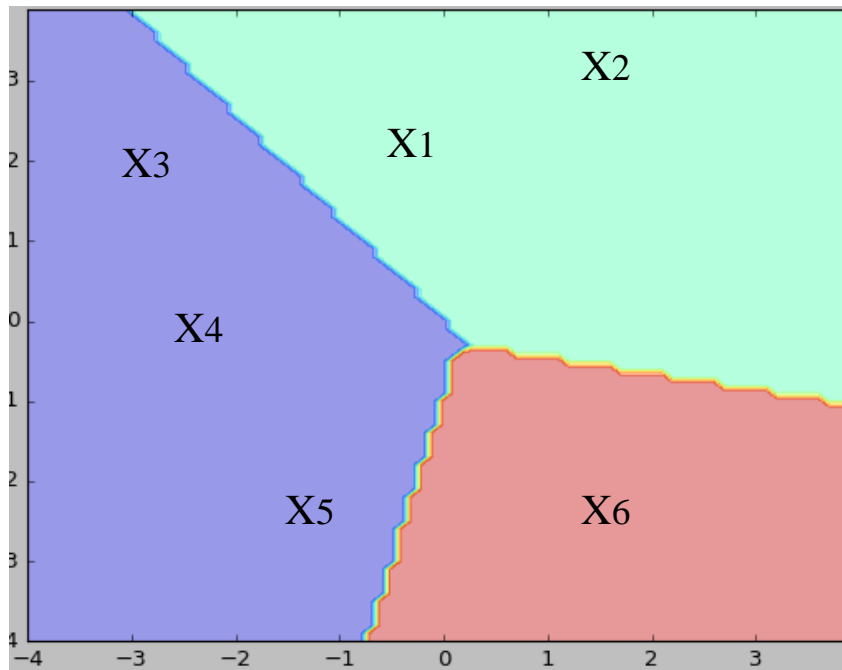
$y(\text{blue dot}) = \text{class 1}$

$y(\text{red dot}) = \text{class 2}$

$y(\text{green dot}) = \text{class 3}$

# Classification

Once training is over, we have a function  $y(x)$  that convert a point  $x$  into a class label



$y(X1) \Rightarrow$  class ■

$y(X2) \Rightarrow$  class ■

$y(X3) \Rightarrow$  class ■

$y(X4) \Rightarrow$  class ■

$y(X5) \Rightarrow$  class ■

$y(X6) \Rightarrow$  class ■

# MNIST

# Example of a classification dataset

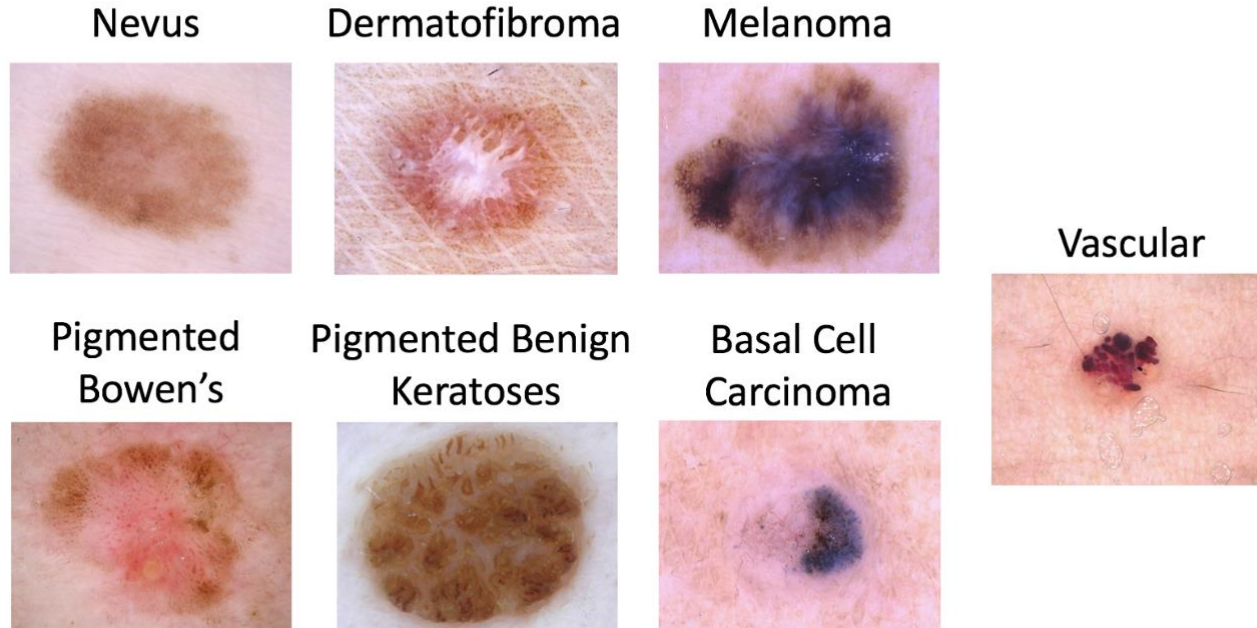
MNIST

- 10 classes
- 70,000 images
  - => 60,000 training
  - => 10,000 test
- Images are in grayscale
  - => 28x28

We can **vectorize these images** and represent it by a vector of size  $28 \times 28 = 784$  **dimensions**.

# Example of a medical classification dataset

ISIC melanoma classification challenge.



<https://challenge.isic-archive.com/landing/2018/>

# Example of a medical classification dataset

ISIC melanoma classification challenge dataset

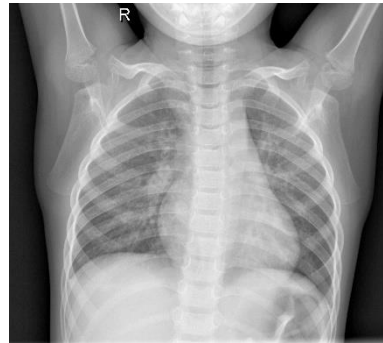
- 7 classes
- 11,527 images,
  - => 10,015 training
  - => 1,512 test
- Each image is in RGB
  - => 628x417x3

We can **vectorize these images** and represent it by a vector of size  $628 \times 417 \times 3 = 785,628$  **dimensions**.

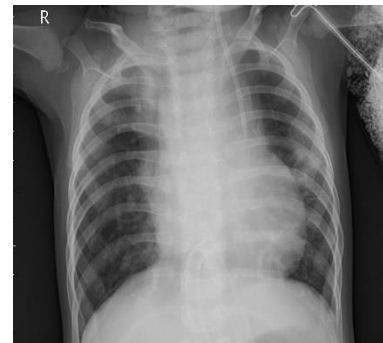
# Example of a medical classification dataset

*Chest X-Ray Pneumonia*

Healthy



Pneumonia



<https://www.kaggle.com/datasets/paultimothymooney/chest-xray-pneumonia>

# Example of a medical classification dataset

*Chest X-Ray Pneumonia*

- 2 classes
- 5,840 images,
  - => 5,216 training
  - => 624 test
- Each image is in grayscale
  - => 336 x 264\*

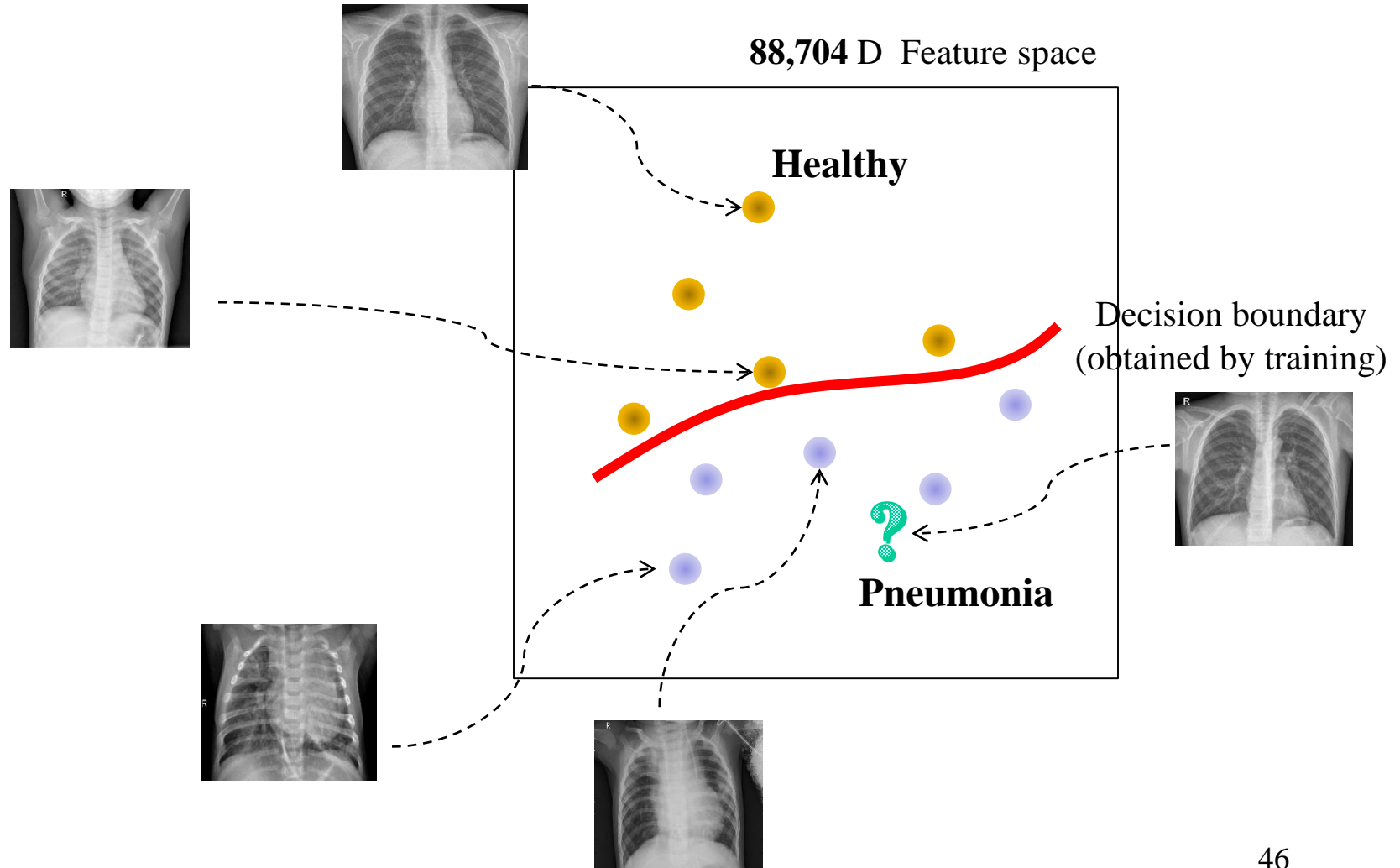
We can **vectorize these images** and represent it by a vector of size  $336 \times 264 = \mathbf{88,704}$  dimensions.

\* Rescaled version



# Supervised learning

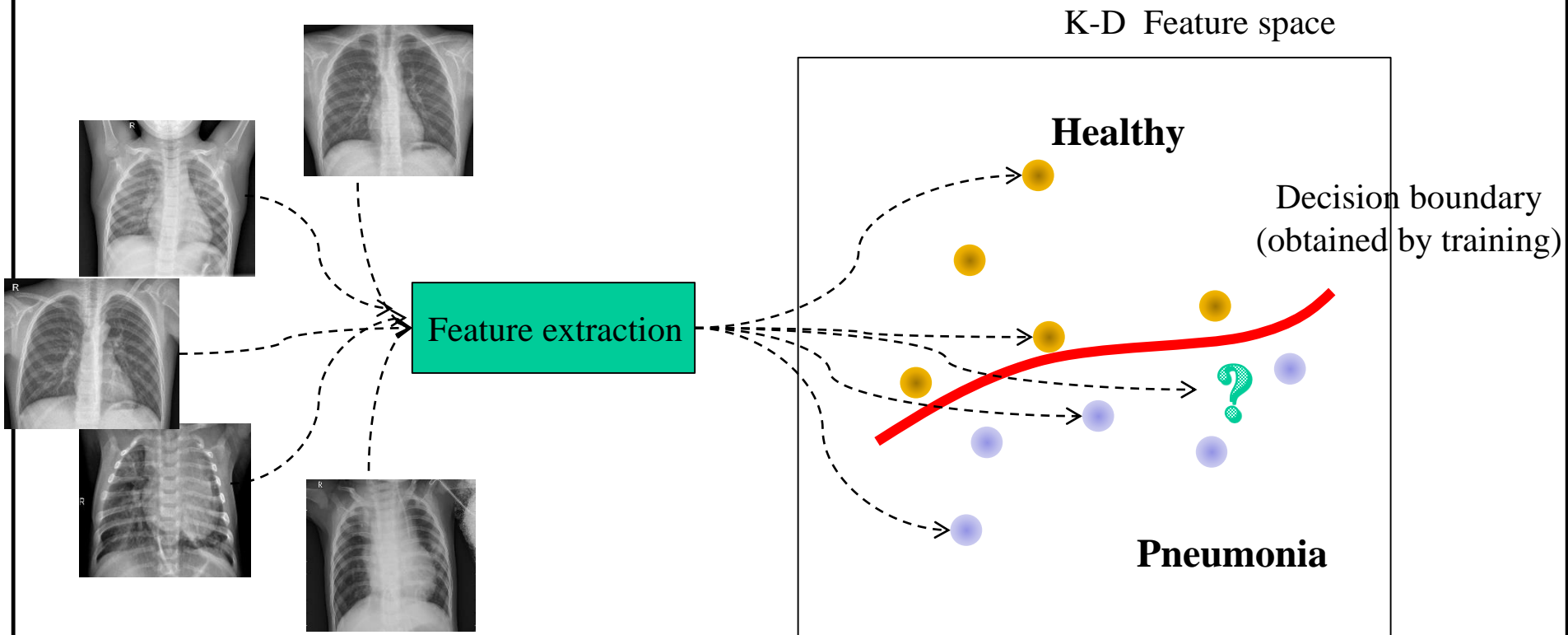
*Chest X-Ray Pneumonia*





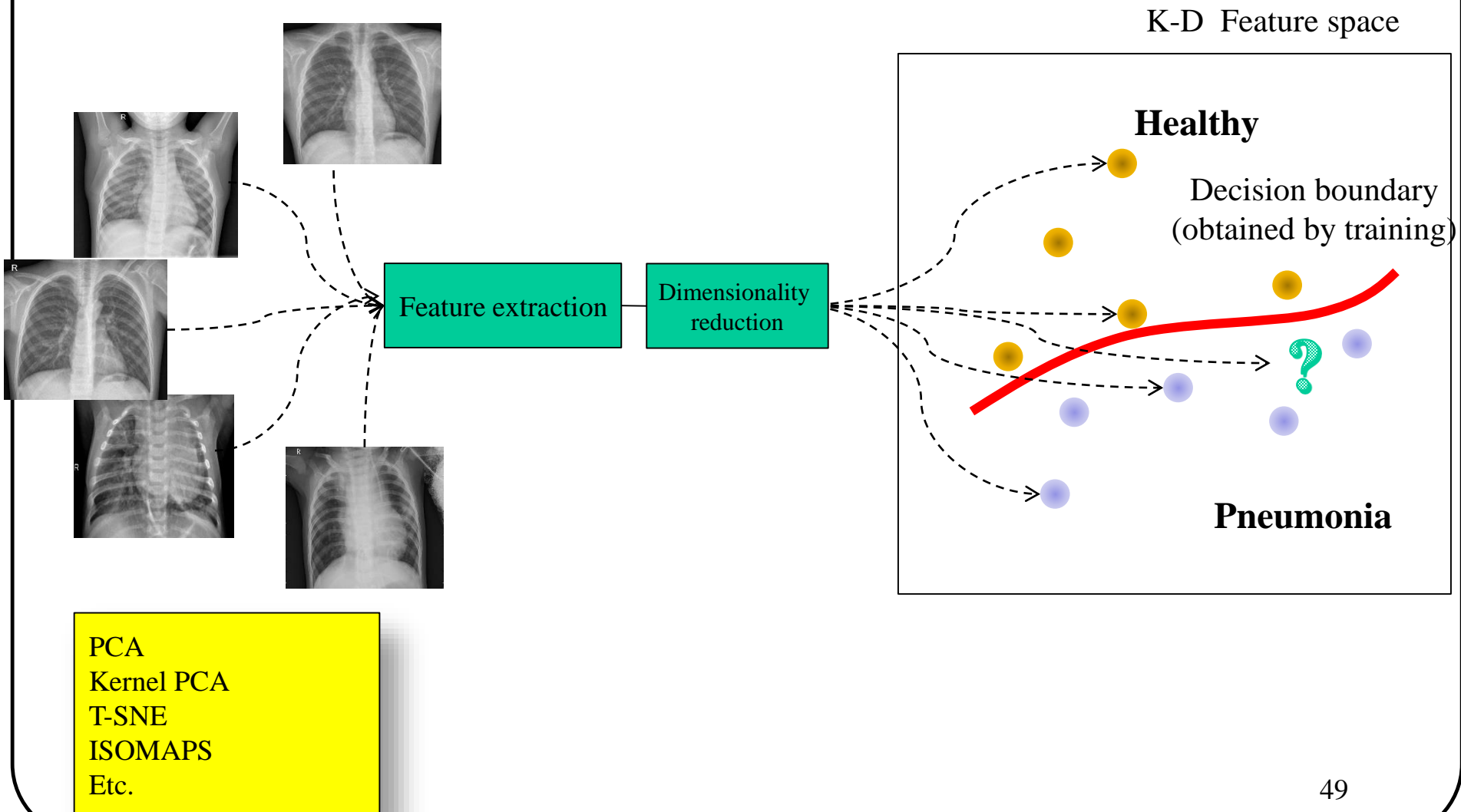
Very large feature spaces (like **88,704 dim**) are problematic.

# Supervised learning



Edge detection  
Histograms  
Filters of all kinds  
Moments  
Etc.

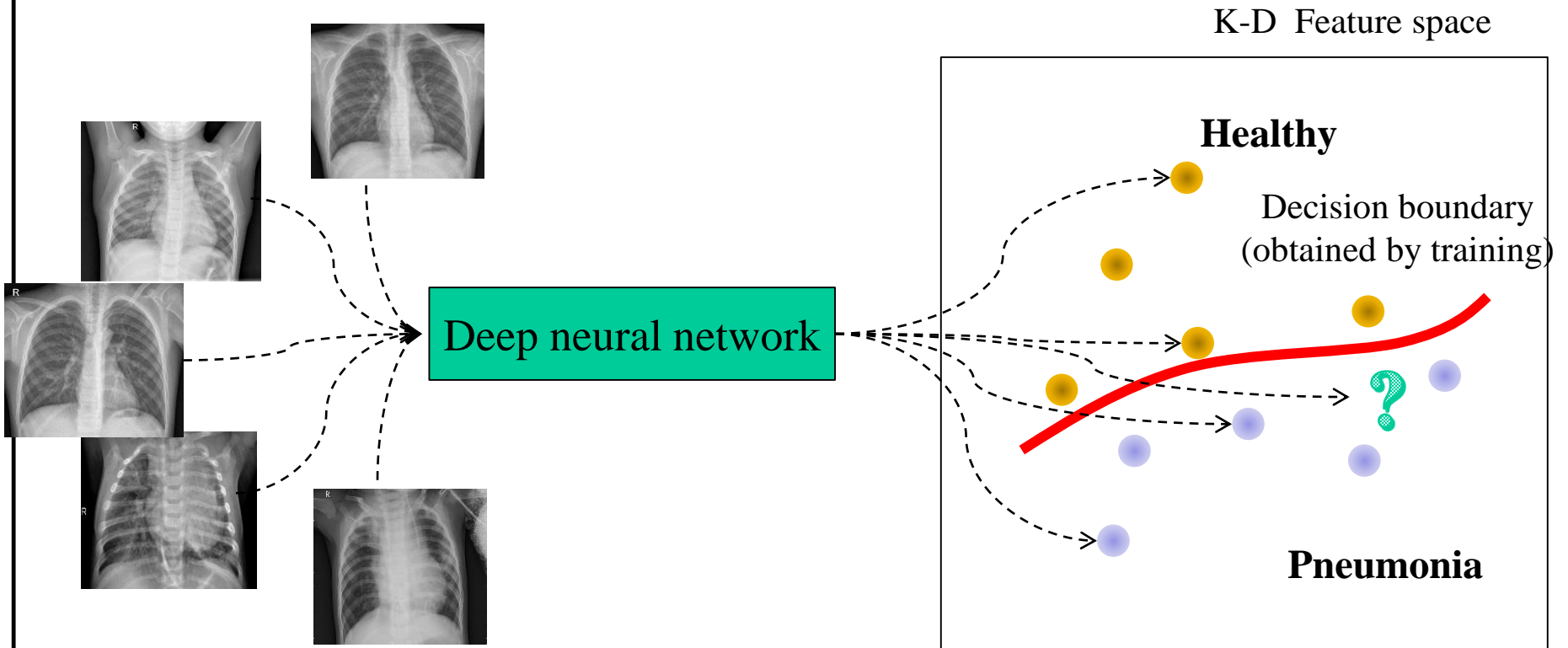
# Supervised learning





Spoiler alert

# In 2022...

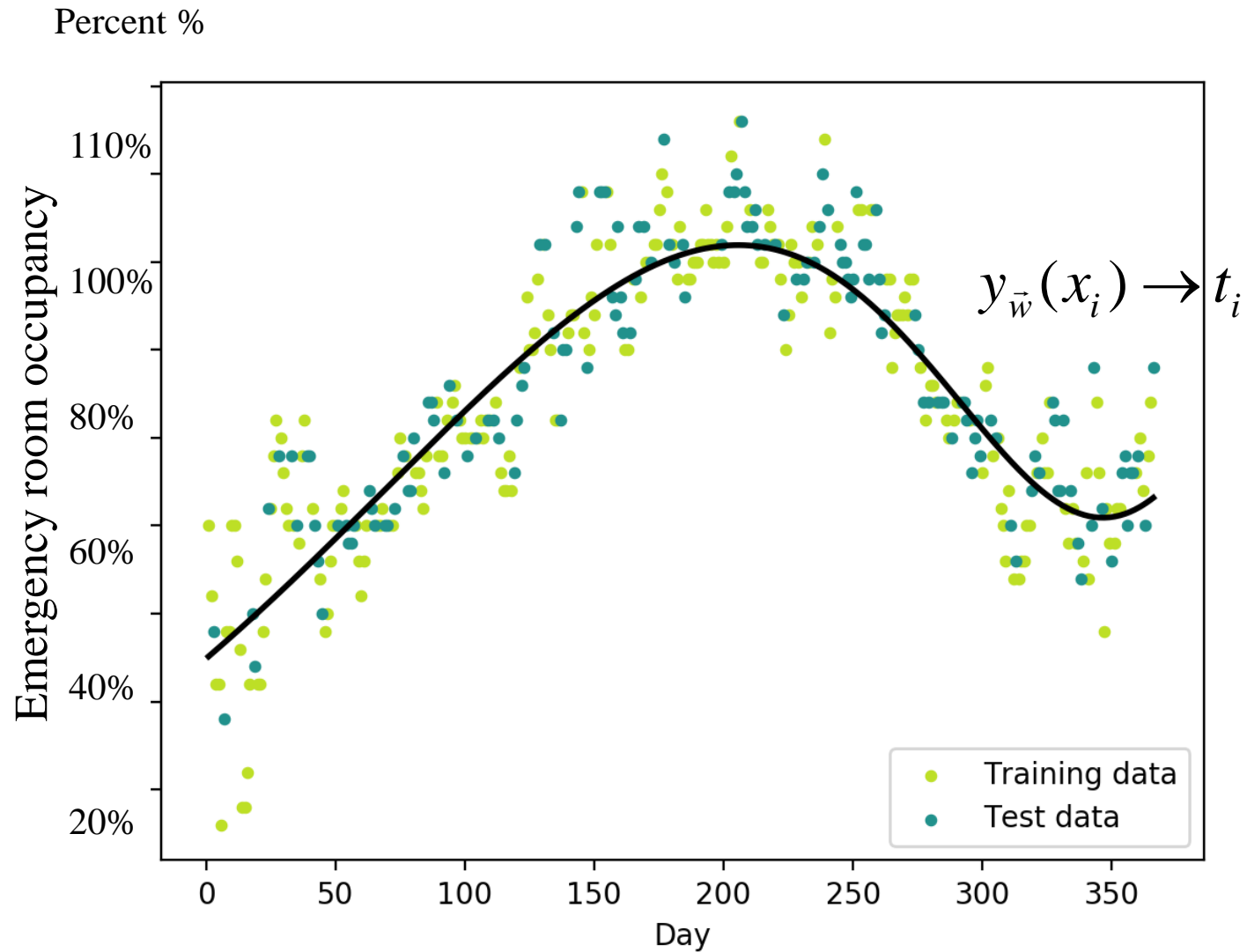


# Supervised learning

## Two main applications

- **Classification** : the target is a class label  $t \in \{1, \dots, K\}$ 
  - Exemple : disease recognition
    - ✓  $\vec{x}$  : vector of medical measures, age, sex, etc.
    - ✓  $t$  : myocardial infarction, dilated cardiomyopathy, hypertrophic cardiomyopathy, normal
  
- **Régression** : the target is a real number  $t \in \mathbb{R}$ 
  - Exemple : prediction of life expectancy
    - ✓  $\vec{x}$  : vector of medical measures, age, sex, etc.
    - ✓  $t$  : number of months before death.

# Example





Formal example:  
**1D regression**

# 1D Regression

## Simple example

### ➤ Data

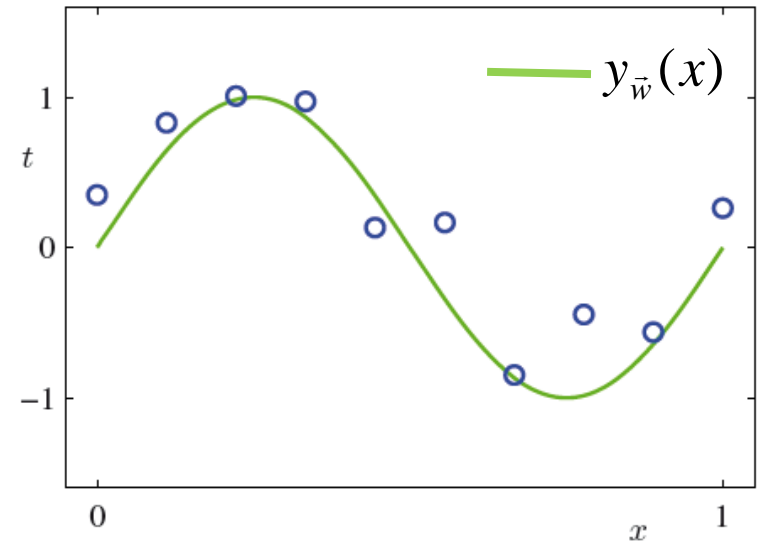
- ✓ input : scalar  $x \in [0,1]$
- ✓ target : scalar  $t \in [-1,1]$

### ➤ Training dataset $D$ contains:

- ✓  $X = (x_1, \dots, x_N)^T$
- ✓  $T = (t_1, \dots, t_N)^T$

### ➤ Goal :

- ✓ Get a prediction  $\hat{t}$  for each new data  $\hat{x}$



# 1D Regression

## Simple example

### ➤ Data

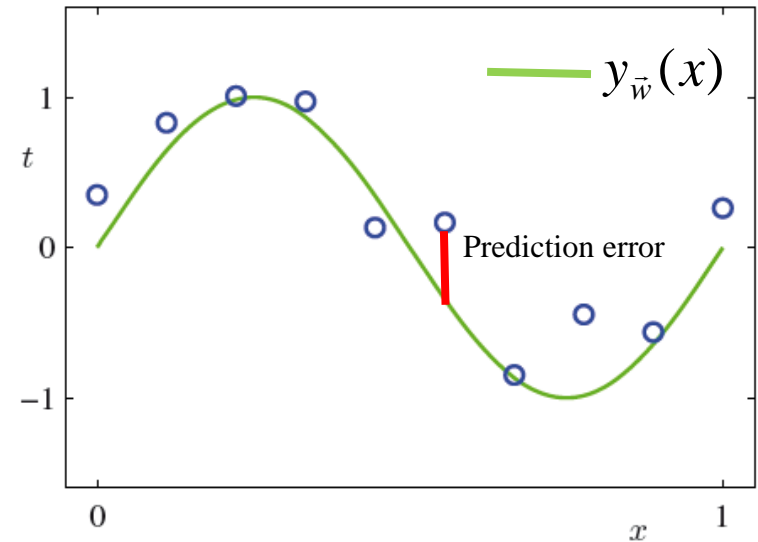
- ✓ input : scalar  $x \in [0,1]$
- ✓ target : scalar  $t \in [-1,1]$

### ➤ Training dataset $D$ contains:

- ✓  $X = (x_1, \dots, x_N)^T$
- ✓  $T = (t_1, \dots, t_N)^T$

### ➤ Goal :

- ✓ Get a prediction  $\hat{t}$  for each new data  $\hat{x}$



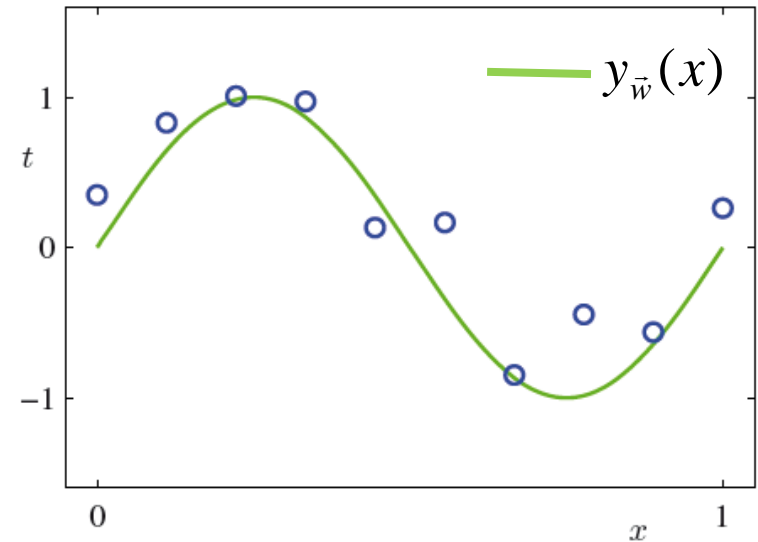
# Polynomial shape

- Let's assume that our data has a **polynomial shape**

$$\begin{aligned}y_{\vec{w}}(x) &= w_0 + w_1x + w_2x^2 + \dots + w_Mx^M \\ &= \sum_{i=0}^M w_i x^i\end{aligned}$$

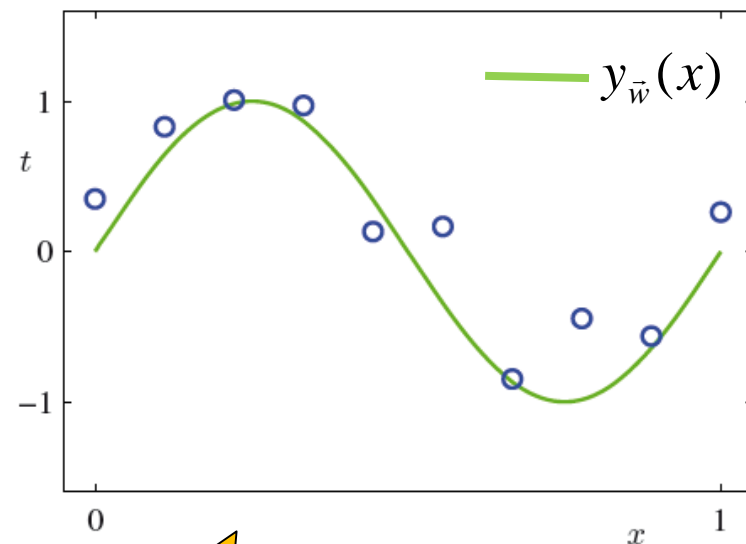
- $y_{\vec{w}}(x)$  is our **model**

- ✓ It implicitly includes our hypotheses on the problem at hand
- ✓ A model always has parameters that one needs to find (here  $\vec{w}$ )



# Unknowns

$$y_{\vec{w}}(x) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M$$
$$= \sum_{i=0}^M w_i x^i$$



Two unknowns

$$\vec{w} \in R^M$$

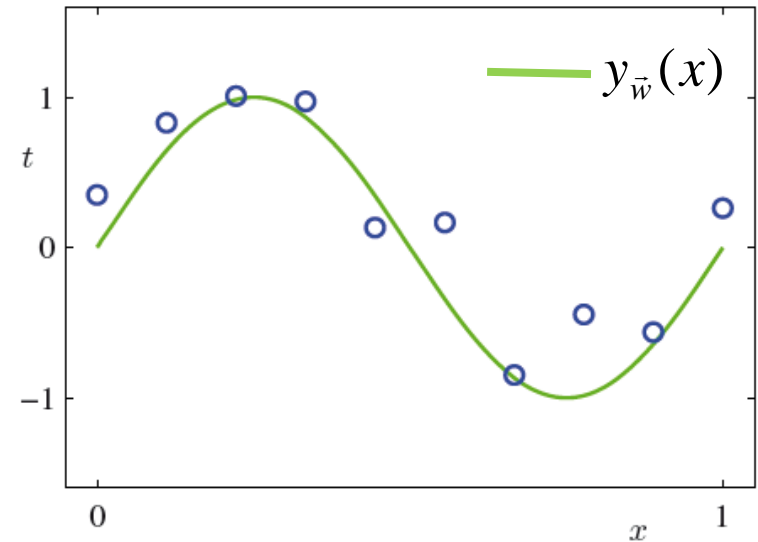
**Parameters**

$$M \in \mathbb{N}^{\geq 0}$$

**Hyper-parameters**

# Training

$$\begin{aligned}y_{\vec{w}}(x) &= w_0 + w_1x + w_2x^2 + \dots + w_Mx^M \\ &= \sum_{i=0}^M w_i x^i\end{aligned}$$



Two unknowns

$$\vec{w} \in R^M$$

$$M \in \mathbb{N}^{\geq 0}$$

**Training = find the « Best »  $w$   
(and sometimes  $M$ ) from  $X$  and  $T$**

# Polynomial regression

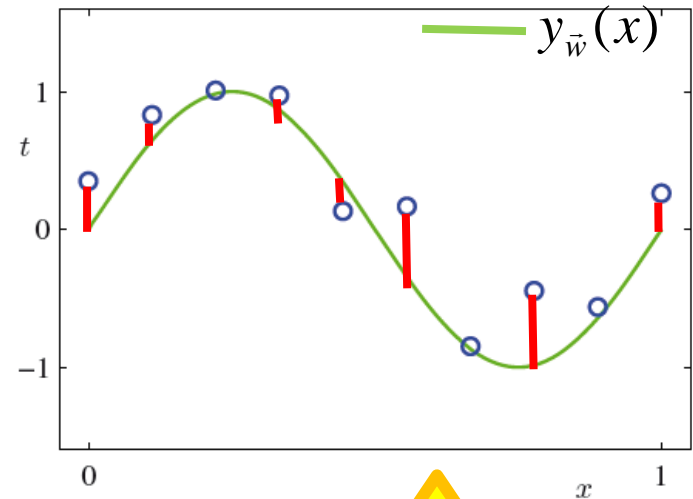
➤ If  $M$  is *given*, how can we find the best  $\vec{w}$ ?

The « best »  $\vec{w}$  is the one that minimizes the error (or the **loss**) on the training data

$$E_D(\vec{w}) = \frac{1}{N} \sum_{n=1}^N (y_{\vec{w}}(x_n) - t_n)^2$$

➤ More formally

$$\vec{w} = \arg \min_{\vec{w}} E_D(\vec{w})$$



mean squared error  
(MSE)

# Hyperparameters

➤ How to find the best  $M$ ?

Hyperparameters cannot be **estimated** with classical **optimization algorithms** (gradient descent, Newton method, Simplex method, etc.) as for  $\vec{w}$

**hyperparameters are often fixed manually.**

**BUT BEWARE!!!** Their value greatly influence the **final results**.

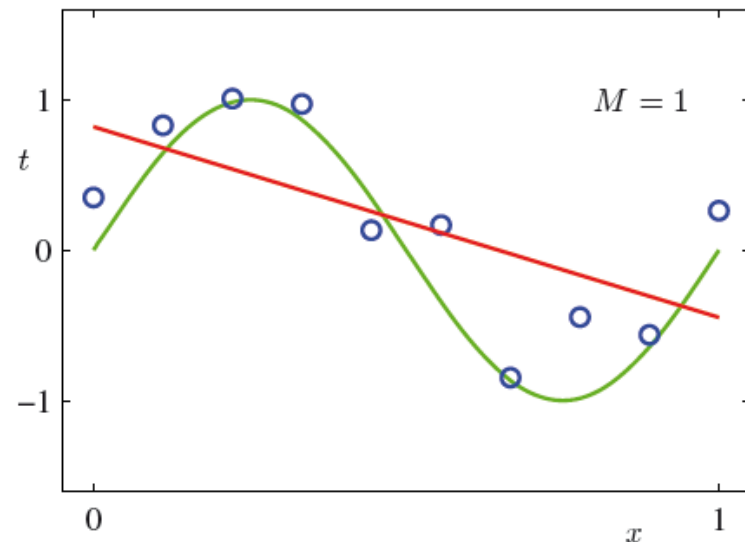
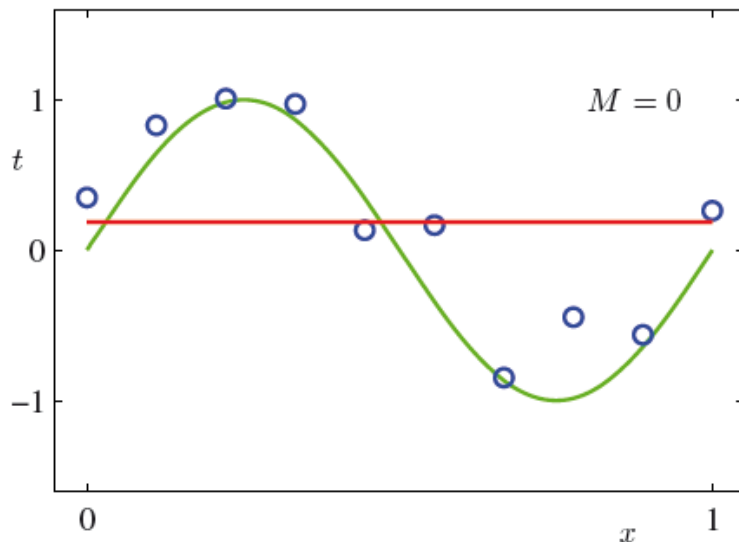


# Underfitting

$$M = 0 \Rightarrow y_{\vec{w}}(x) = w_0$$

$$M = 1 \Rightarrow y_{\vec{w}}(x) = w_0 + w_1 x$$

A small  $M$  gives a simplistic model that will **underfit** the data.

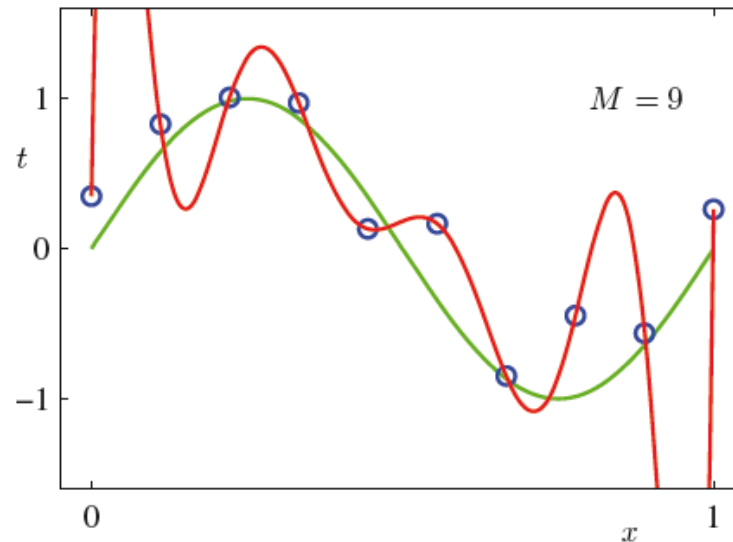


$$E_D(\vec{w}) \Rightarrow \text{large}$$

$$E_{D_{test}}(\vec{w}) \Rightarrow \text{large}$$

# Overfitting

A large  $M$  gives a model that « **learn by heart** » and thus **overfit** training data

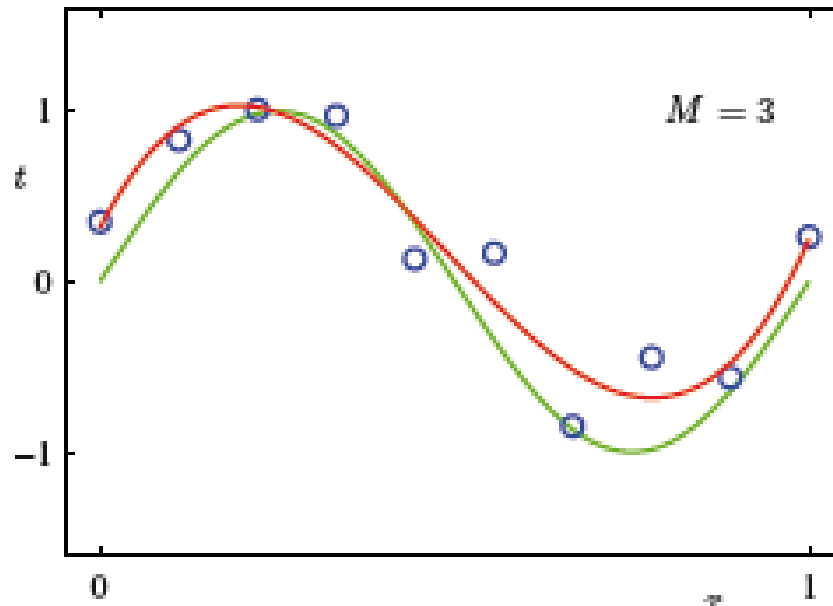


$$E_D(\vec{w}) \Rightarrow \text{VERY low}$$

$$E_{D_{test}}(\vec{w}) \Rightarrow \text{Large}$$

# Over- and underfitting

Need for an **intermediate value** for which the training and the testing errors are low



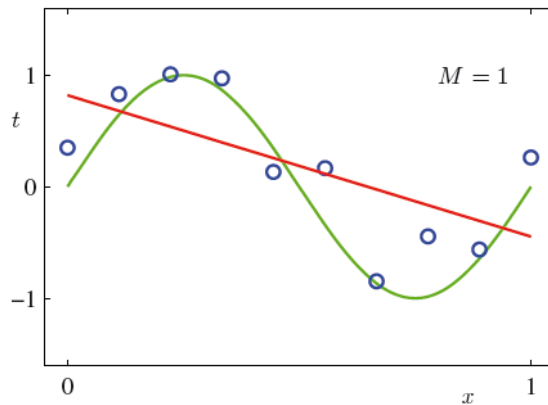
$$E_D(\vec{w}) \Rightarrow \text{low}$$

$$E_{D_{test}}(\vec{w}) \Rightarrow \text{low}$$

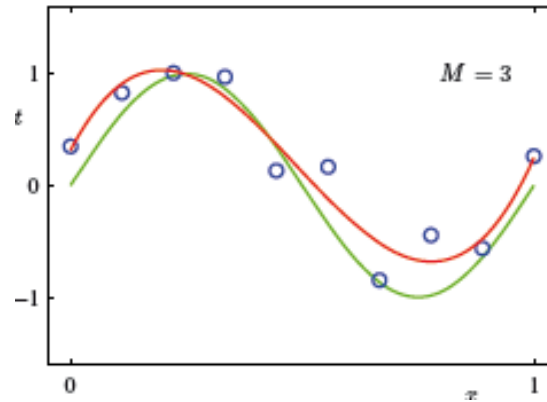
# Hyperparameters often control the **capacity** of a model

**Capacity:** ability of a model to fit the training data

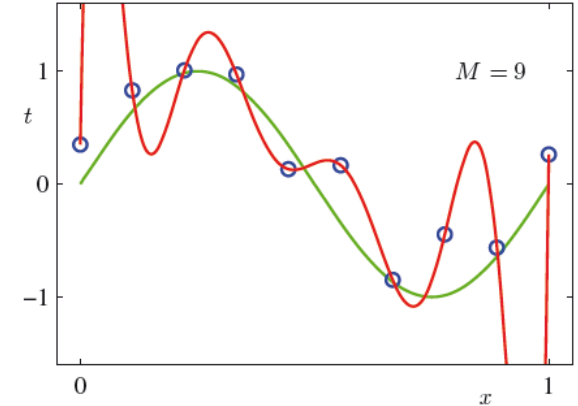
Low capacity



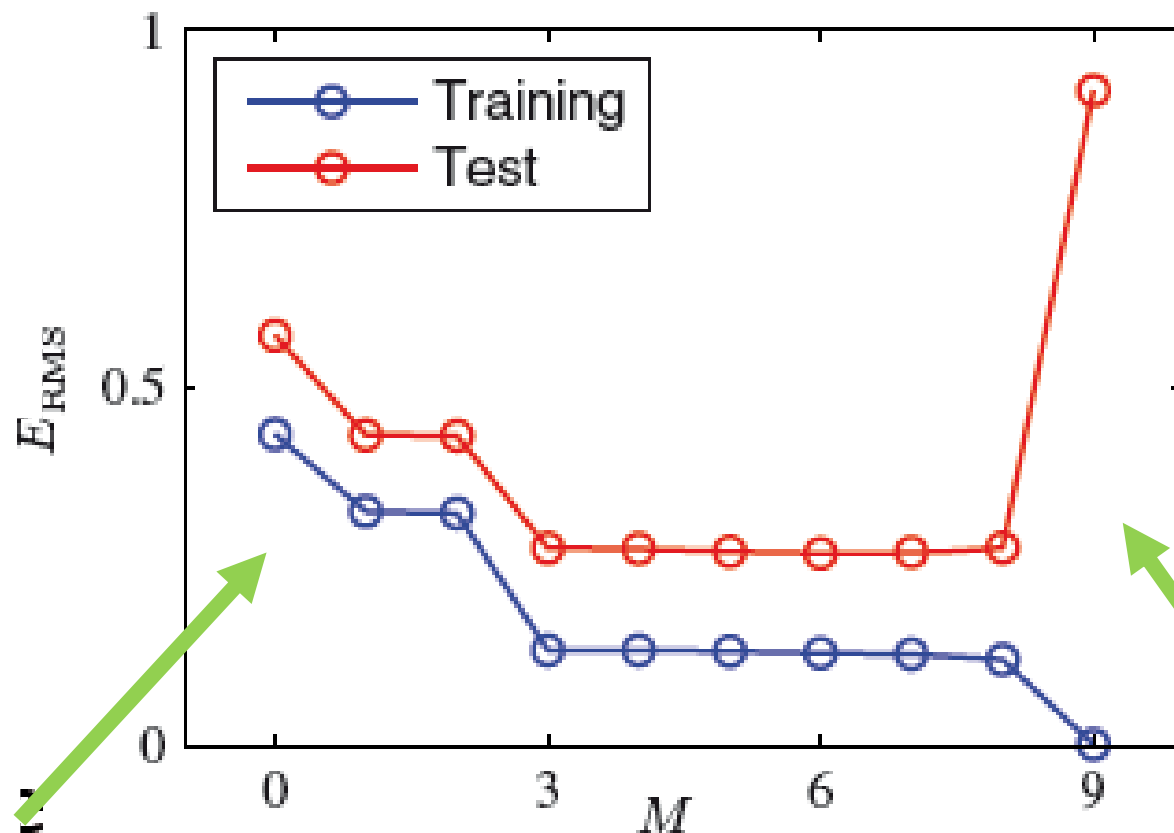
Medium capacity



Large capacity



# Over- and underfitting



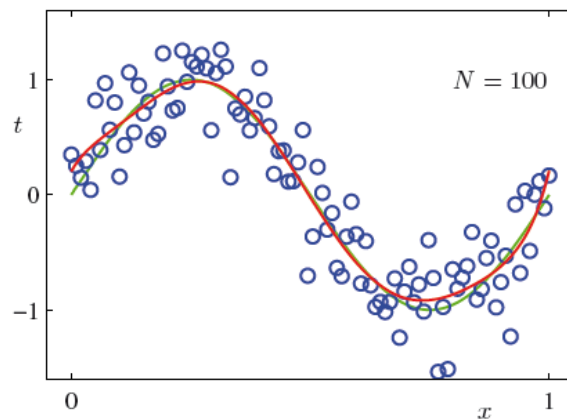
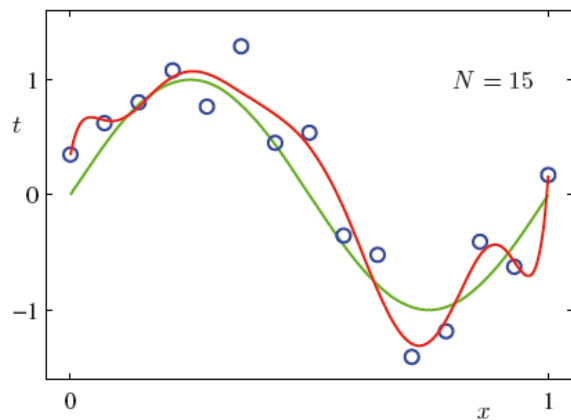
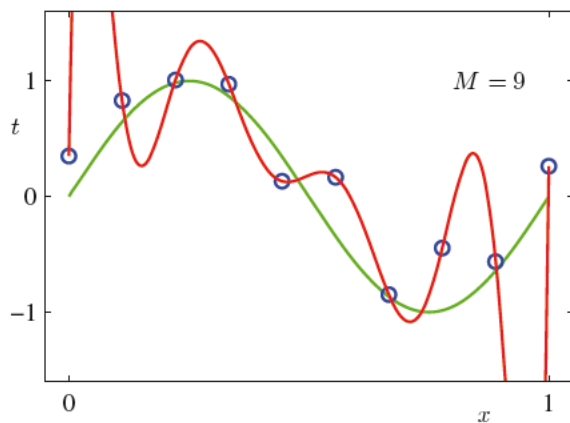
**Underfitting**

$$E_{RMS} = \sqrt{\frac{E(\vec{w})}{N}}$$

**Overfitting**

# Generalization

The more data you have, the better a high capacity model will generalize.



If we have a **large  $M$**  and **few training data**, how do we prevent our model from overfitting?



# Regularization

Parameter values  $\vec{w}$  for different  $M$  **without** regularization

	$M = 0$	$M = 1$	$M = 3$	$M = 9$
$w_0$	0.19	0.82	0.31	0.35
$w_1$		-1.27	7.99	232.37
$w_2$			-25.43	-5321.83
$w_3$			17.37	48568.31
$w_4$				-231639.30
$w_5$				640042.26
$w_6$				-1061800.52
$w_7$				1042400.18
$w_8$				-557682.99
$w_9$				125201.43



# Regularization

To prevent over-fitting

1. Choose a small « M »
2. Reduce capacity by **regularization**

Exemple : penalise the **L2 norm**

$$E_D(\vec{w}) = \frac{1}{N} \sum_{n=1}^N (t_n - y_{\vec{w}}(\vec{x}))^2 + \lambda \|\vec{w}\|^2$$

Constante that controls  
regularization

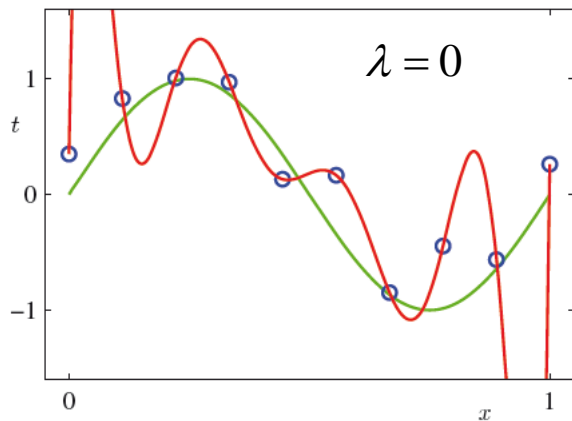
$$\|\vec{w}\|^2 = \vec{w}^T \vec{w} = w_0^2 + w_1^2 + \dots + w_M^2$$

Ridge model

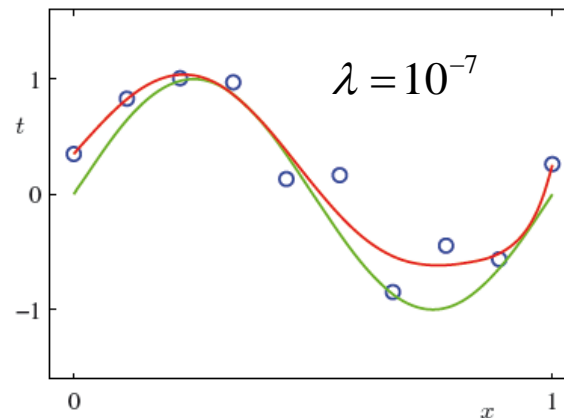
# Regularization

Strong regularization= less capacity

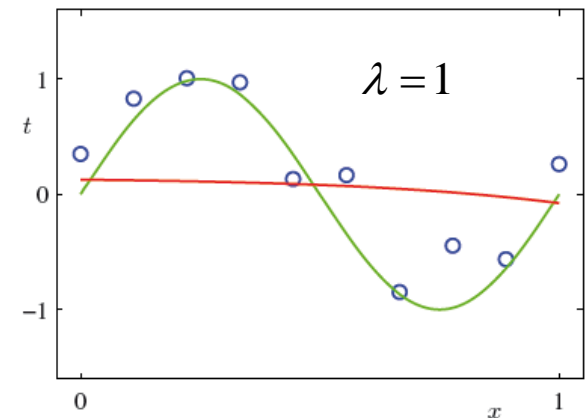
M=9



M=9

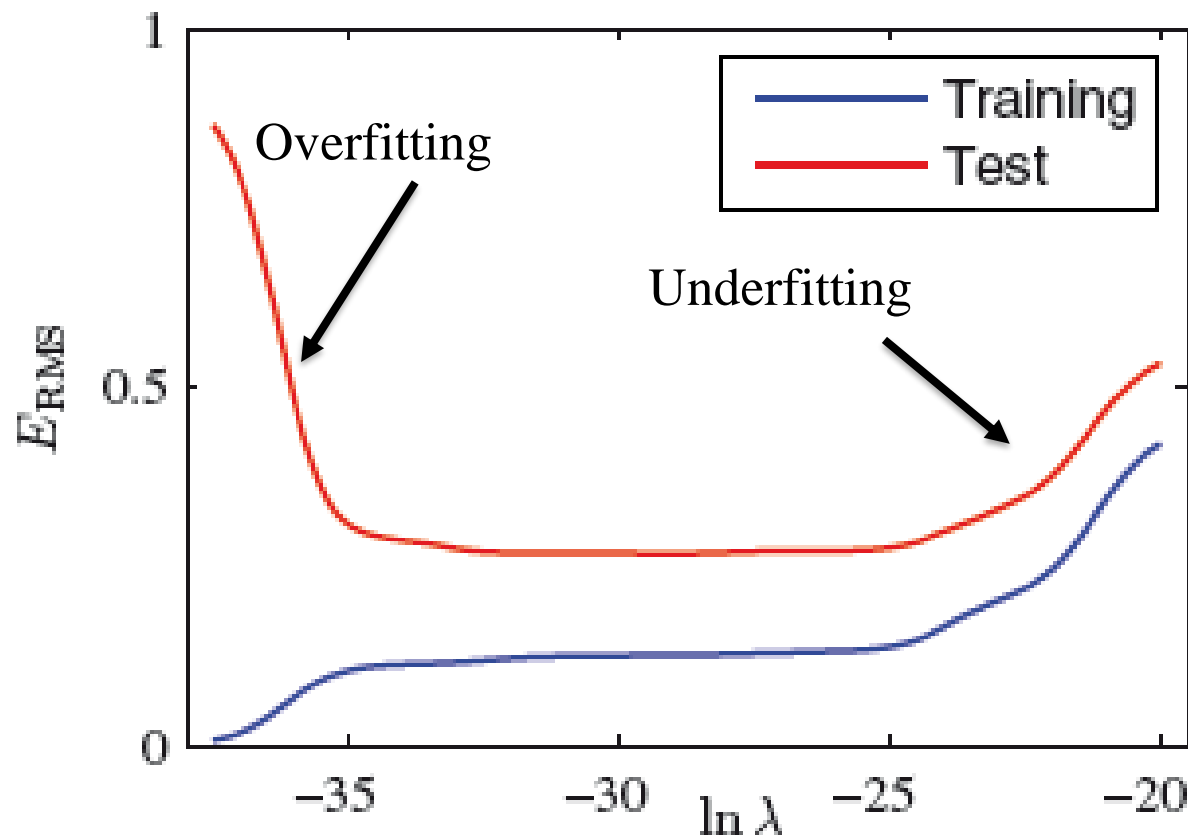


M=9



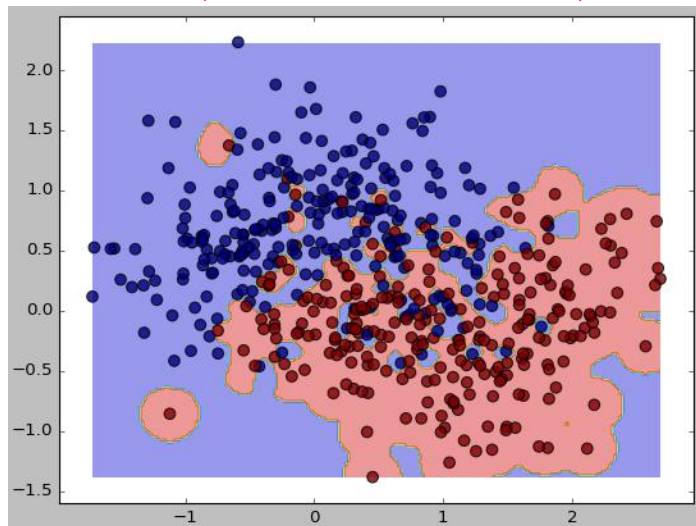
# Regularization

Influence on the training and testing error

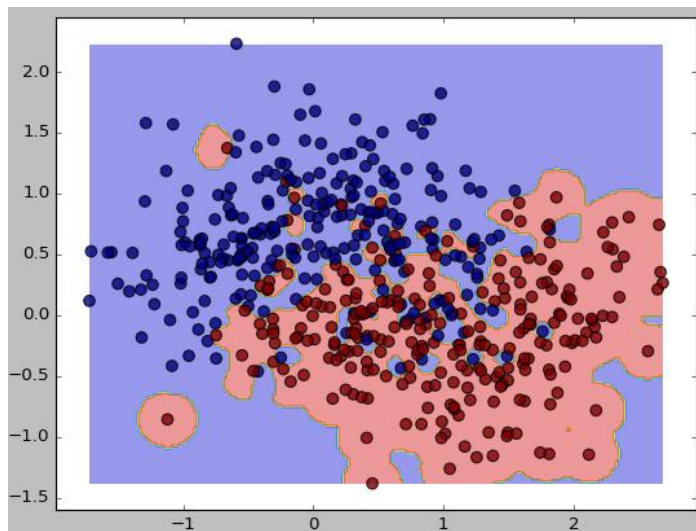


Over- and under-fitting  
also influence classification

## Overfitting (Classification)

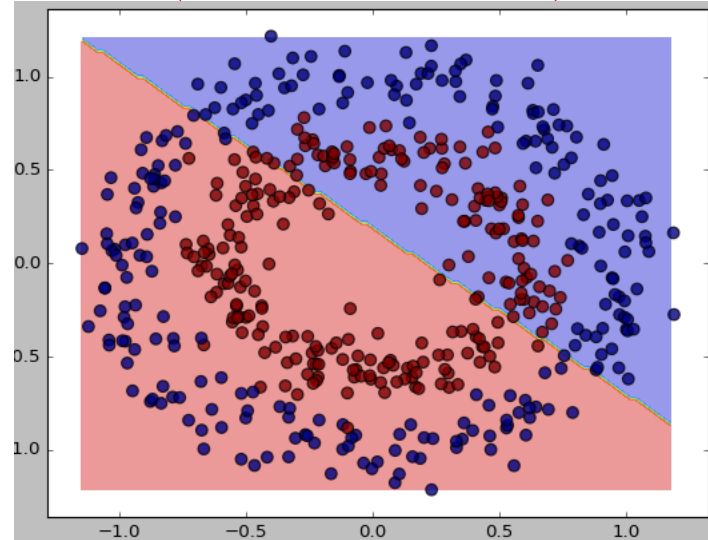


Training accuracy = 99.6%

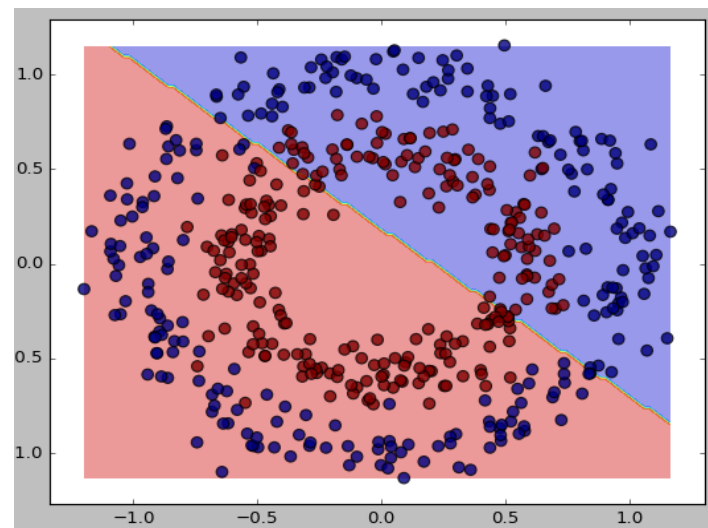


Testing accuracy = 78%

## Underfitting (Classification)

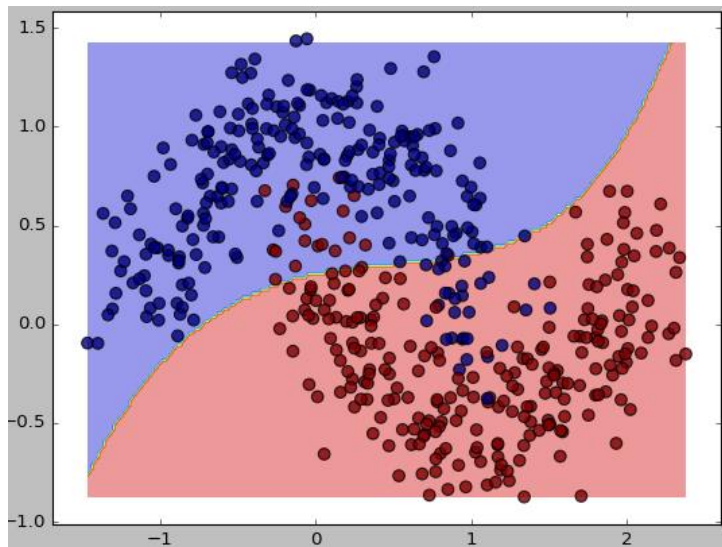


Training accuracy = 52.2%

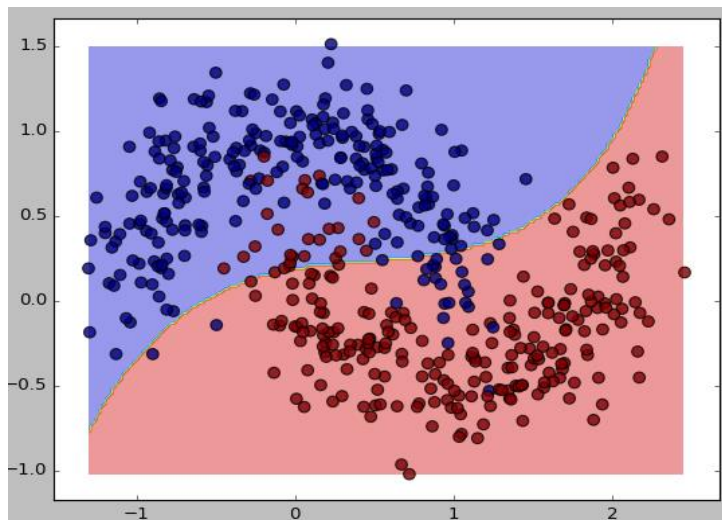


Testing accuracy = 51.2%

Could be better...

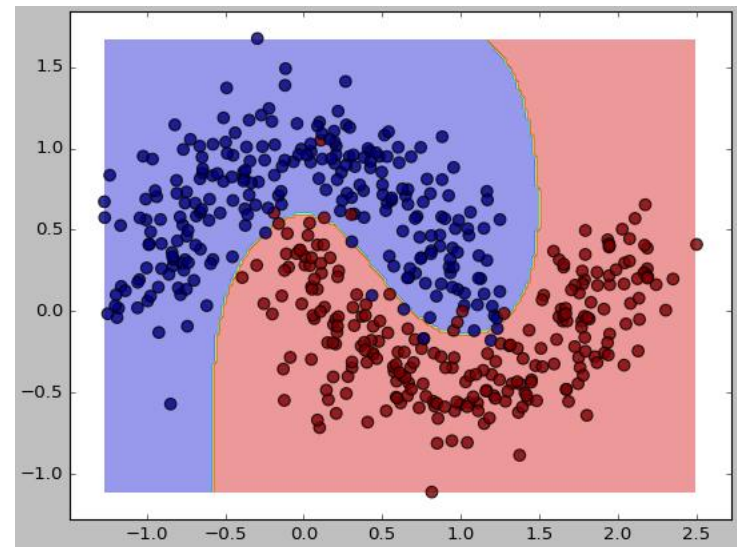


Training accuracy = 82%

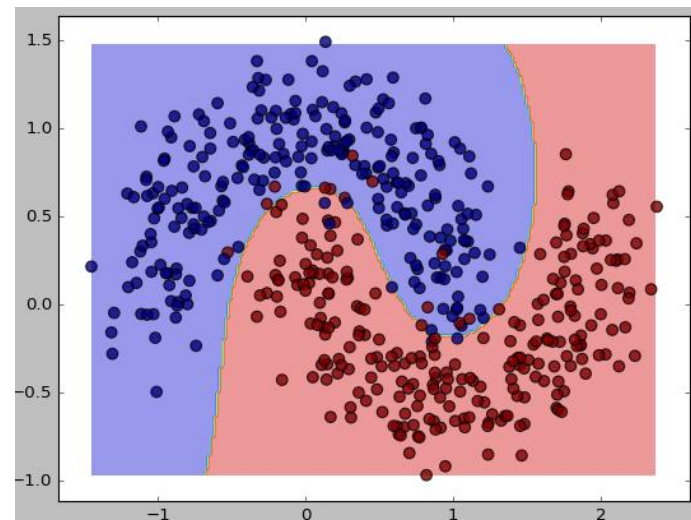


Testing accuracy = 80%

Wonderful !!!



Training accuracy = 97.8%



Testing accuracy = 96.2%

$$E_D(\vec{w}) = \sum_{n=1}^N (y_{\vec{w}}(x_n) - t_n)^2 + \lambda \|\vec{w}\|^2$$

$$\|\vec{w}\|^2 = \vec{w}^T \vec{w} = w_0^2 + w_1^2 + \dots + w_M^2$$

# Model selection

How to find the right hyper-parameters?

$M$  and  $\lambda$

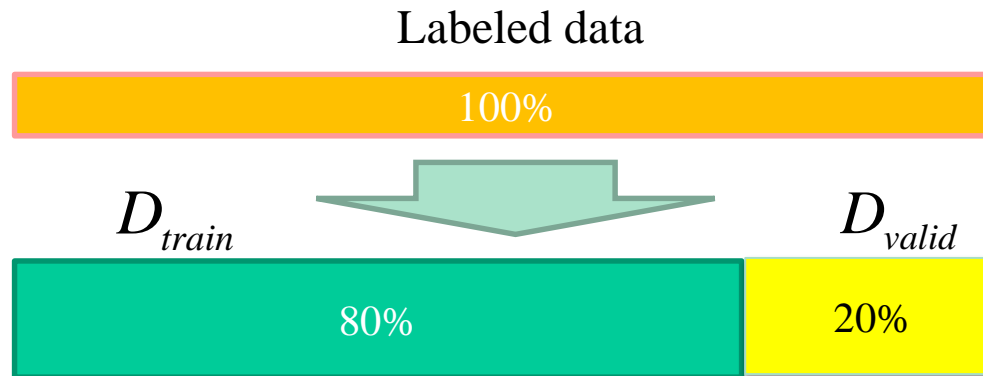
# How to find the right $M$ and the right $\lambda$ ?

- **Very bad idea** : choose randomly
- **Bad idea** : take many  $(M, \lambda)$  and keep the one with the lowest training error
  - overfitting
- **Bad idea** : take many  $(M, \lambda)$  and keep the one with the lowest test error
  - $D_{test}$  should NEVER be used to train a model
- **Good solution** : take many  $(M, \lambda)$  and keep the one with the lower **validation error**



# Cross-validation

1- Randomly divide data in 2 groups



2- FOR  $M$  from  $M_{\min}$  to  $M_{\max}$   
FPR  $\lambda$  from  $\lambda_{\min}$  to  $\lambda_{\max}$

Train the model on  $D_{train}$   
Compute error on  $D_{valid}$

3- Keep  $(M, \lambda)$  with the lowest **validation error**

# *k-fold cross-validation*

```
FOR  $M$  from  $M_{\min}$  to  $M_{\max}$   
  FOR  $\lambda$  from  $\lambda_{\min}$  to  $\lambda_{\max}$   
    FOR  $j$  from 0 to  $K$ 
```

Divide the labeled data in 2 groups  $\Rightarrow D_{train} D_{valid}$

Train the model on  $D_{train}$   
Compute error on  $D_{valid}$

Keep  $(M, \lambda)$  with the lowest **mean validation error**

# K-fold cross-validation with $K = 10$

Mean validation error

STD

2.832	(+/-0.116)	for {'regression': 'poly', 'M': 3, 'lambda': 0.01}
1.854	(+/-0.072)	for {'regression': 'poly', 'M': 3, 'lambda': 0.1}
1.910	(+/-0.065)	for {'regression': 'poly', 'M': 3, 'lambda': 1}
1.902	(+/-0.077)	for {'regression': 'poly', 'M': 3, 'lambda': 10}
2.844	(+/-0.101)	for {'regression': 'poly', 'M': 4, 'lambda': 0.01}
2.864	(+/-0.089)	for {'regression': 'poly', 'M': 4, 'lambda': 0.1}
1.910	(+/-0.065)	for {'regression': 'poly', 'M': 4, 'lambda': 1}
1.894	(+/-0.086)	for {'regression': 'poly', 'M': 4, 'lambda': 10}
2.848	(+/-0.080)	for {'regression': 'poly', 'M': 5, 'lambda': 0.01}
1.904	(+/-0.064)	for {'regression': 'poly', 'M': 5, 'lambda': 0.1}
0.916	(+/-0.069)	for {'regression': 'poly', 'M': 5, 'lambda': 1}
1.870	(+/-0.072)	for {'regression': 'poly', 'M': 5, 'lambda': 10}
2.846	(+/-0.090)	for {'regression': 'poly', 'M': 6, 'lambda': 0.01}
2.906	(+/-0.062)	for {'regression': 'poly', 'M': 6, 'lambda': 0.1}
1.904	(+/-0.075)	for {'regression': 'poly', 'M': 6, 'lambda': 1}
2.858	(+/-0.112)	for {'regression': 'poly', 'M': 6, 'lambda': 10}

**BEST!**

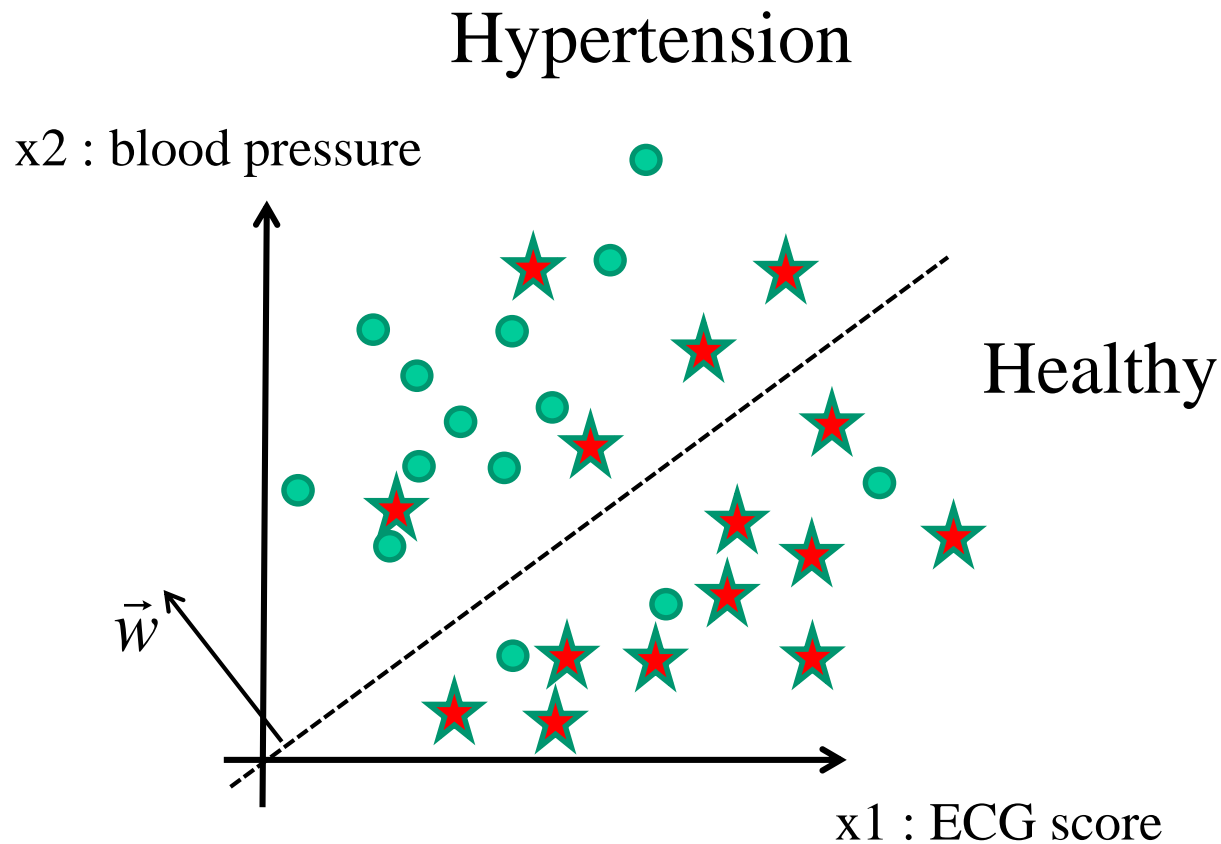
**M=5,  
 $\lambda = 1$**

# In short

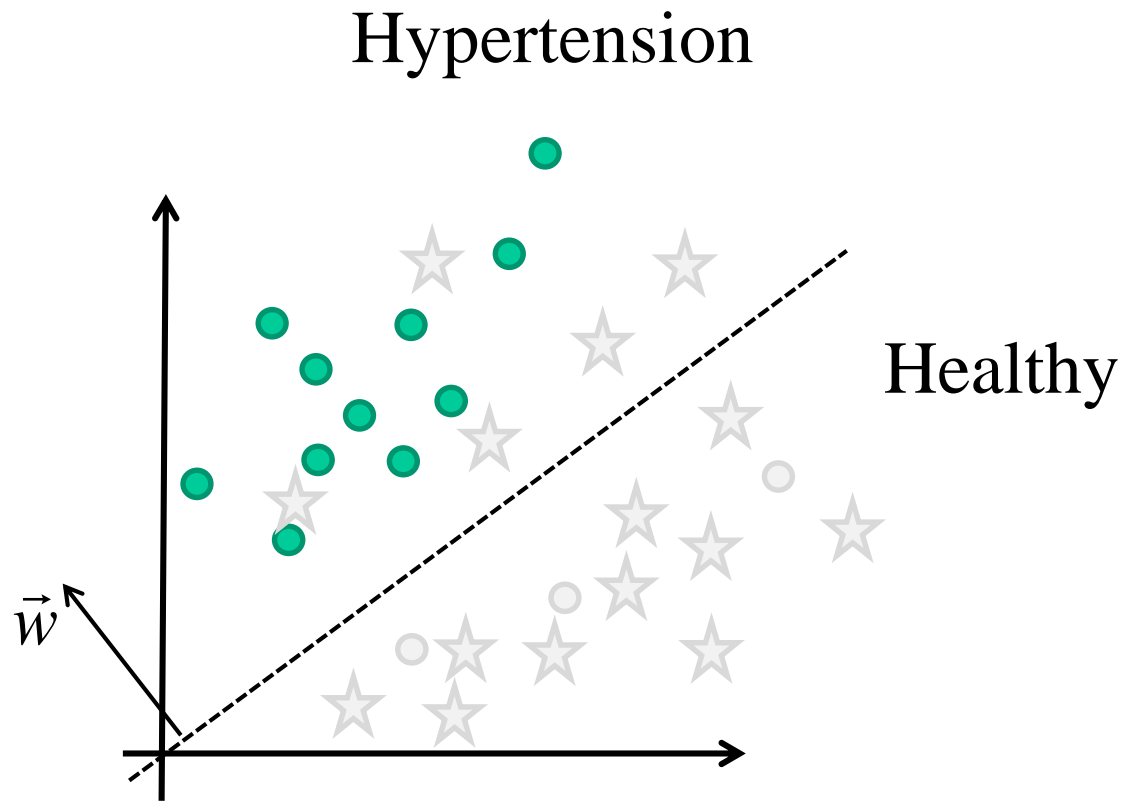
- ✓ The goal is to **train a model** on a **training dataset** with good **generalization** capabilities
- ✓ Has **hyper-parameters** that control the **capacity** of the model, choisis à l'aide d'une procédure de **sélection de modèle**
- ✓ mesure sa performance de **généralisation** sur un **ensemble de test**
- ✓ Aura une meilleure performance de généralisation si la **quantité de données d'entraînement augmente**
- ✓ Peut souffrir de **sous-apprentissage** (pas assez de capacité) ou de **sur-apprentissage** (trop de capacité)

# Evaluation metrics

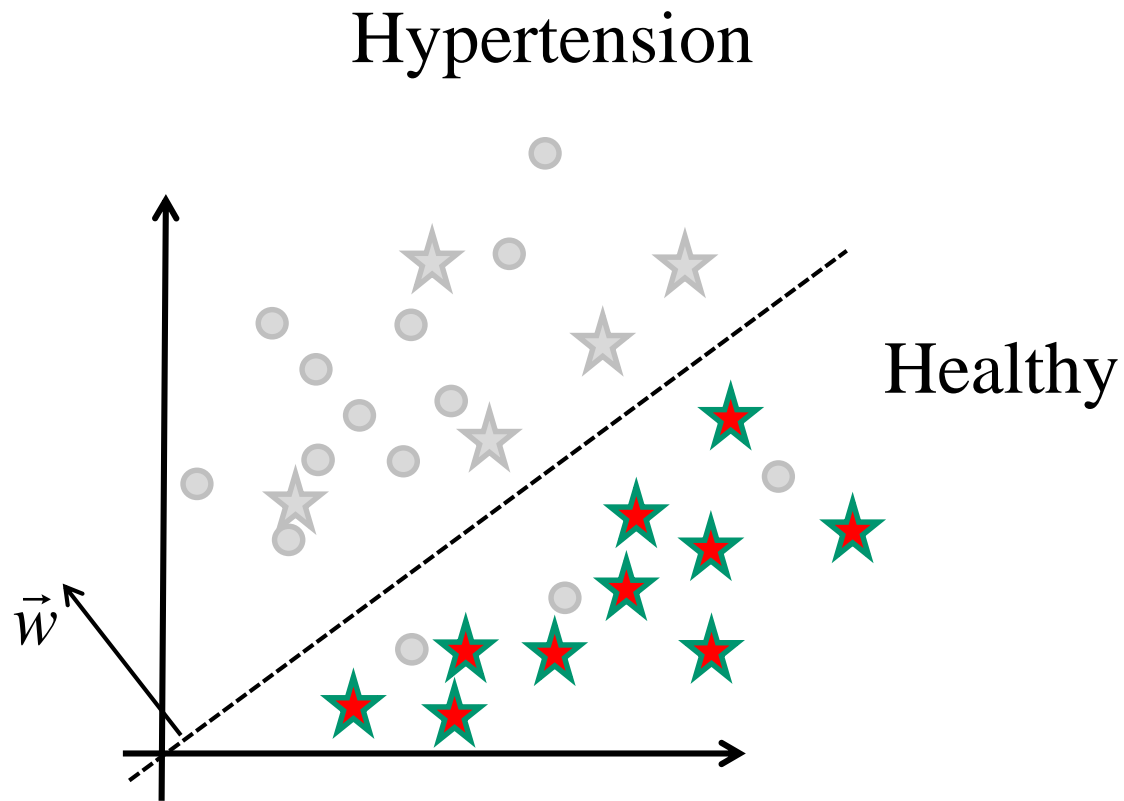
# How to evaluate a model?



# True positive (11)

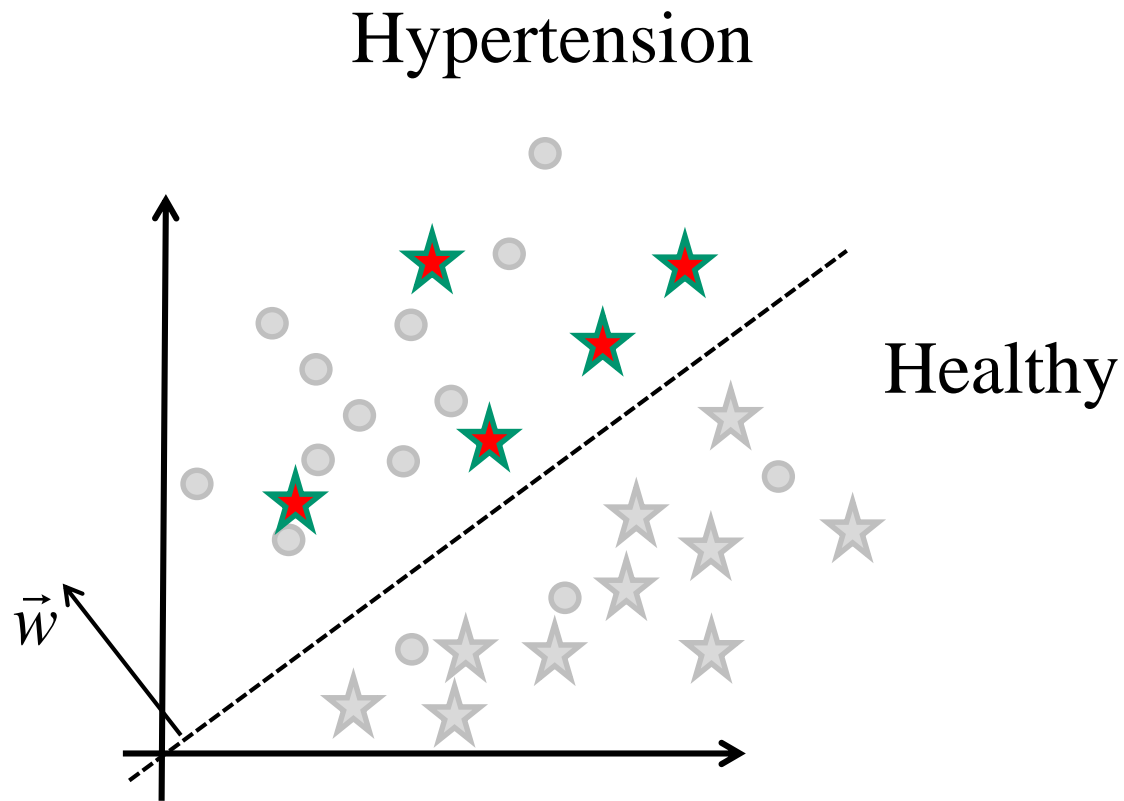


# True negative (10)

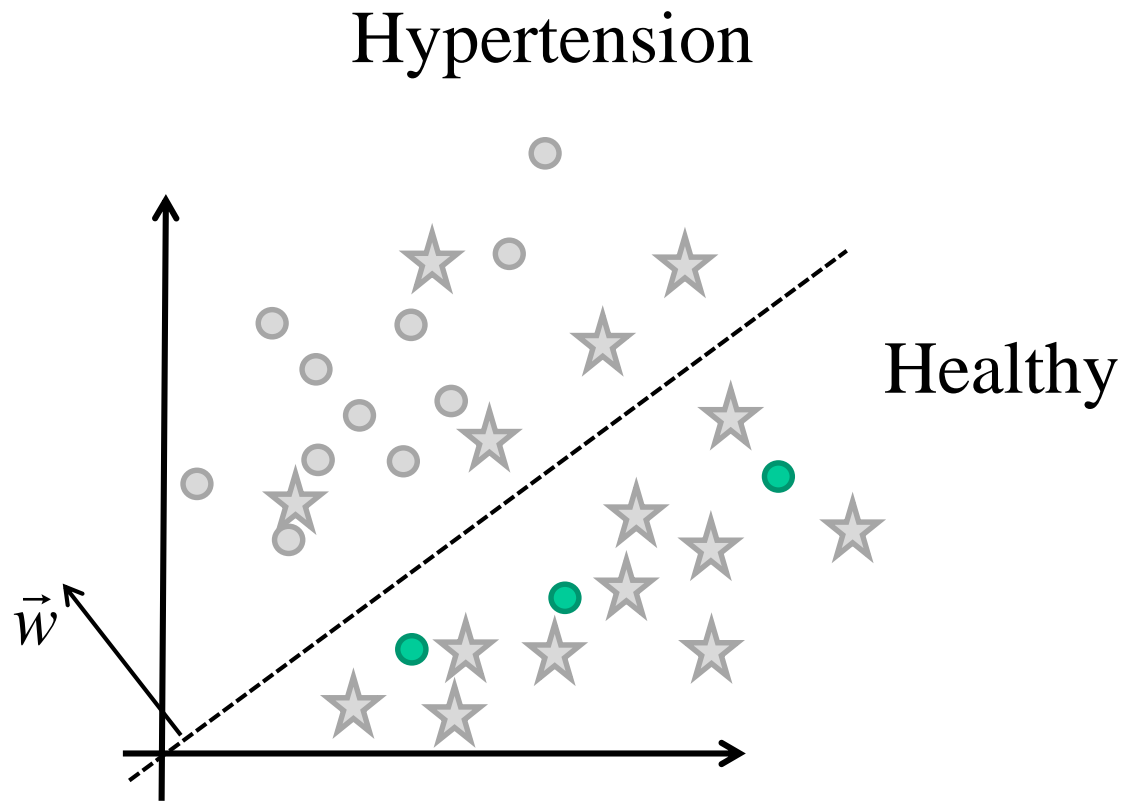




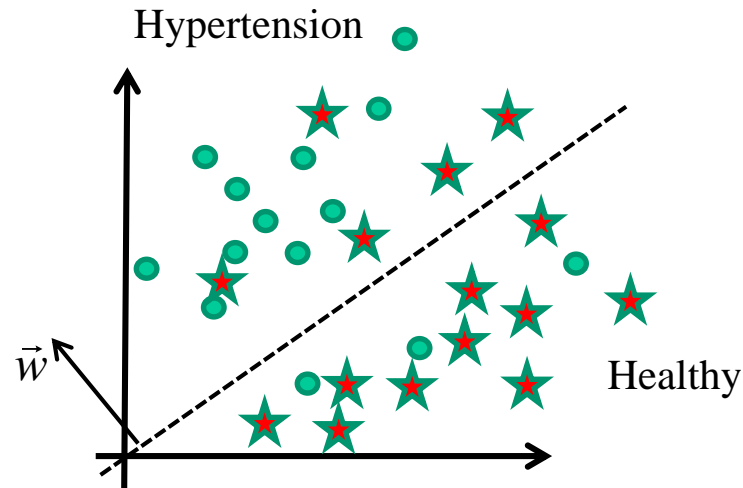
# False positive (5)



## False negative (3)



# Confusion matrix



		Ground truth	
		Positive	Negative
Model prediction	Positive	TP = 11	FP=5
	Negative	FN=3	TN=10

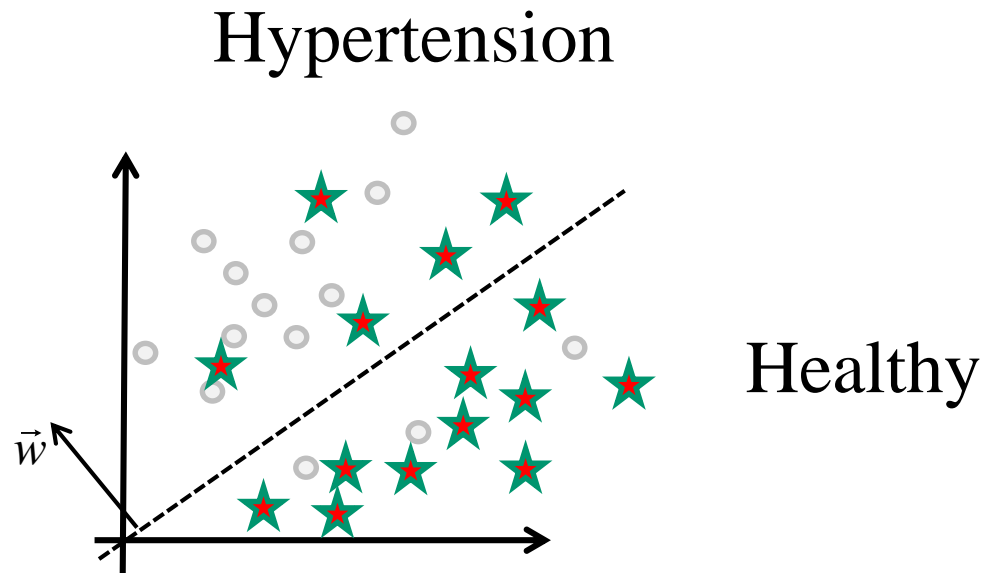
TN + FP = 15 = TOTAL # negative

TP + FN = 14 = TOTAL # positive

TP + FP = 16 = TOTAL # of patients classified +1

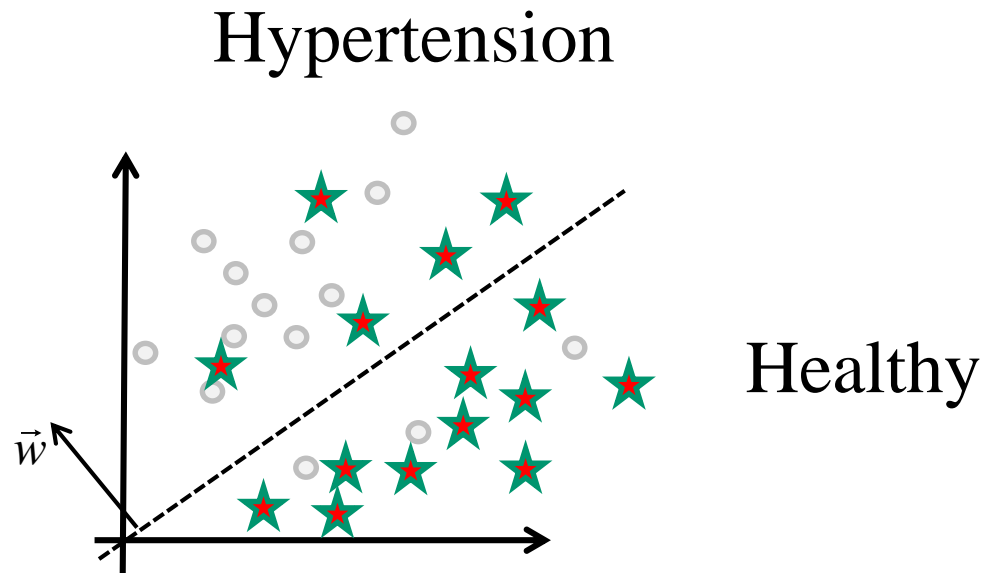
FN + TN = 13 = TOTAL # of patients classified -1

# False positive rate



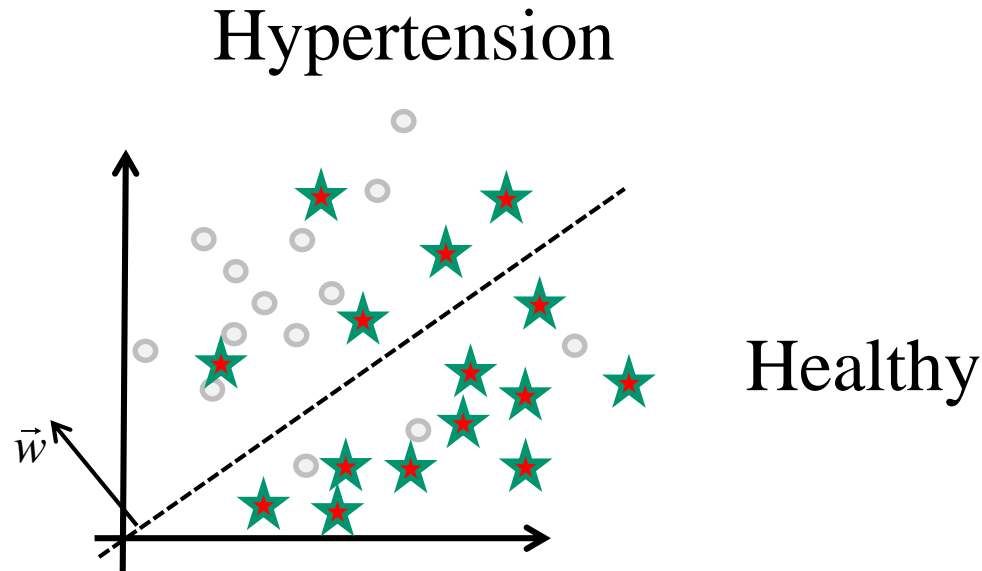
$$\text{FPR} = \text{FP} / (\text{FP} + \text{TN}) = 5/15$$

# Specificity



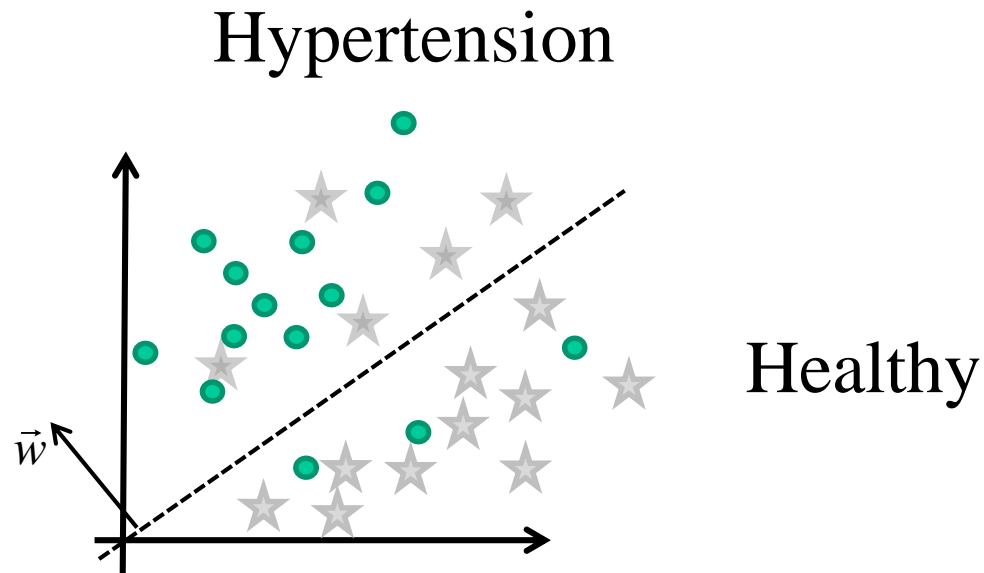
$$Sp = TN/(FP+TN)=11/15$$

# Specificity



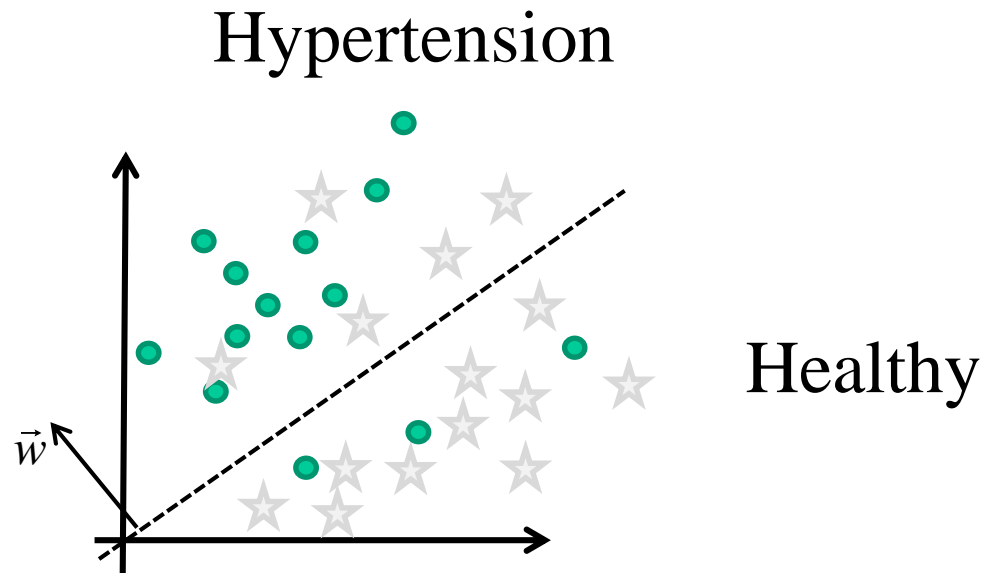
$$Sp = TN/(FP+TN)=1-FNR=11/15$$

# False negative rate



$$\mathbf{FNR} = \text{FN}/(\text{FN}+\text{TP}) = 3/14$$

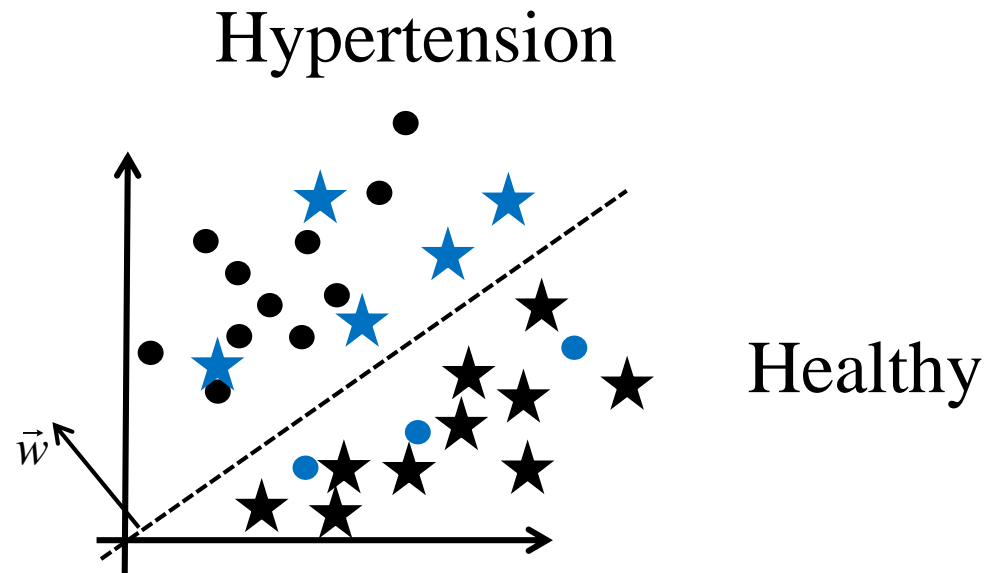
# Recall



$$\mathbf{Re} = \text{TP}/(\text{FN}+\text{TP})=1-\text{FNR}=11/14$$

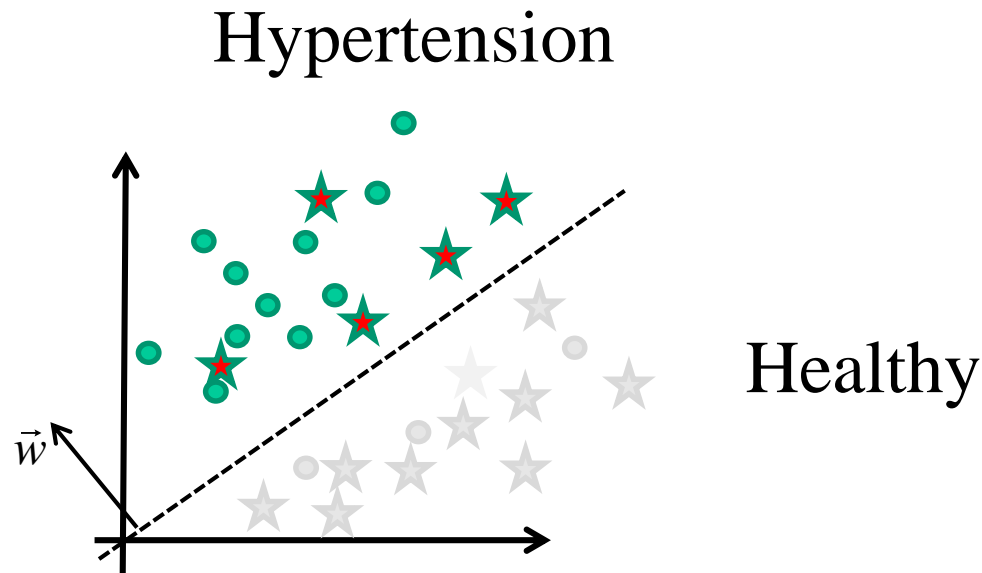


# Accuracy



Rate of good classification =  $(TP+TN)/(FP+FN+TP+TN) = 21/29$

# Precision



$$\mathbf{Pr} = \text{TP}/(\text{TP}+\text{FP}) = 11/16$$

# In short

		Ground truth	
		Positive	Negative
Model prediction	Positive	TP = 11	FP=5
	Negative	FN=3	TN=10

TN + FP = 15 = TOTAL # negative

TP + FN = 14 = TOTAL # positive

TP + FP = 16 = TOTAL # of patients classified +1

FN + TN = 13 = TOTAL # of patients classified -1

False positive rate =  $FP/(FP+TN) = 5/15$

False negative rate =  $FN/(FN+TP) = 3/14$

Specificity (**Sp**) =  $TN/(FP+TN)=1-FPR=10/15$

Recall (**Re**) =  $TP/(TP+TN)=11/14$

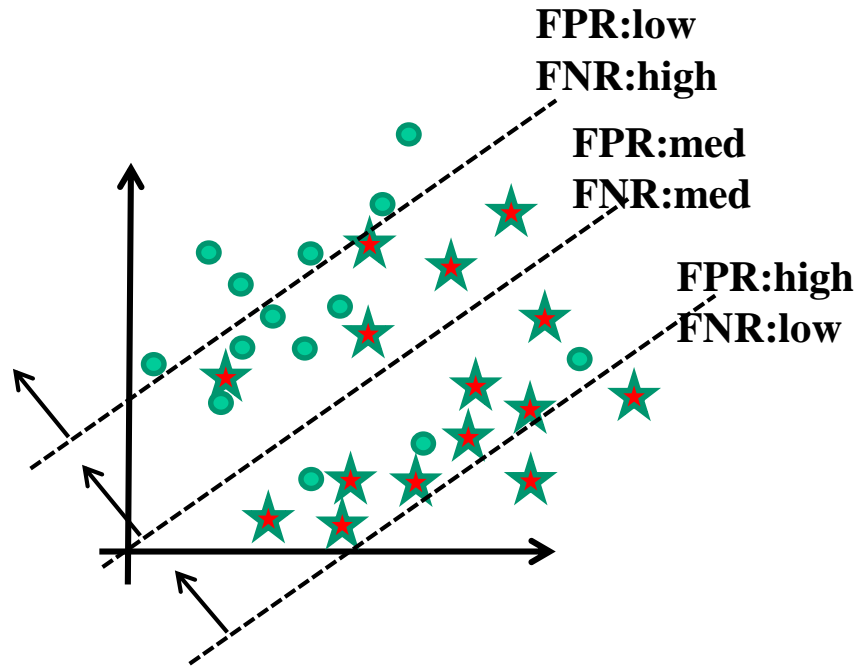
Precision (**Pr**) =  $TP/(TP+FP) = 11/16$

Accuracy =  $(TP+TN)/(FP+FN+TP+TN) = 21/29$

**F-measure** =  $2 * (Re*Pr)/(Pr+Re)=0.73$

# Different thresholds, different results

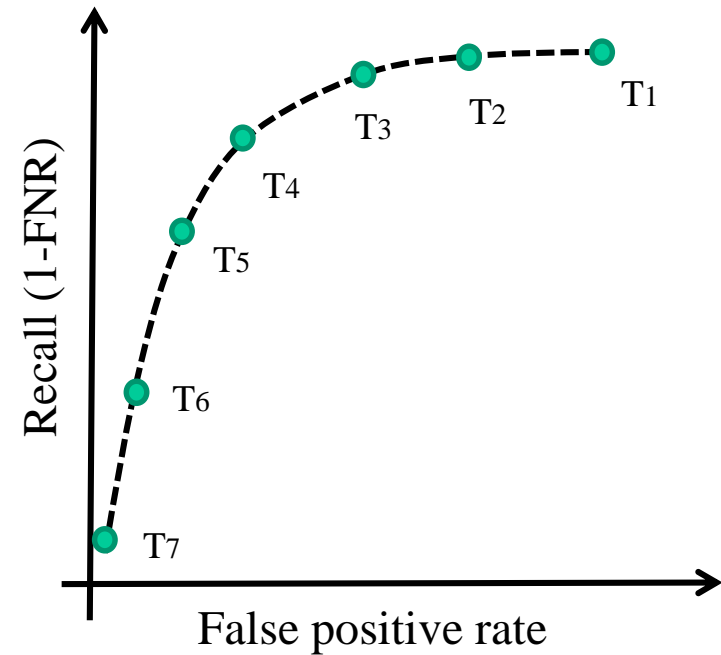
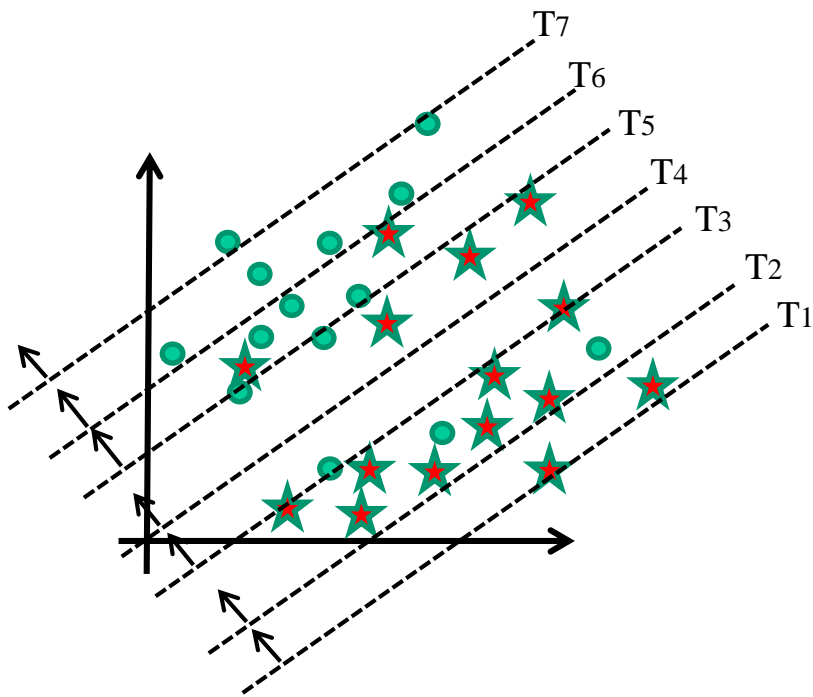
$P(x)$



# ROC curves

Compute Recall and FPR for **different thresholds**

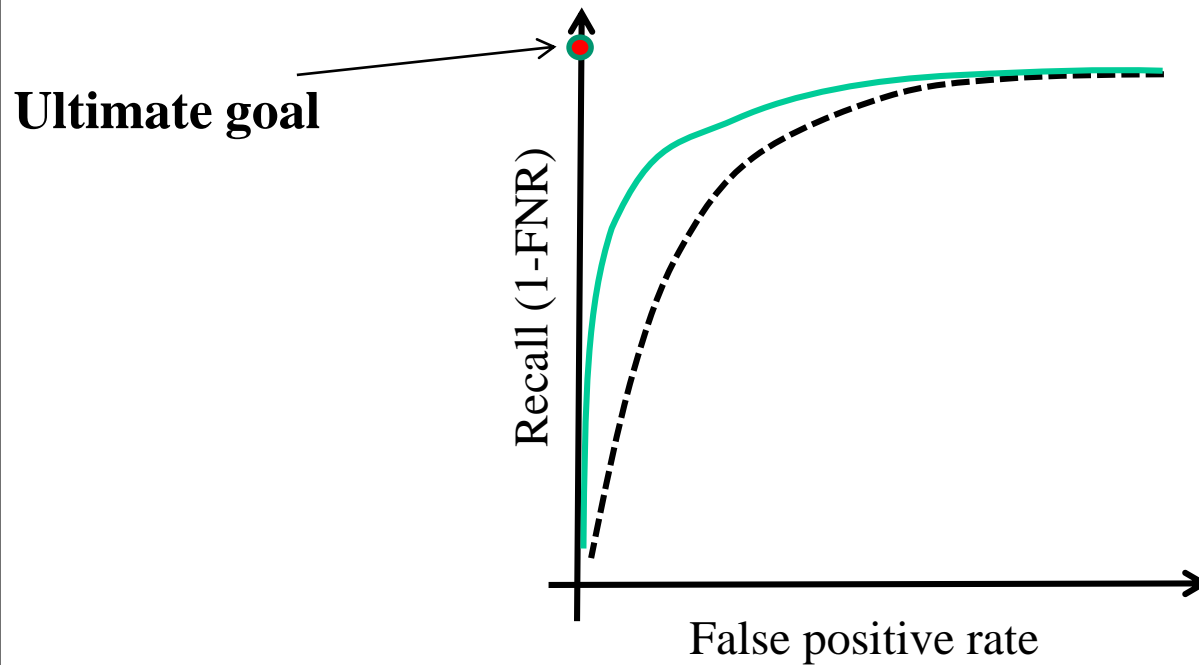
$P(x)$



# ROC curves

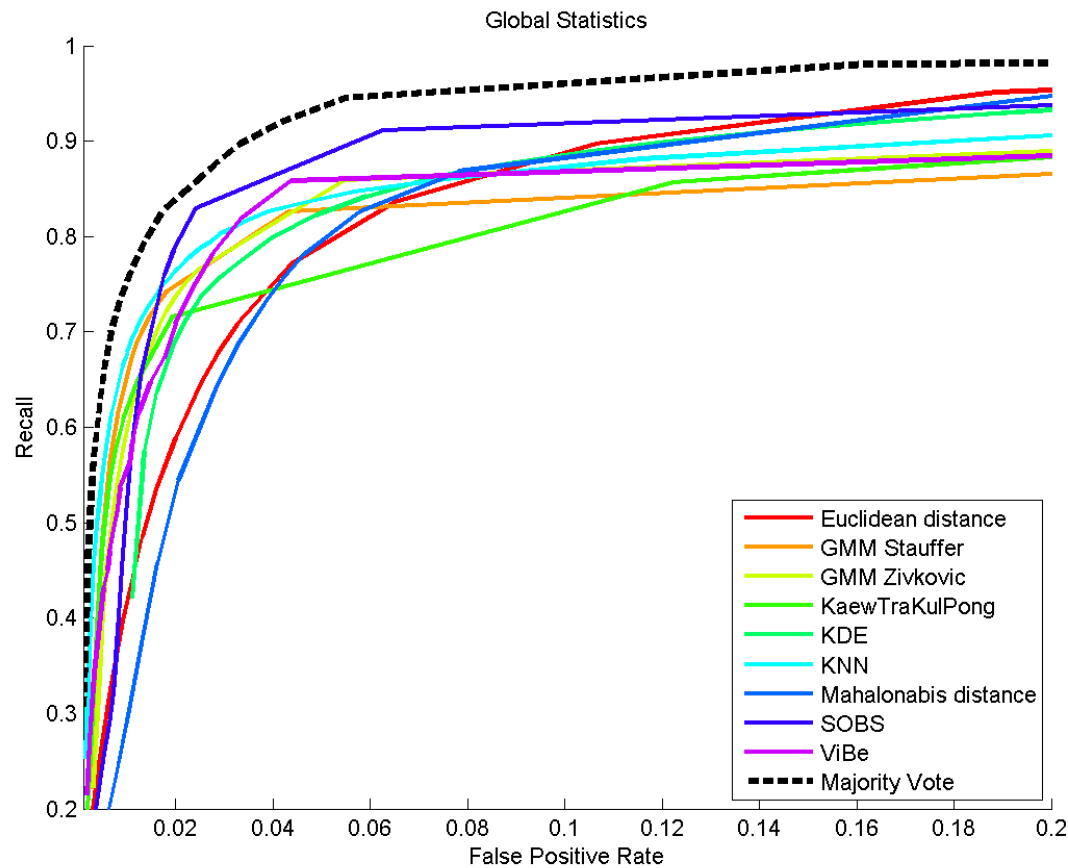
Good way for comparing methods

Which method is best?



# ROC curves

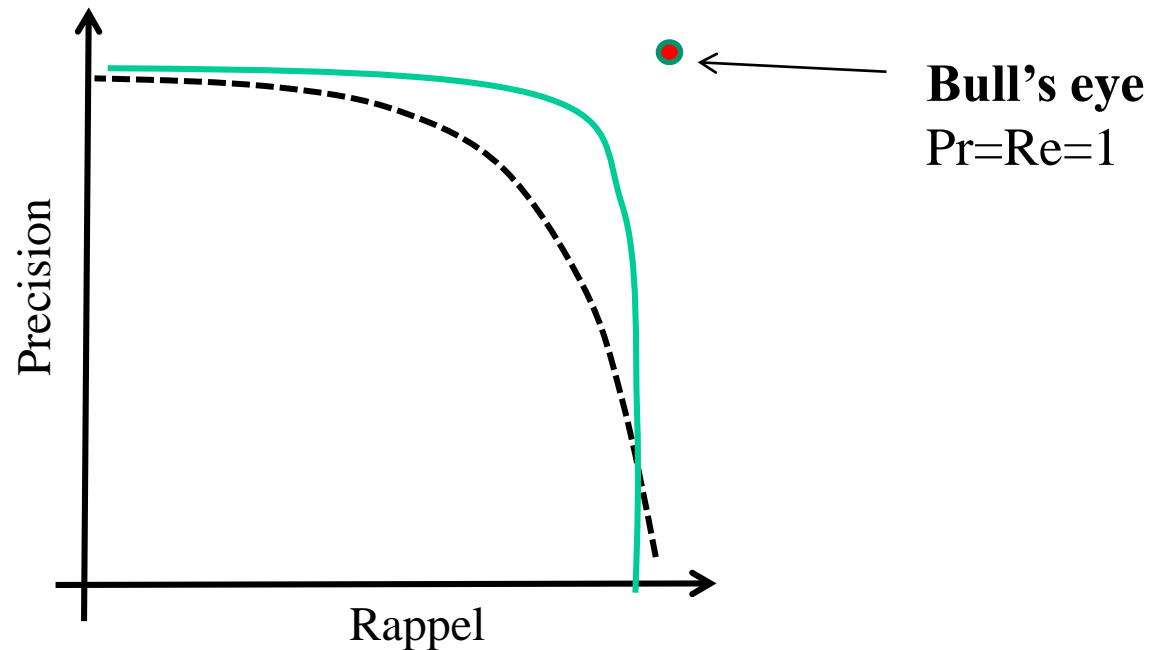
**Example :** 10 motion detection methods



# Precision recall curve

Same spirit that the ROC curve

Which one is best?





# Model ensembling

# Why use only one model?

Does combining several model works well?

In practice ... **yes!**

Combining several models is often called **ensembling**

# Why use only one model?

## Combining what?

- Several **different** models
- The **same model** trained with different **hyperparameters**
- The **same model** trained on different **data**.

# Typically 2 ways of combining models

- *Bagging* : good for models with a **large capacity**
- *Boosting* : good for models with a **low capacity**

# Typically 2 ways of combining models

- *Bagging* : good for models with a **large capacity**
- *Boosting* : good for models with a **low capacity**

# Bagging

- A simple way of combining several models

- FOR  $i = 1$  to  $m$

- Train model  $y_{w,i}(\vec{x})$

- **Ensemble** the  $m$  models

- For **regression**

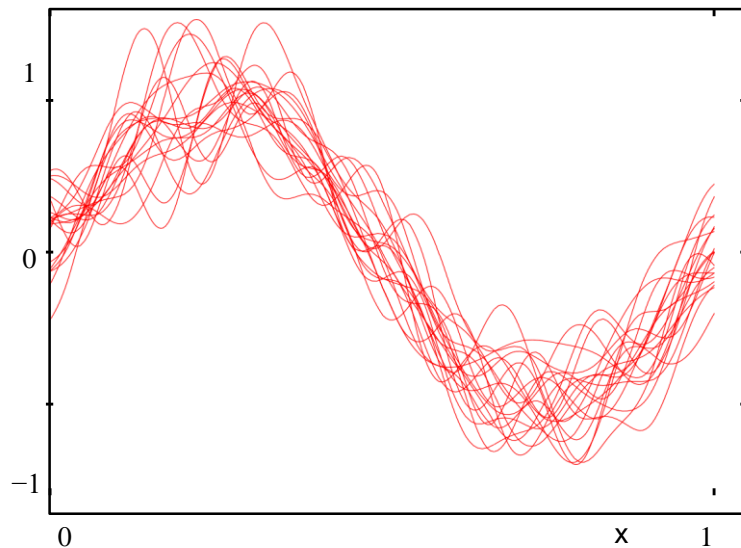
- $y_{COM}(\vec{x}) = \frac{1}{m} \sum_{i=1}^m y_{w,i}(\vec{x})$

- For a **classification**

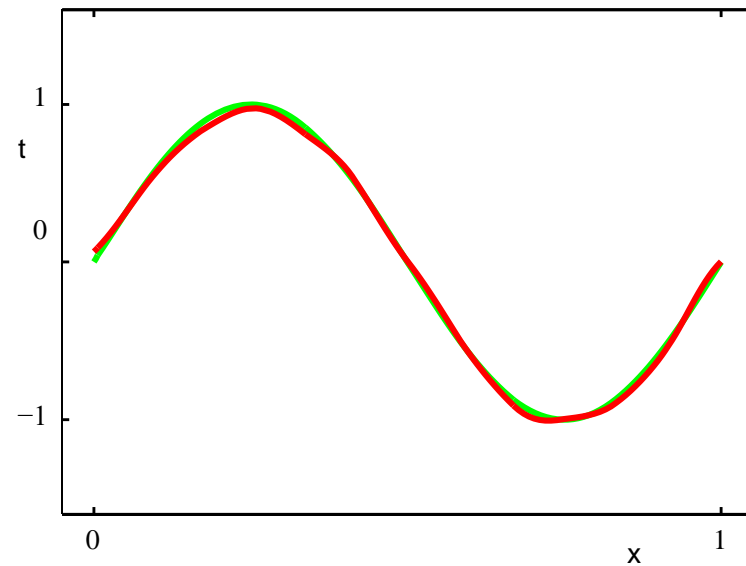
- **majority vote**

# Ex: polynomial regression $M=25$

100 models trained on  
100 different training sets



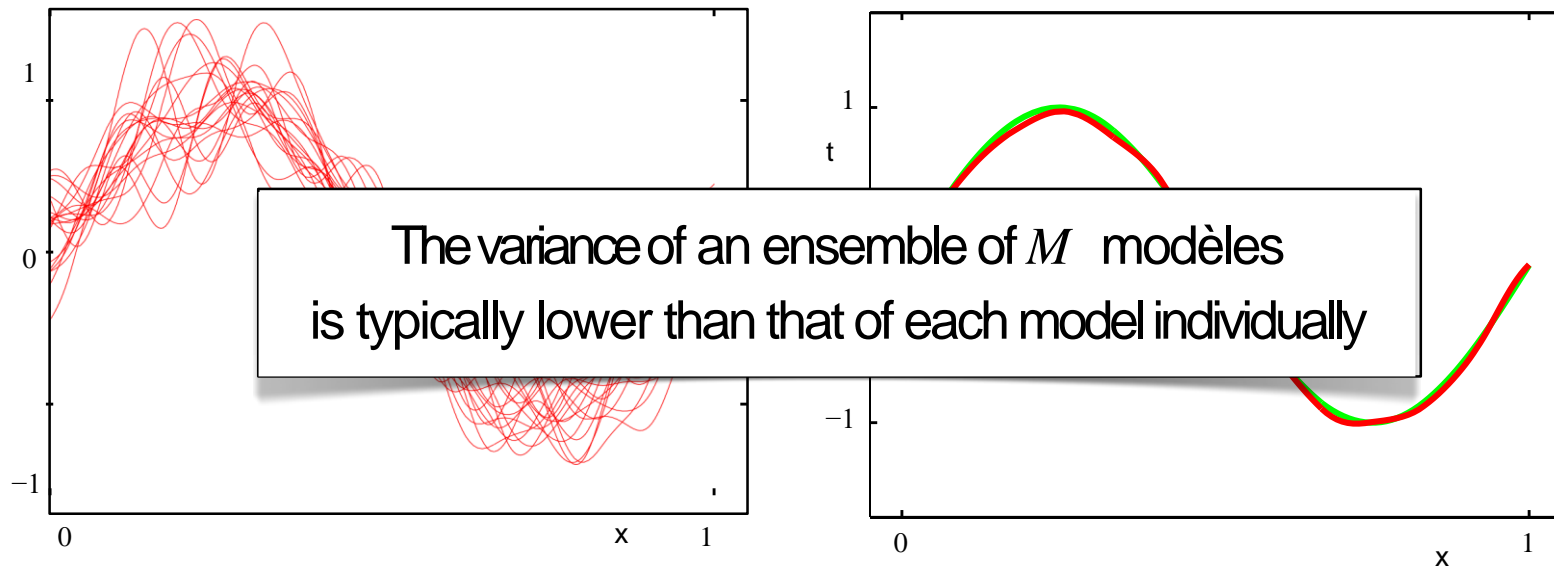
Ensemble of 100 models  
 $V_s$   
True model



# Ex: polynomial regression $M=25$

100 models trained on  
100 different training sets

Ensemble of 100 models  
 $V_s$   
True model





# Bootstrap

What if you do not have enough data for building 100 trainingsets? What can you do???

Solution 1 : **Data augmentation**

Solution 2 : **Bootstrapping**.

FOR j from 1 to 100 DO

$$D_{bootstrap} = \{ \quad \}$$

FOR N iterations

- Choose randomly a natural number between 1 and N

$$- D_{bootstrap} = D_{bootstrap} \cup \{(\vec{x}_n, t_n)\}$$

$$D_{train,j} = D_{bootstrap}$$

Train j-th model

# In short we have seen

- Supervised vs unsupervised learning
- Regression vs Classification
- Linear vs non-linear models
- Parameters vs hyper-parameters
- Over vs Underfitting
- Cross-validation
- Metrics
- Ensembling – Bootstrapping

Thank you!