



Montreal, July 8–12

DLM I 2024

Basics of deep learning part 1

By

Pierre-Marc Jodoin



UNIVERSITÉ DE
SHERBROOKE

Before we start

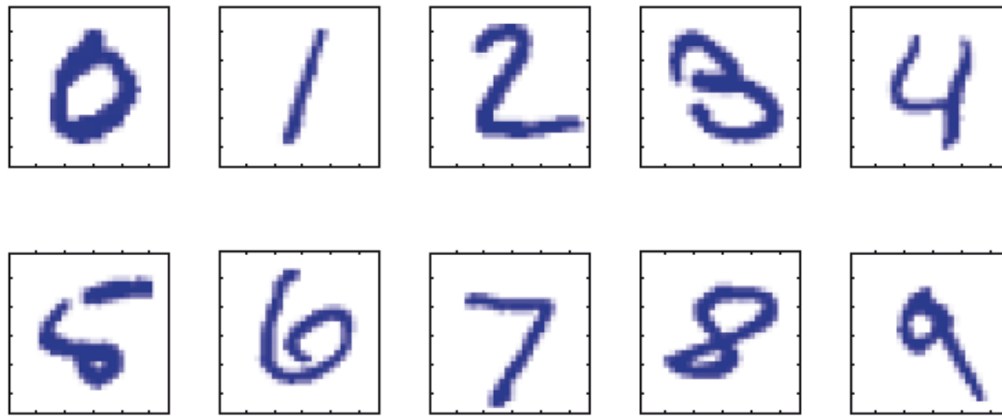


jodoin.github.io/dlmi2024

What is machine learning?



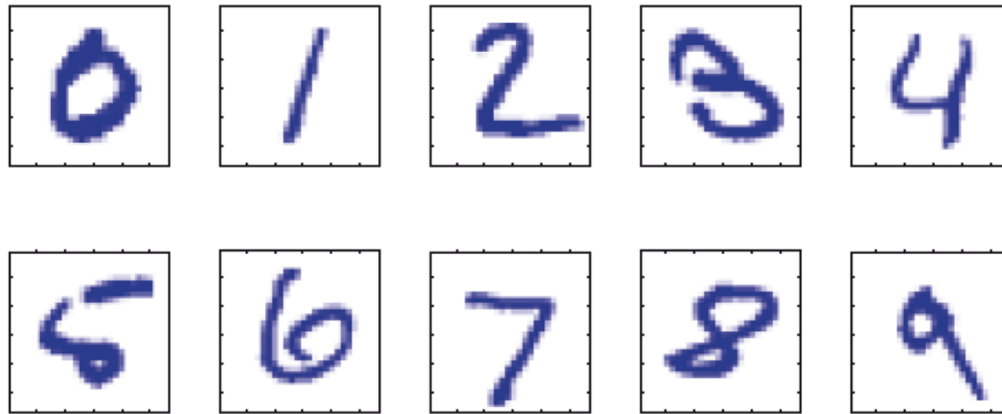
Question : how can one recognize hand written digits?



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?

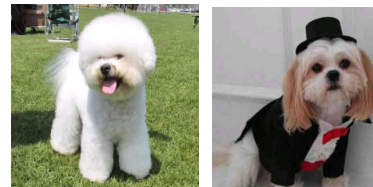


Answer : ~~Design your own rules?~~ **Wrong**

➤ Bad generalization



➤ Often difficult



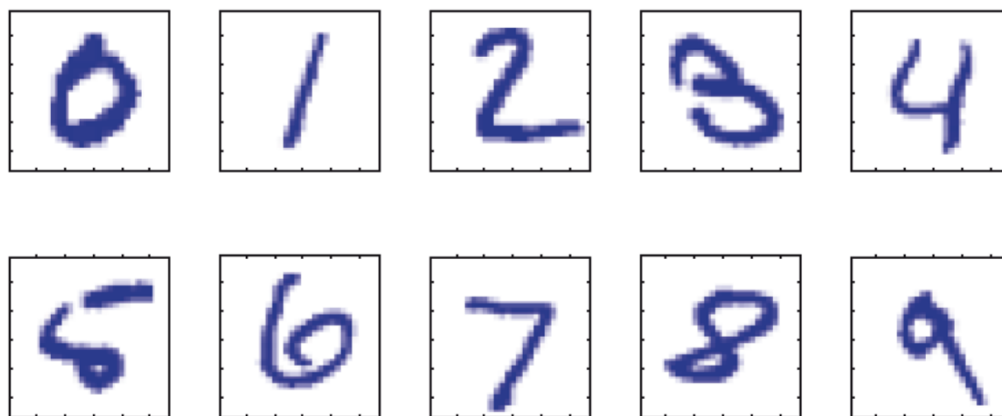
Vs



Dogs

Birds

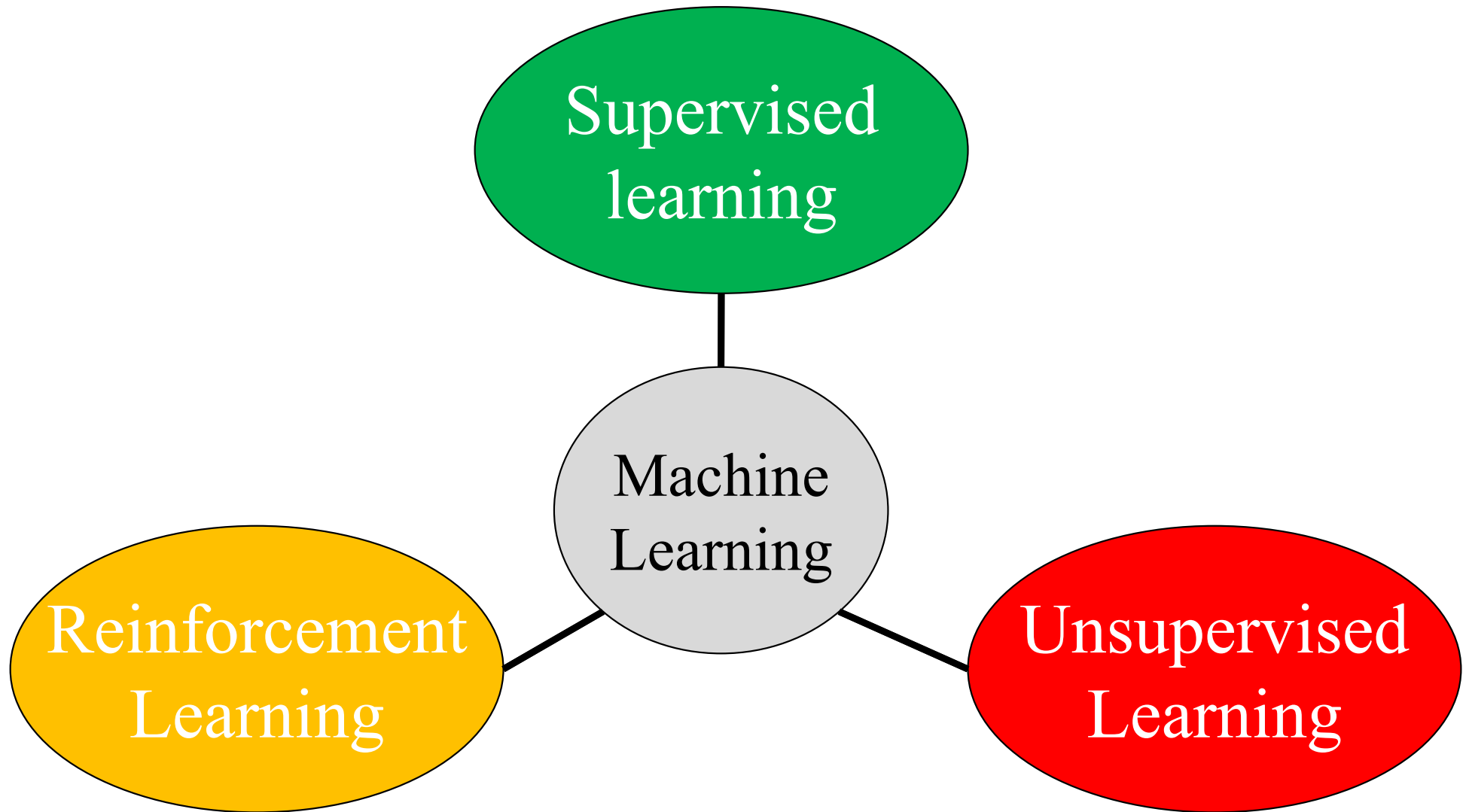
Question : how can one recognize hand written digits?



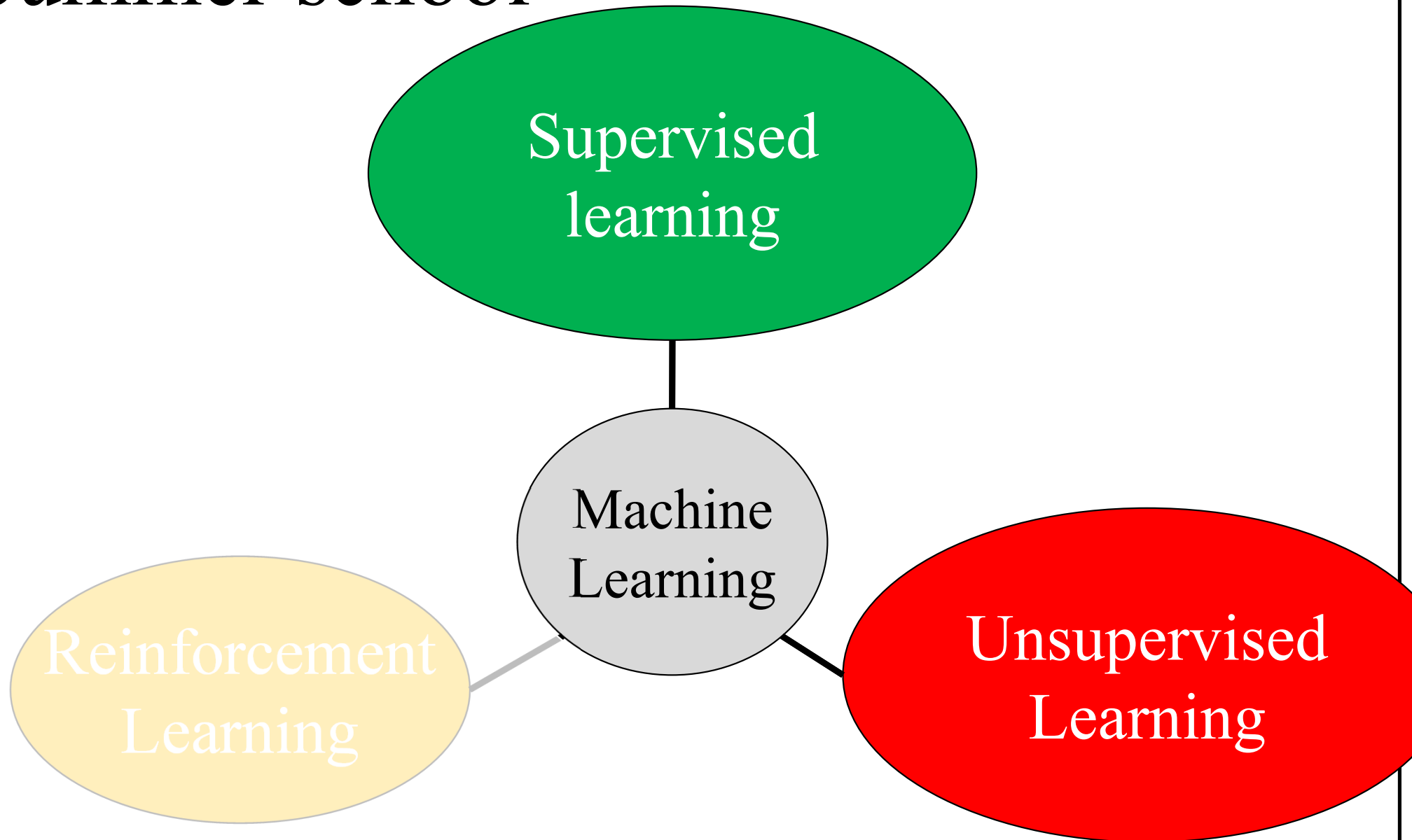
Answer : Let the computer « **learn** » the rules

➤ *Main goal of machine learning*

Three large families

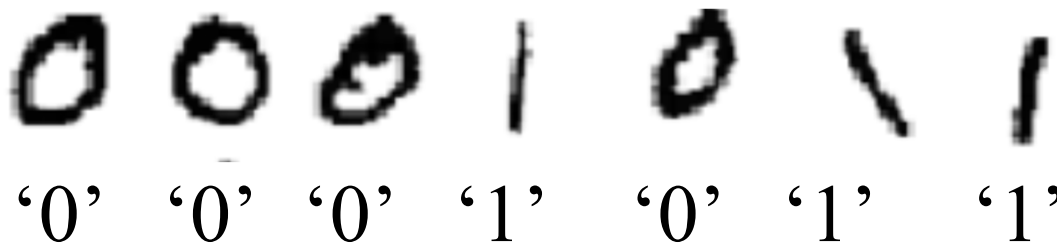


Summer school

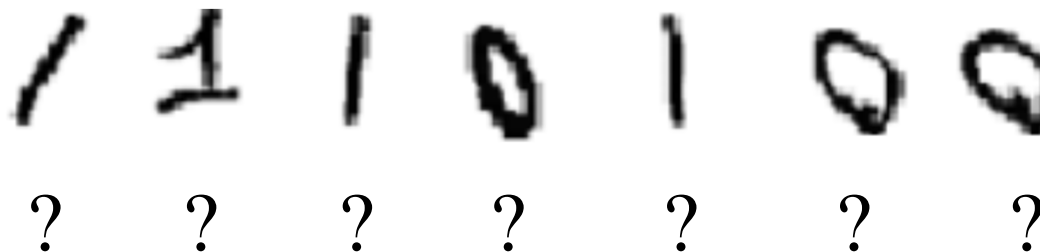


Supervised learning

Provide the algorithm with **annotated training data**

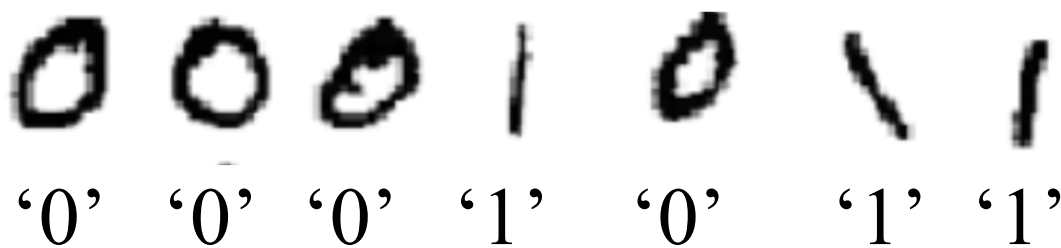


...and the algorithm returns a function capable of **generalizing** on new data



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

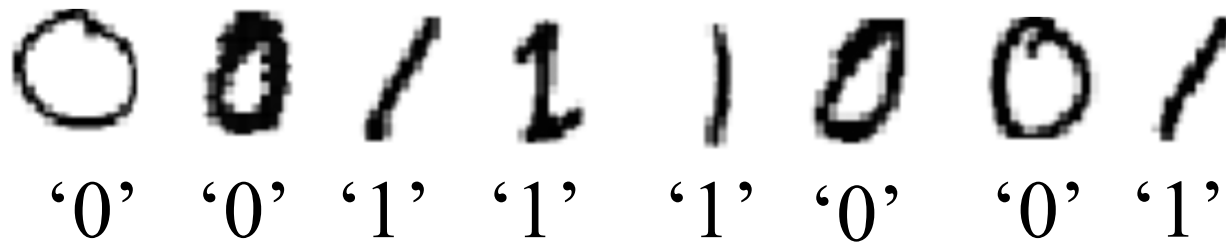
the goal **is to learn** a function that may predict t_i given \vec{x}_i

$$y_W(\vec{x}_i) \rightarrow t_i$$

where W are the **parameters** of the model.

Supervised learning

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



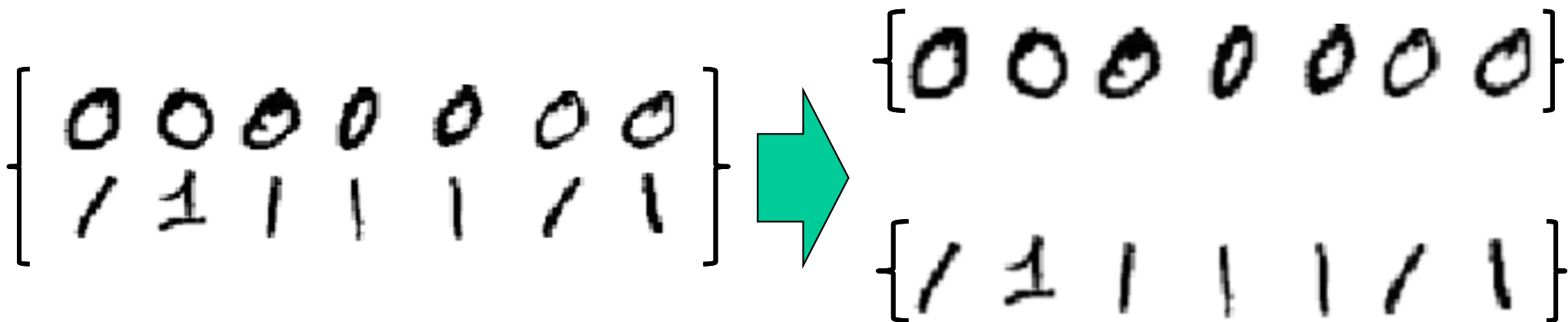
Supervised learning

Unsupervised learning

Unsupervised learning

When no target is explicitly provided

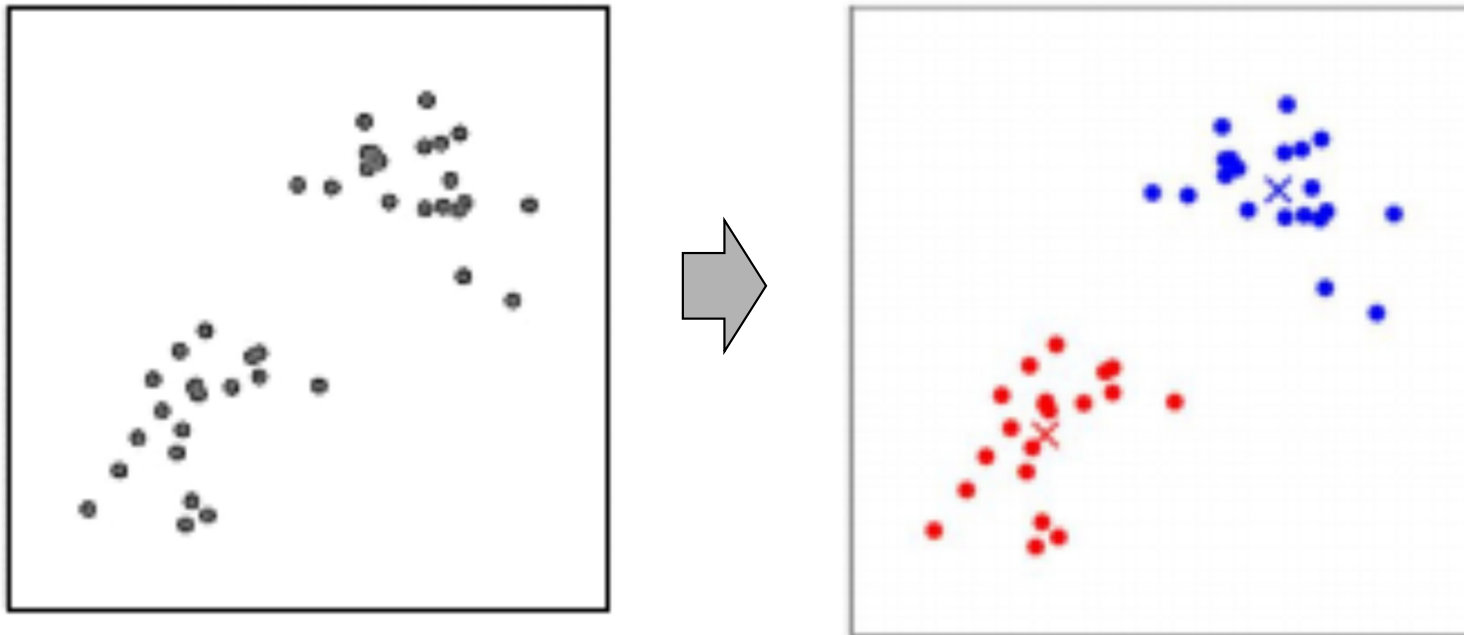
➤ E.g. data *clustering*



Unsupervised learning

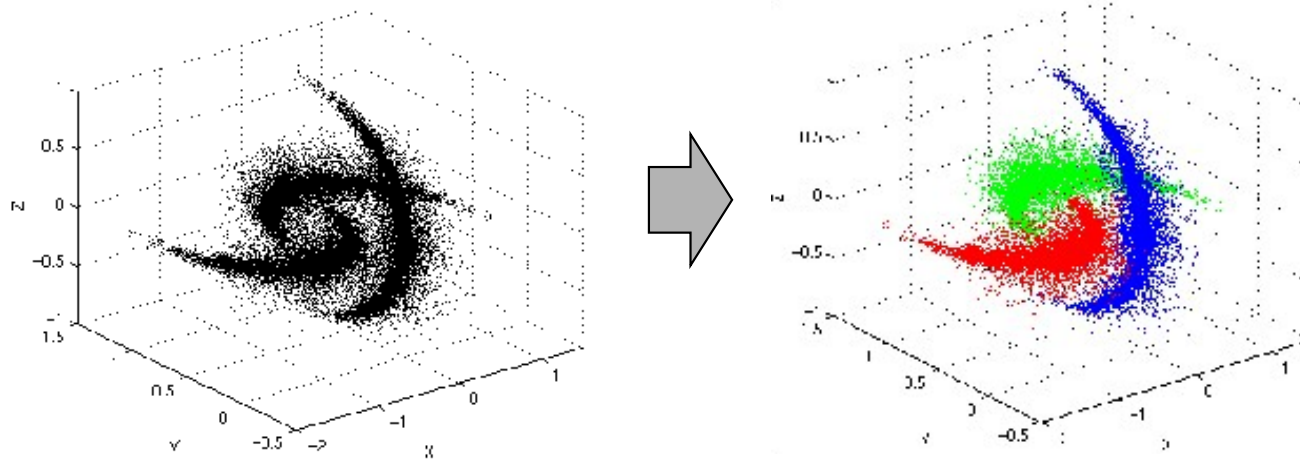
When no target is explicitly provided

➤ E.g. *data clustering*



Unsupervised learning

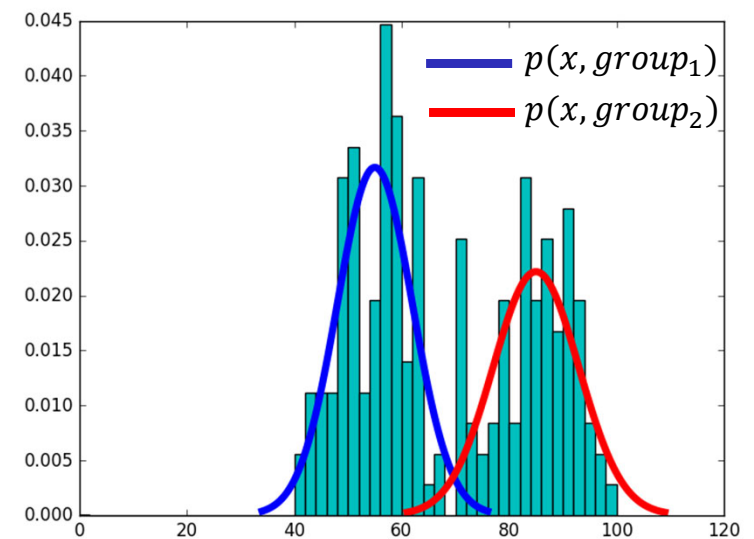
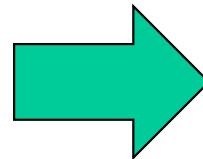
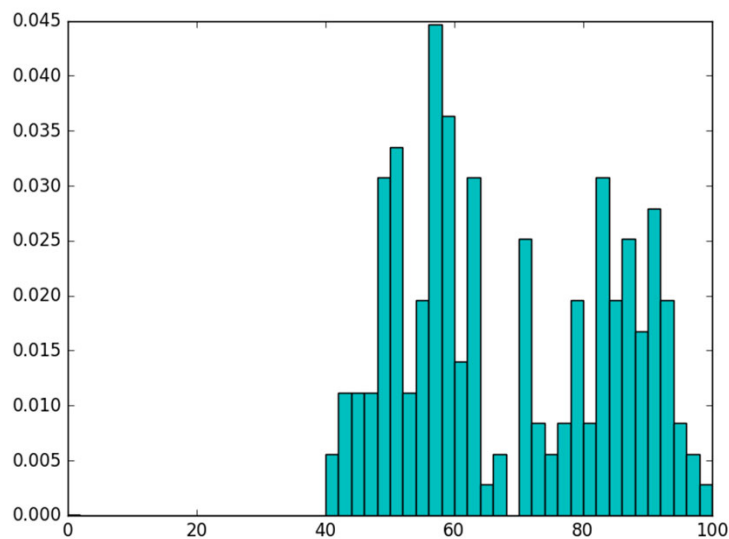
No limit to dimensionality. Could be 3D, 4D,...100kD



Unsupervised learning

Probability density function estimation

Example : find two groups of patients following a memory test



Supervised vs non-supervised

Supervised learning : there is a target

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



**Main topic of
the school**

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
Semi-supervised learning
Graph Neural Nets
Transformers
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

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

Supervised learning

Classification vs regression

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}

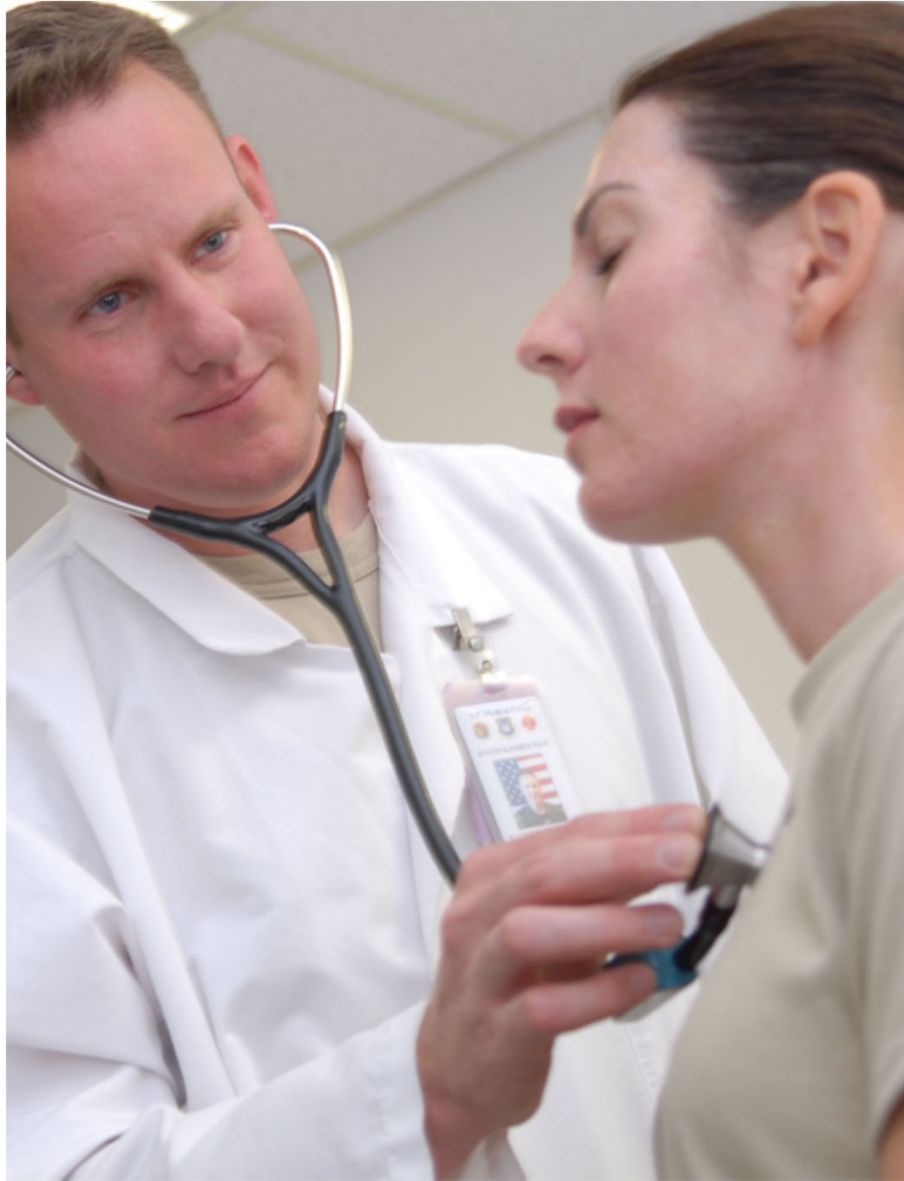
- **Regression** : the target is a real number $t \in \mathbb{R}$
 - Exemple : prediction of life expectancy
 - ✓ \vec{x} : vector of medical measures, age, sex, etc.
 - ✓ t : number of months before death.

Supervised learning

Two main applications

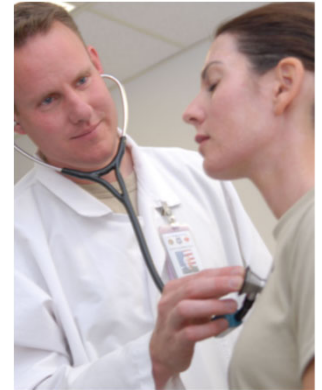
- **Classification** : the target is a class label $t \in \{1, \dots, K\}$
 - Exemple : disease recognition
 - ✓ \vec{x} : vector of medical measures, age, sex, etc.
 - ✓ t : myocardial infarction, dilated cardiomyopathy, hypertrophic cardiomyopathy, normal
- **Régression** : the target is a real number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

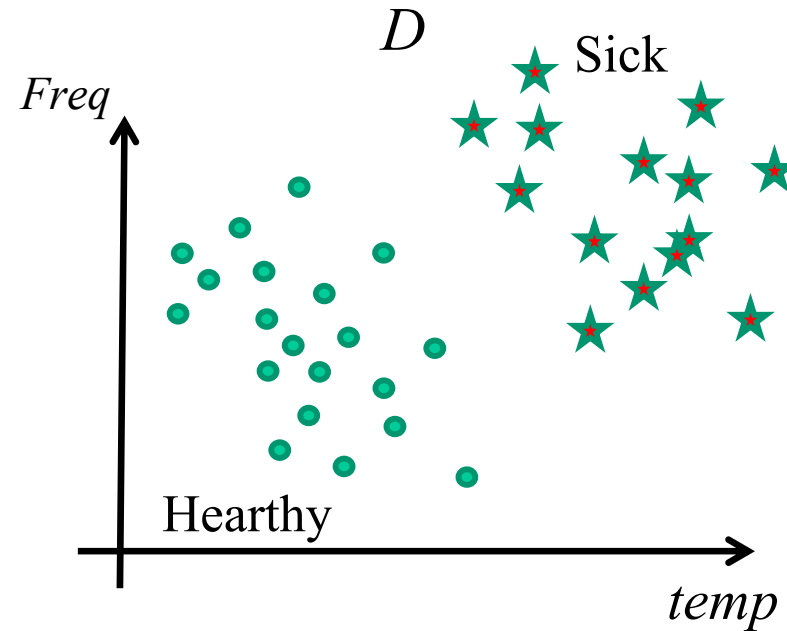
Simple example of binary classification



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

\vec{x} t

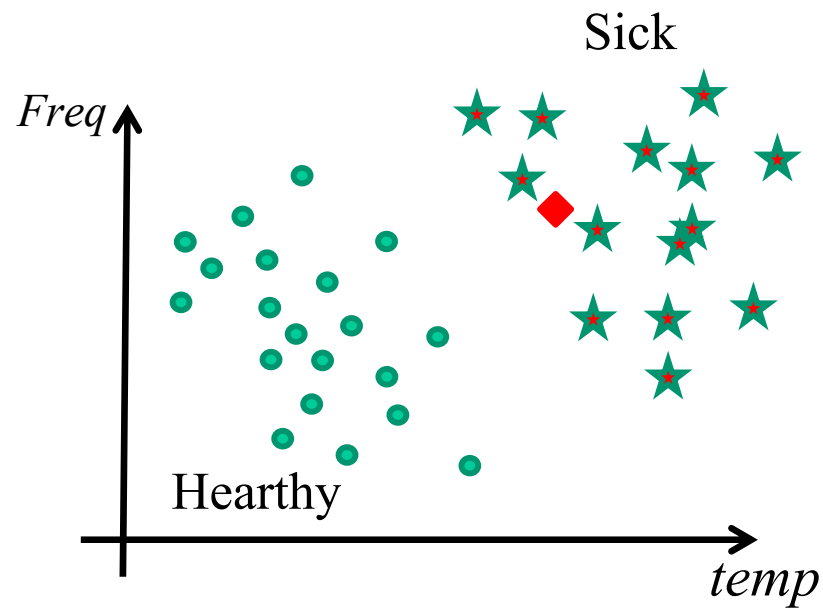


Simple example of binary classification

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



From Wikimedia Commons
the free media repository

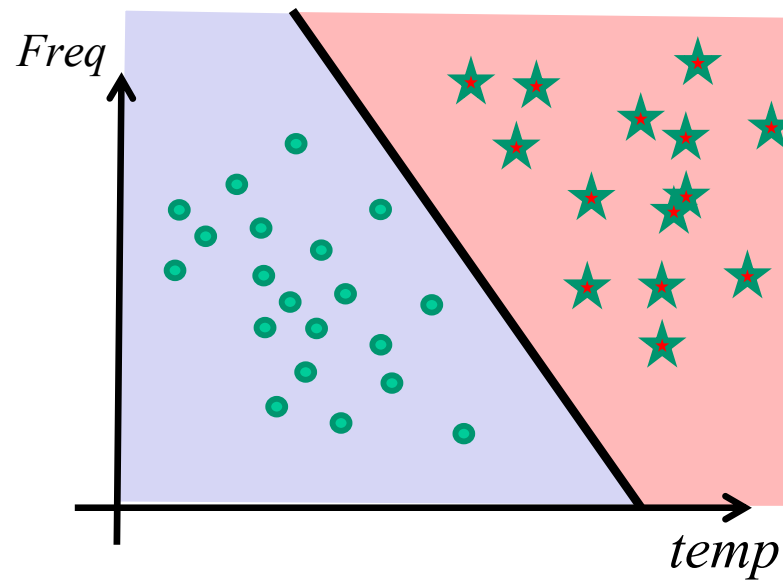


Solution



From Wikimedia Commons
the free media repository

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

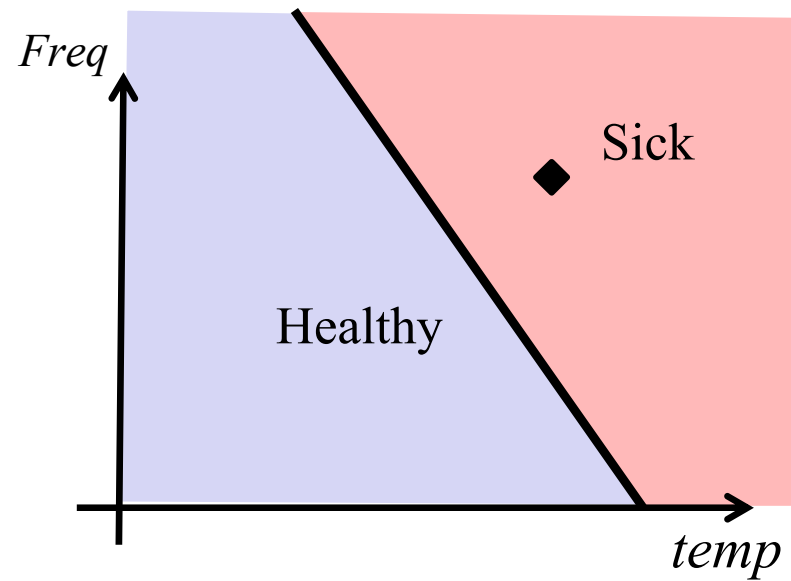


Solution



From Wikimedia Commons
the free media repository

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

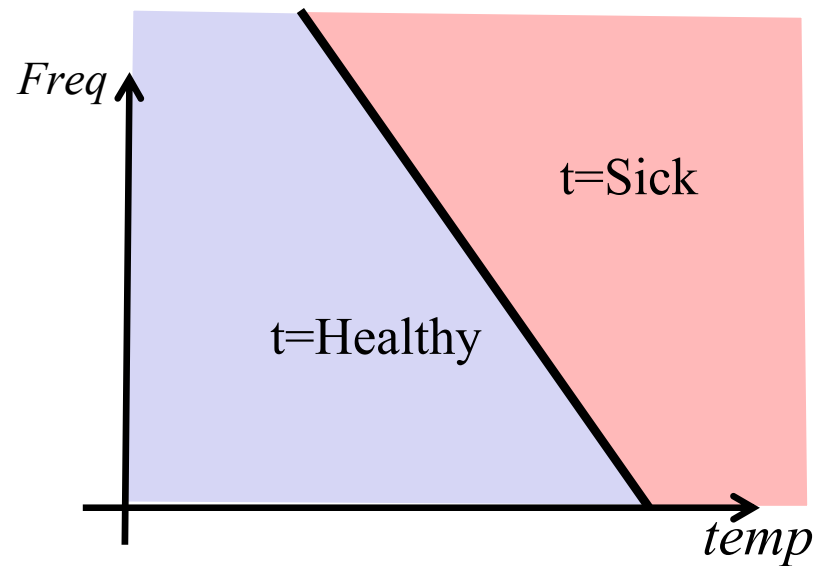


More formally

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

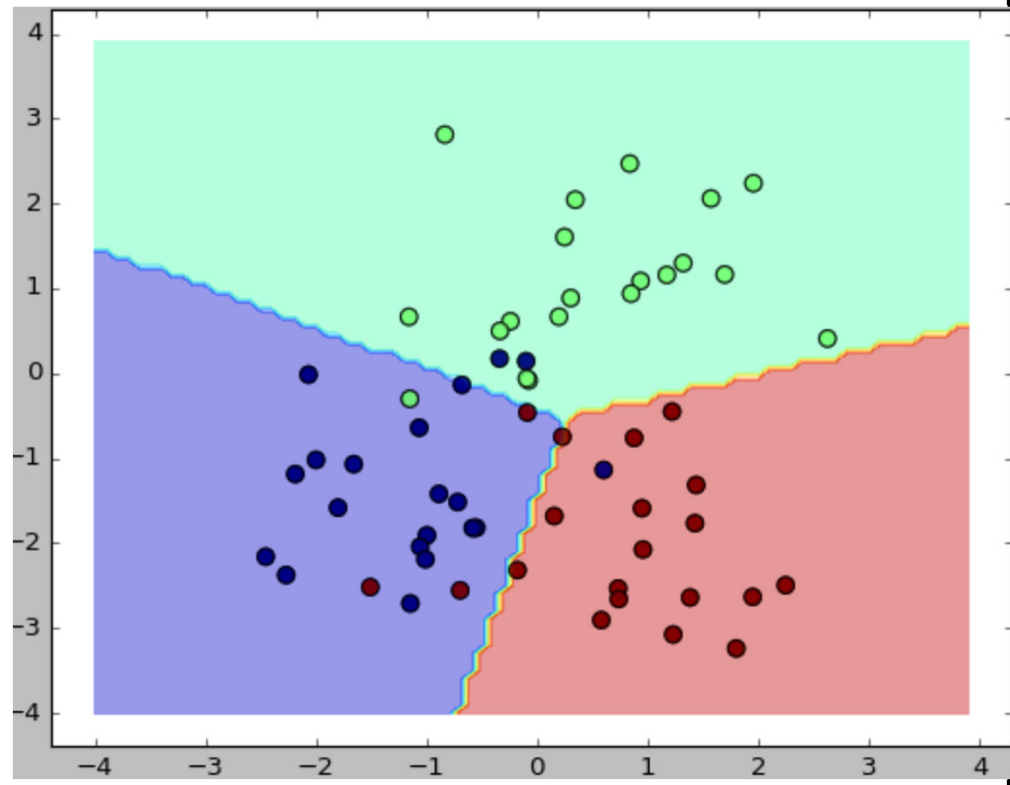
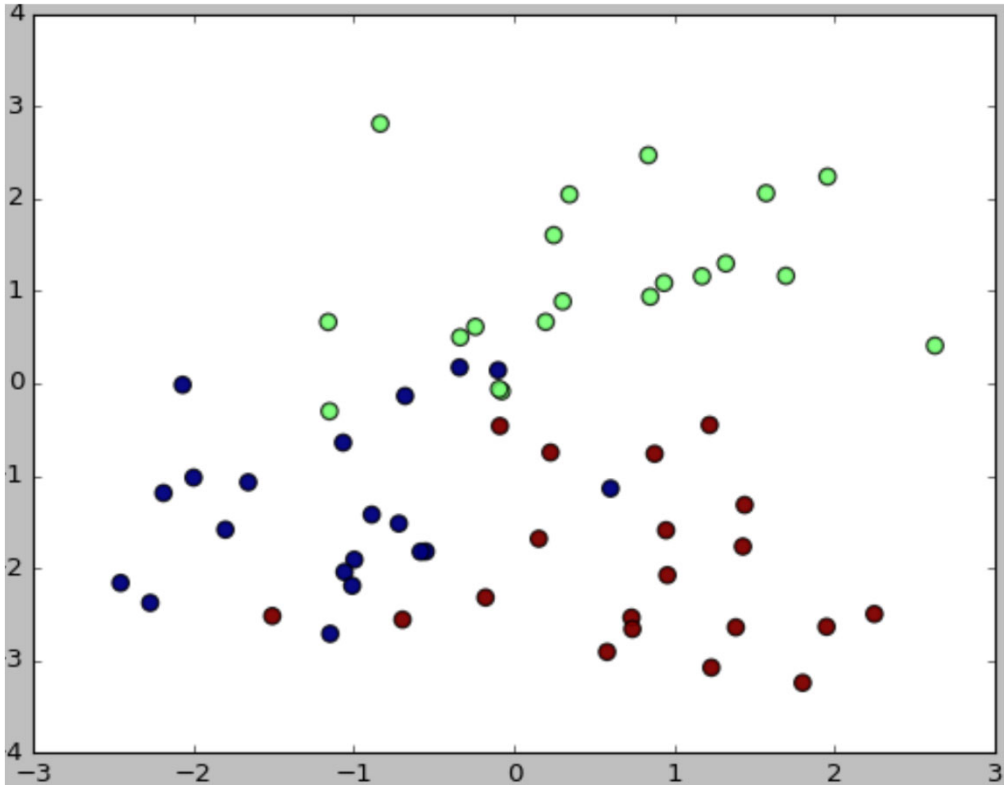


From Wikimedia Commons
the free media repository



Classification

3-class case



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

Once training is over

$y_w(\bullet) = \text{class 1}$

$y_w(\bullet) = \text{class 2}$

$y_w(\bullet) = \text{class 3}$

Example of a classification dataset

MNIST



Example of a classification dataset

MNIST

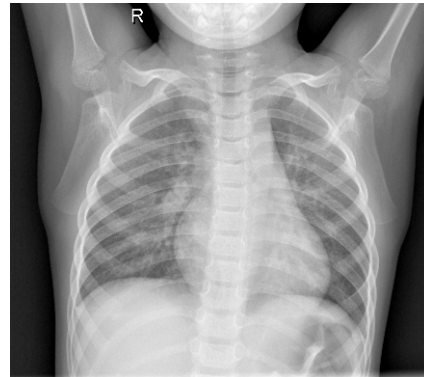
- 10 classes
- 70,000 images
 - => 60,000 training
 - => 10,000 test
- Images are in grayscale
 - => 28x28

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

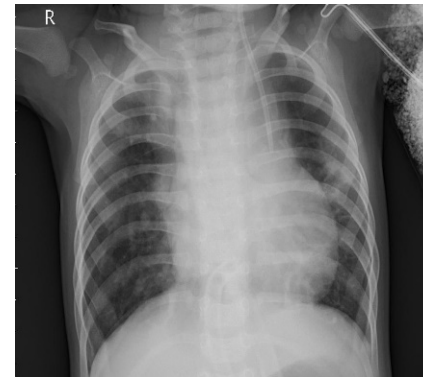
Example of a medical classification dataset

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

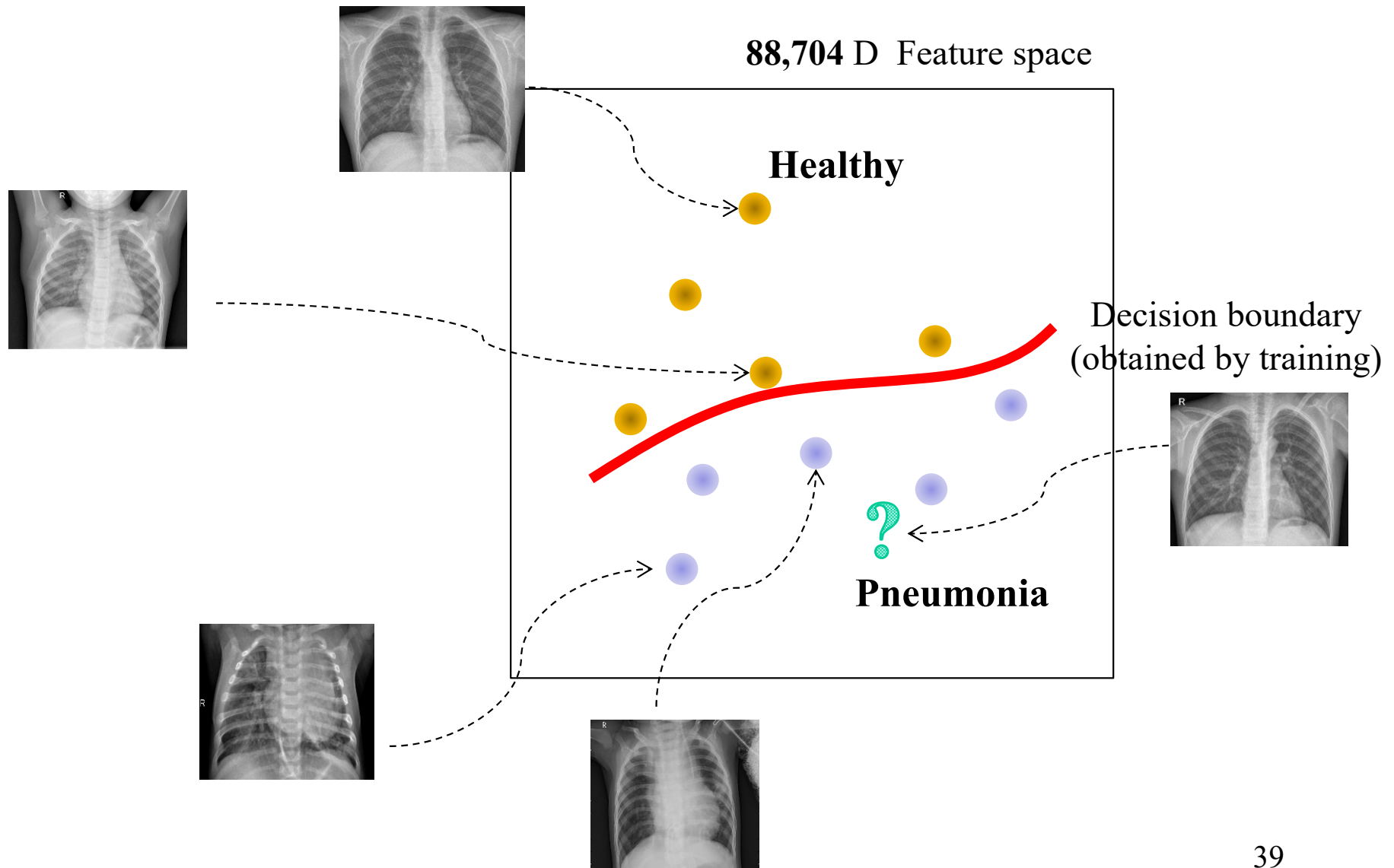
- 2 classes
- 5,840 images,
 - => 5,216 training
 - => 624 test
- Each image is in grayscale
 - => 336 x 264*

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

* Rescaled version

Supervised learning

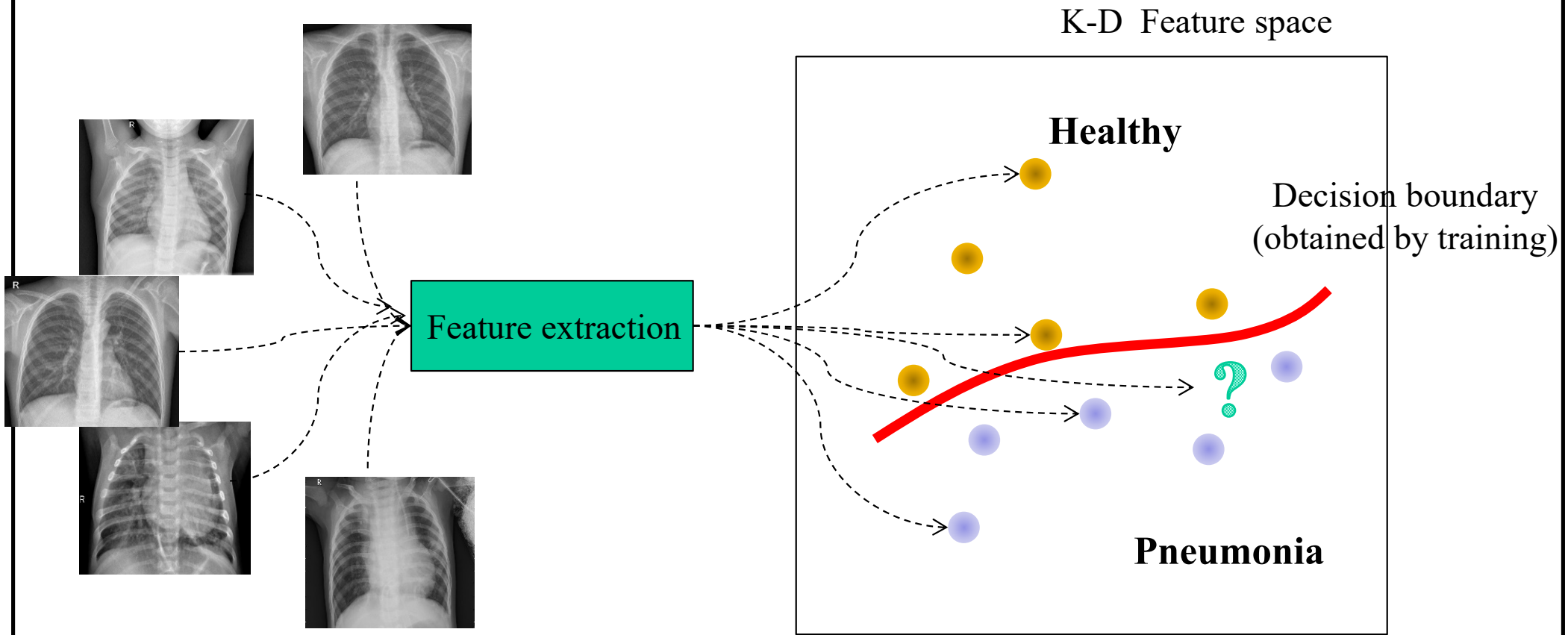
Chess X-Ray Pneumonia





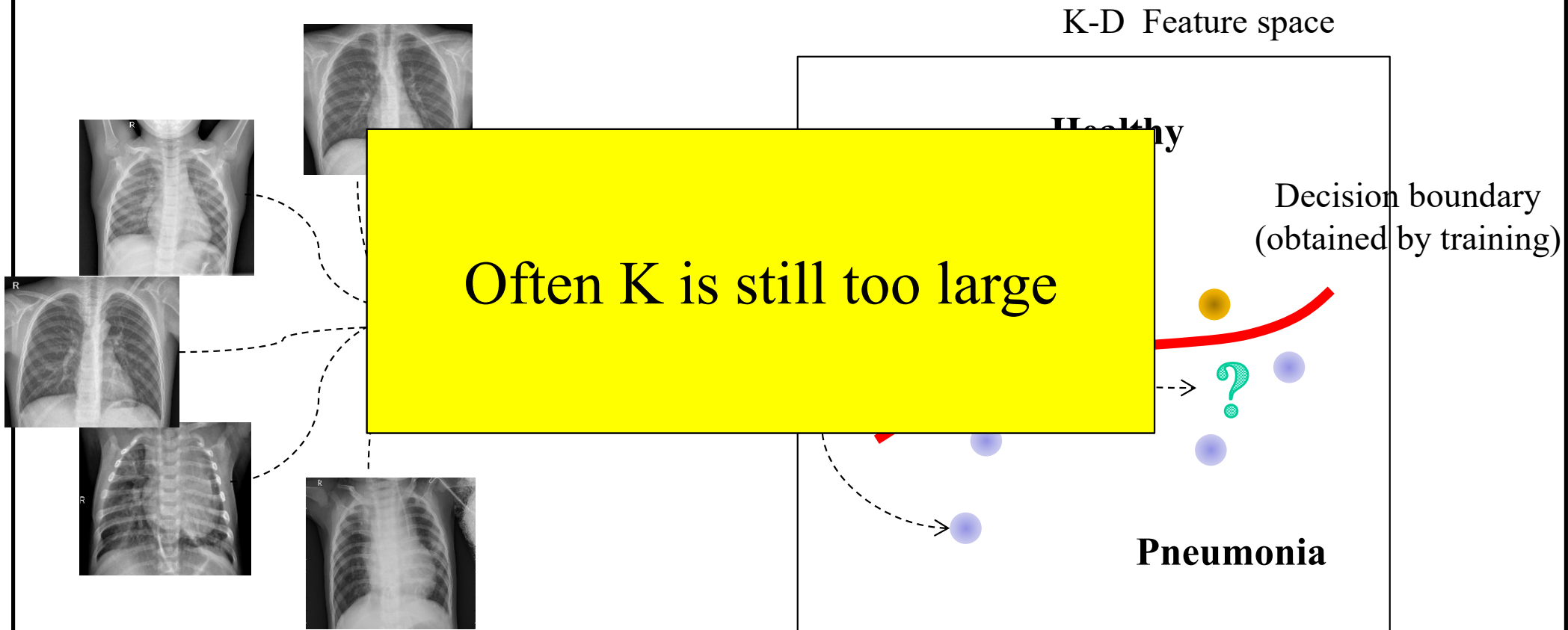
Very large feature spaces (like **88,704 dim**) are **problematic**.

Supervised learning

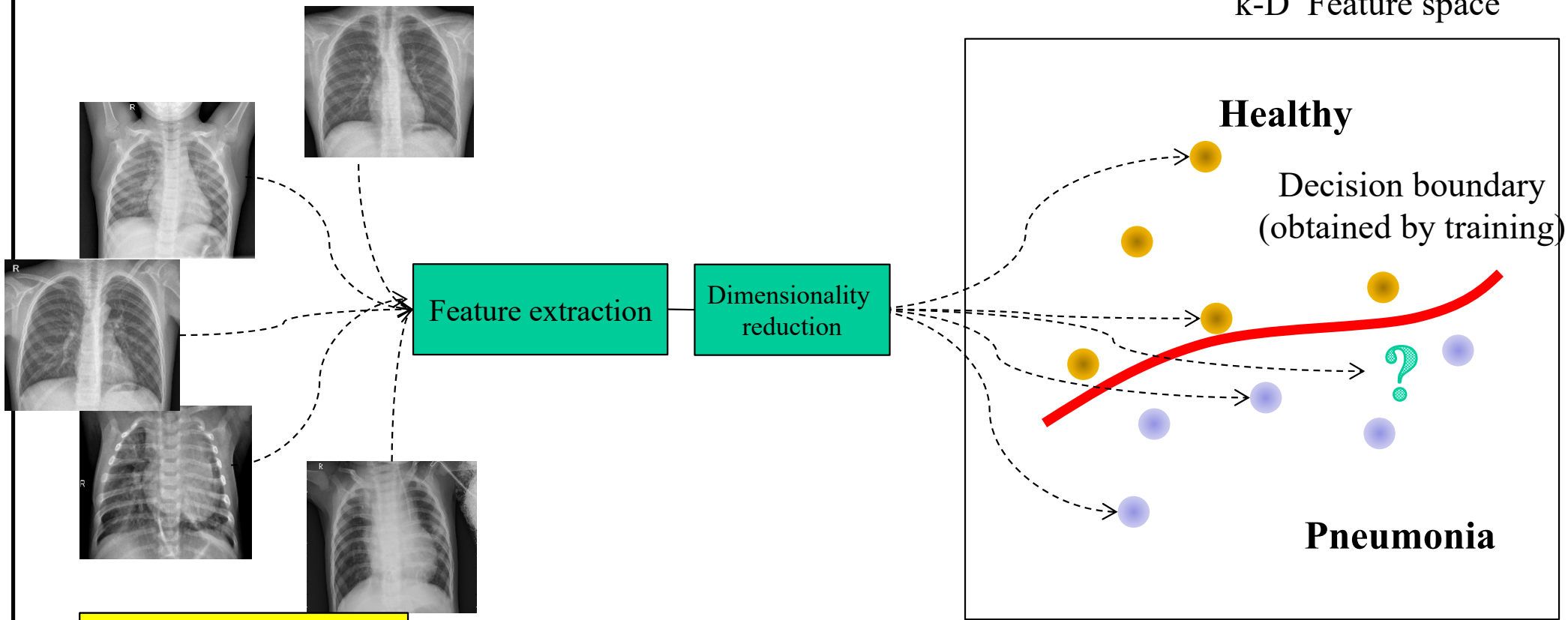


- Edge detection
- Histograms
- Filters of all kinds
- Moments
- Etc.

Supervised learning



Supervised learning

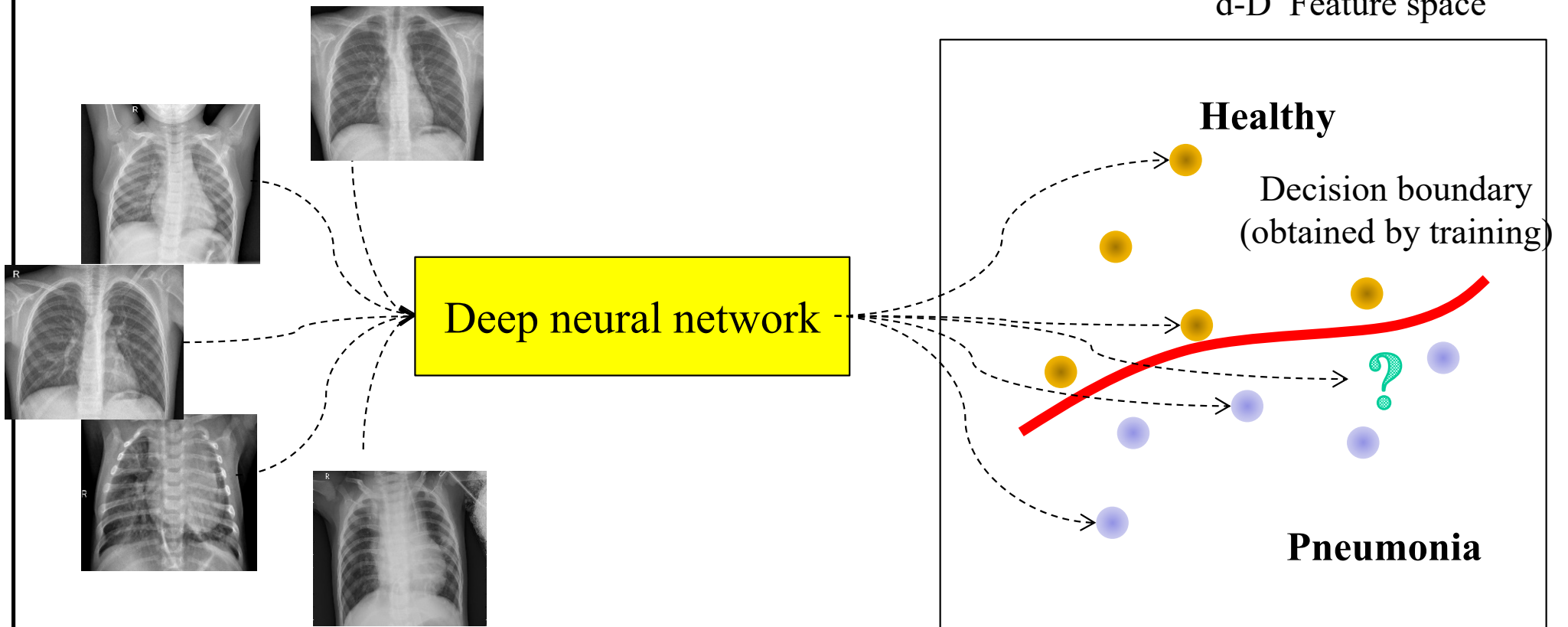


- Feature selection
- PCA
- Kernel PCA
- LLE
- ISOMAPS
- Etc.



Spoiler alert

In 2024...



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
- **Regression** : 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

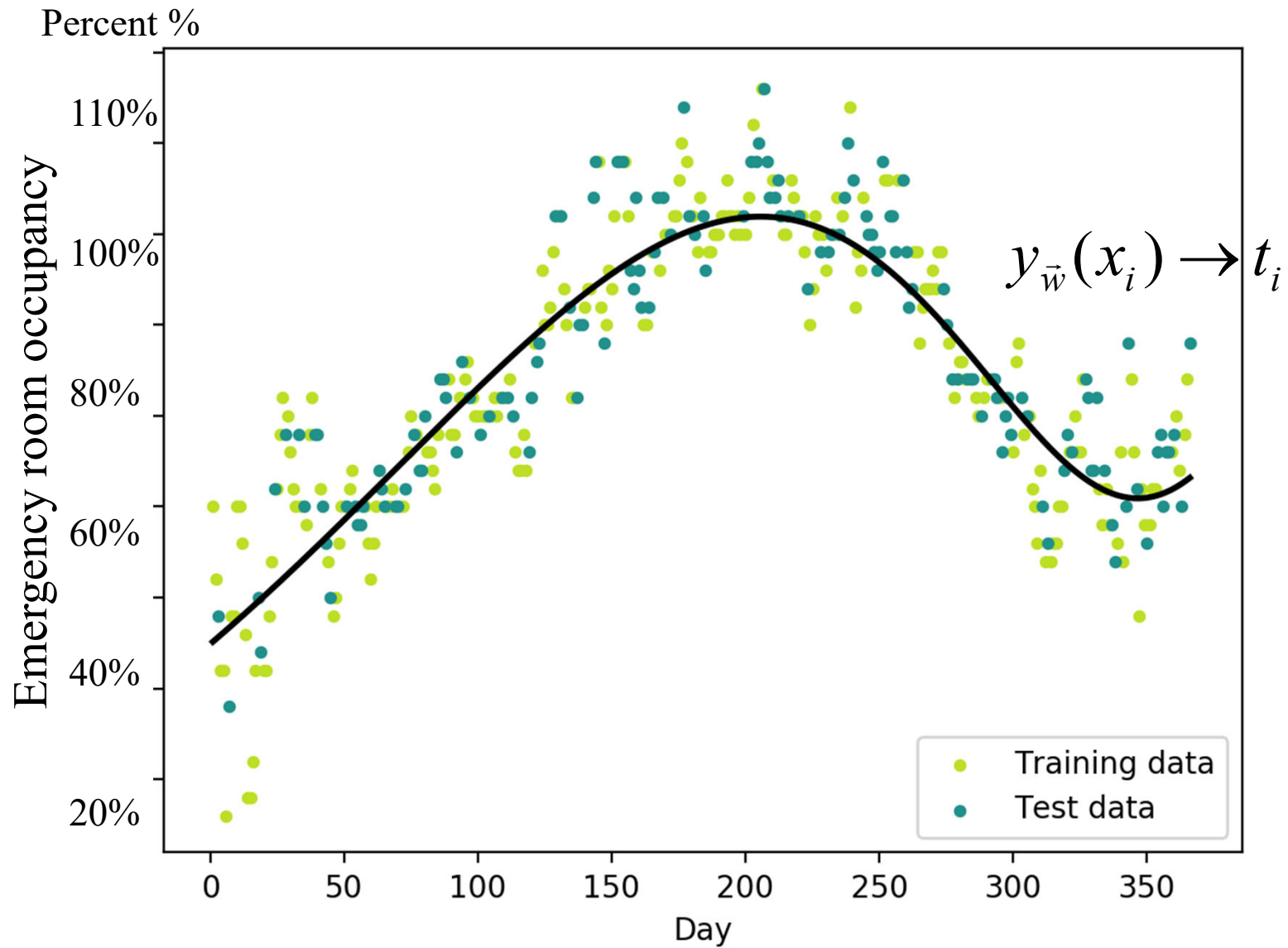
(Linear)

Medical Cost Personal Datasets



Example

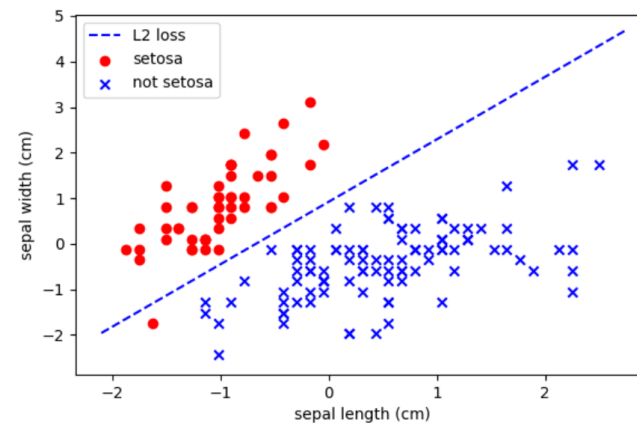
(Nonlinear)



Linear models



<https://rpubs.com/koki25ando/medicalcost>



<https://winder.ai/403-linear-classification/>

Deep neural nets

...

linear models



Vs
?

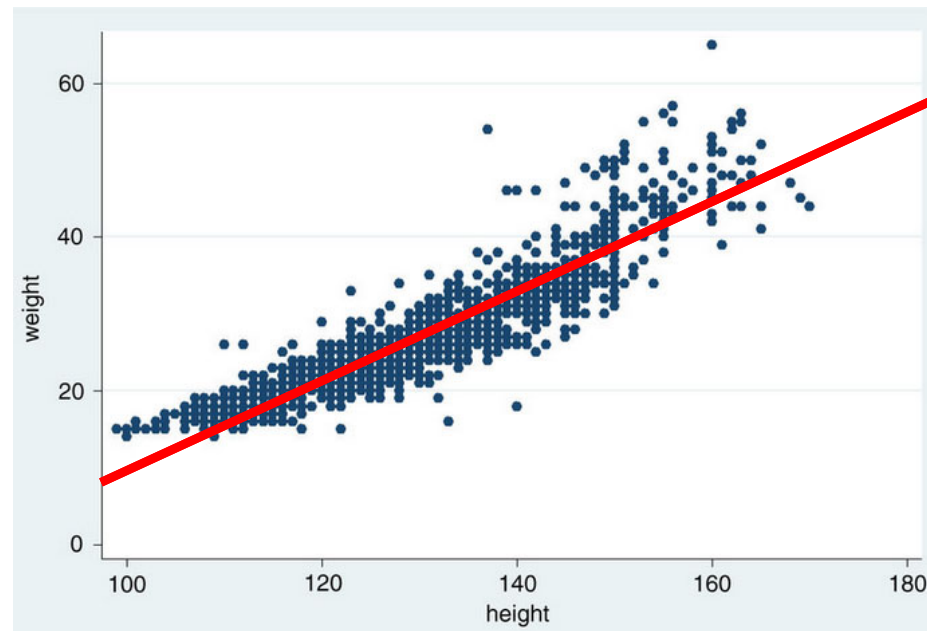


Linear models are to deep neural nets what transistors are to modern processors



Linear models are still relevant

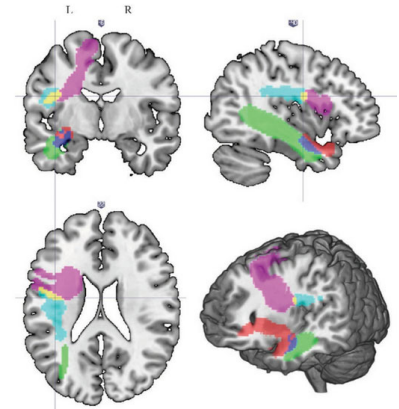
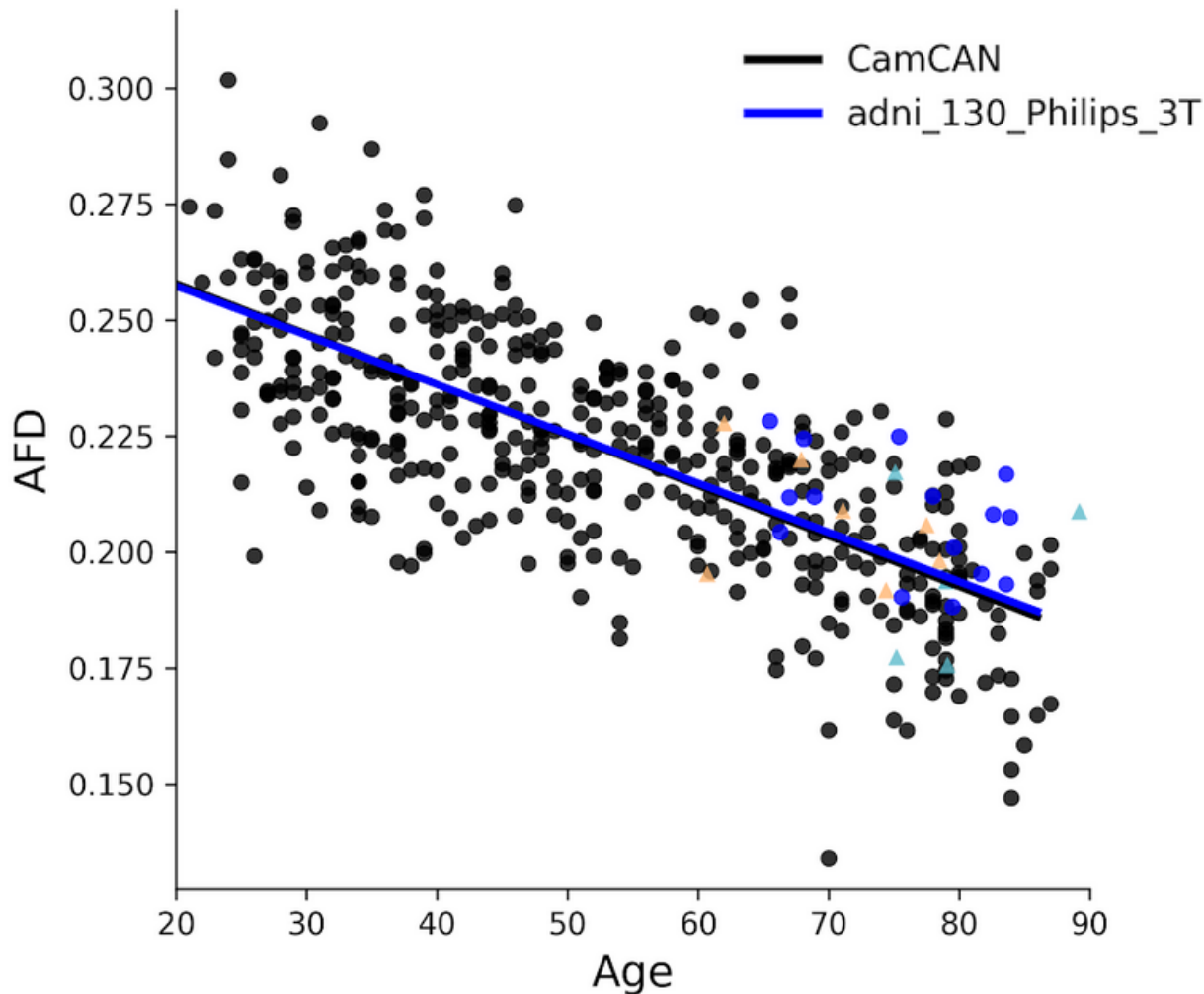
1,694 children surveyed in Tanzania.



Nordin P, Poggensee G, Mtweve S, Krantz I. **From a weighing scale to a pole: a comparison of two different dosage strategies in mass treatment of Schistosomiasis haematobium.** Glob Health Action. 2014

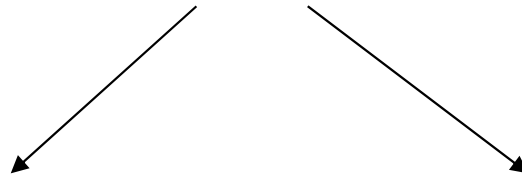
Linear models are still relevant

Apparent Fiber Density
in the white matter



<https://commons.wikimedia.org/>

Linear models



Classification

Regression

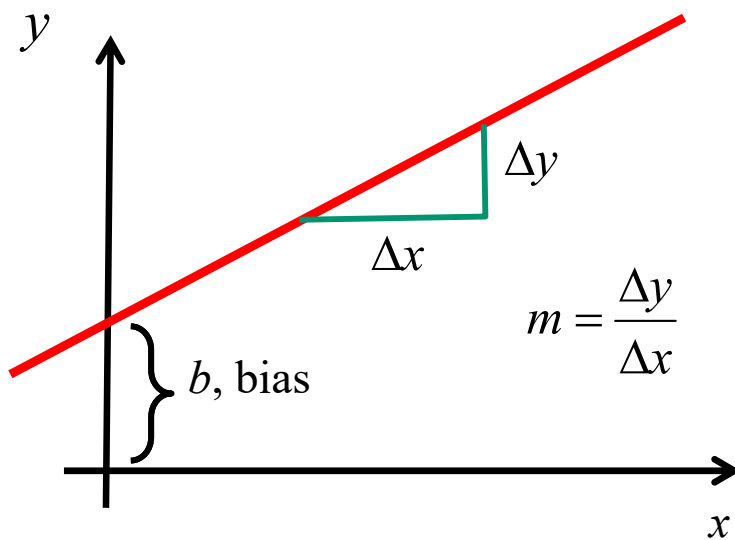
Linear models



Classification

Regression

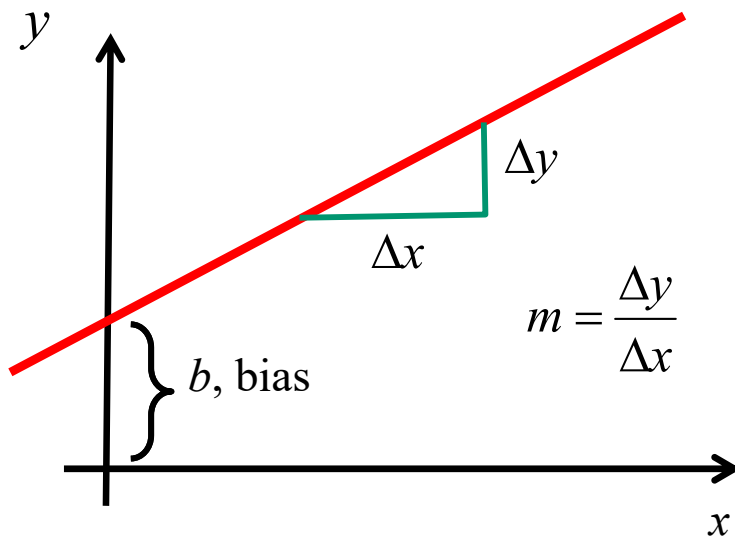
Definition ... a line!



$$y = mx + b$$

slope \uparrow bias \uparrow

Definition ... a line!



$$y = mx + b$$

$$y = \frac{\Delta y}{\Delta x} x + b$$

$$y\Delta x = \Delta y x + b\Delta x$$

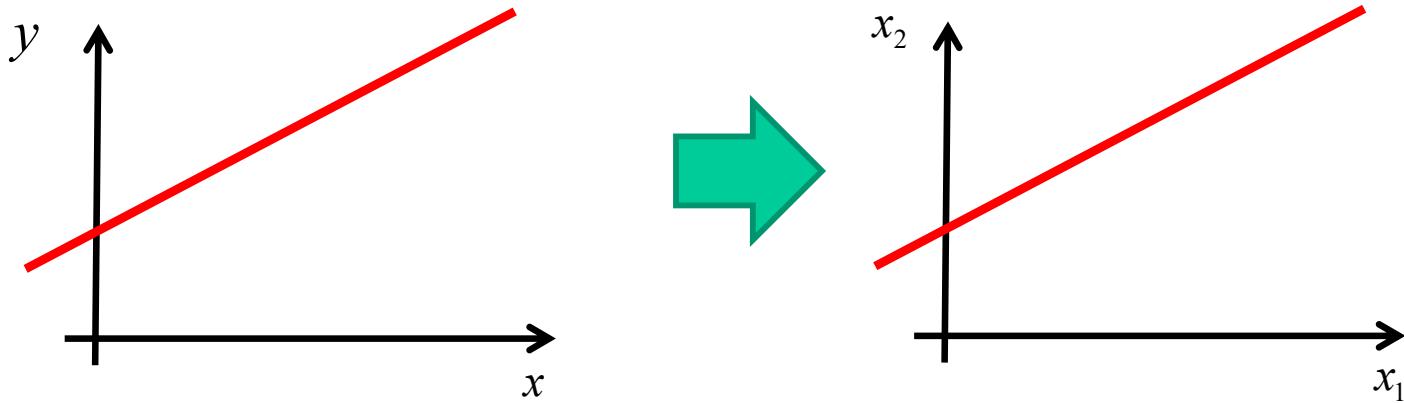
$$0 = \Delta y x - \Delta x y + b\Delta x$$

Rename variables

$$0 = \underbrace{\Delta yx}_{w_1} - \underbrace{\Delta xy}_{w_2} + \underbrace{b\Delta x}_{w_0}$$

Rename variables

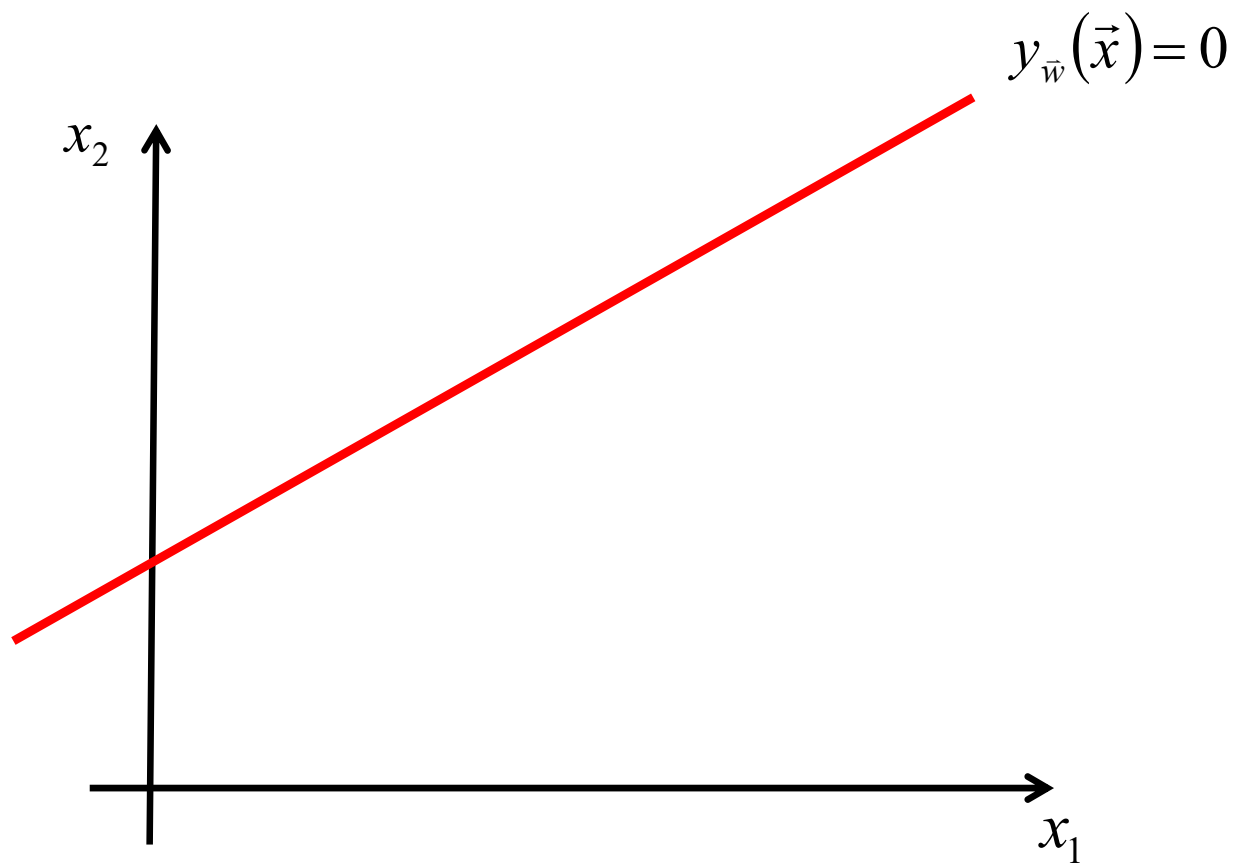
$$0 = w_1x + w_2y + w_0$$



$$0 = w_1x_1 + w_2x_2 + w_0$$

Implicit function

$$y_{\vec{w}}(\vec{x}) = w_1 x_1 + w_2 x_2 + w_0$$



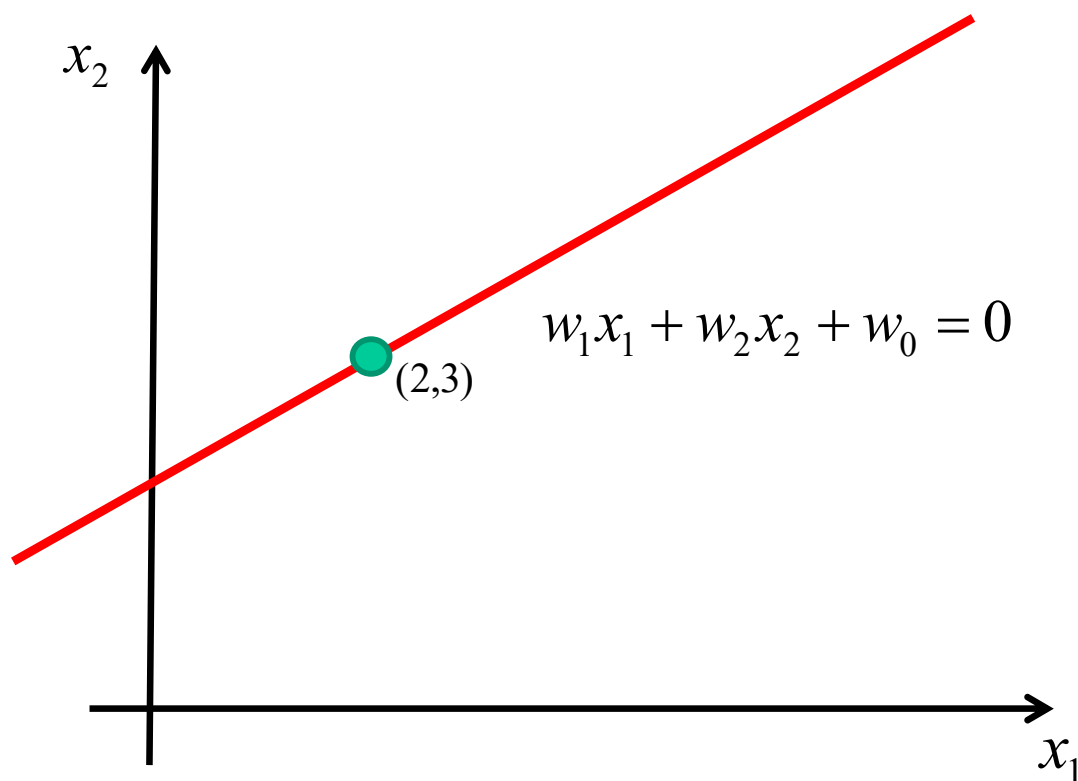
Implicit function

$$y_{\vec{w}}(\vec{x}) = w_1 x_1 + w_2 x_2 + w_0$$

$$w_1 = 1.0$$

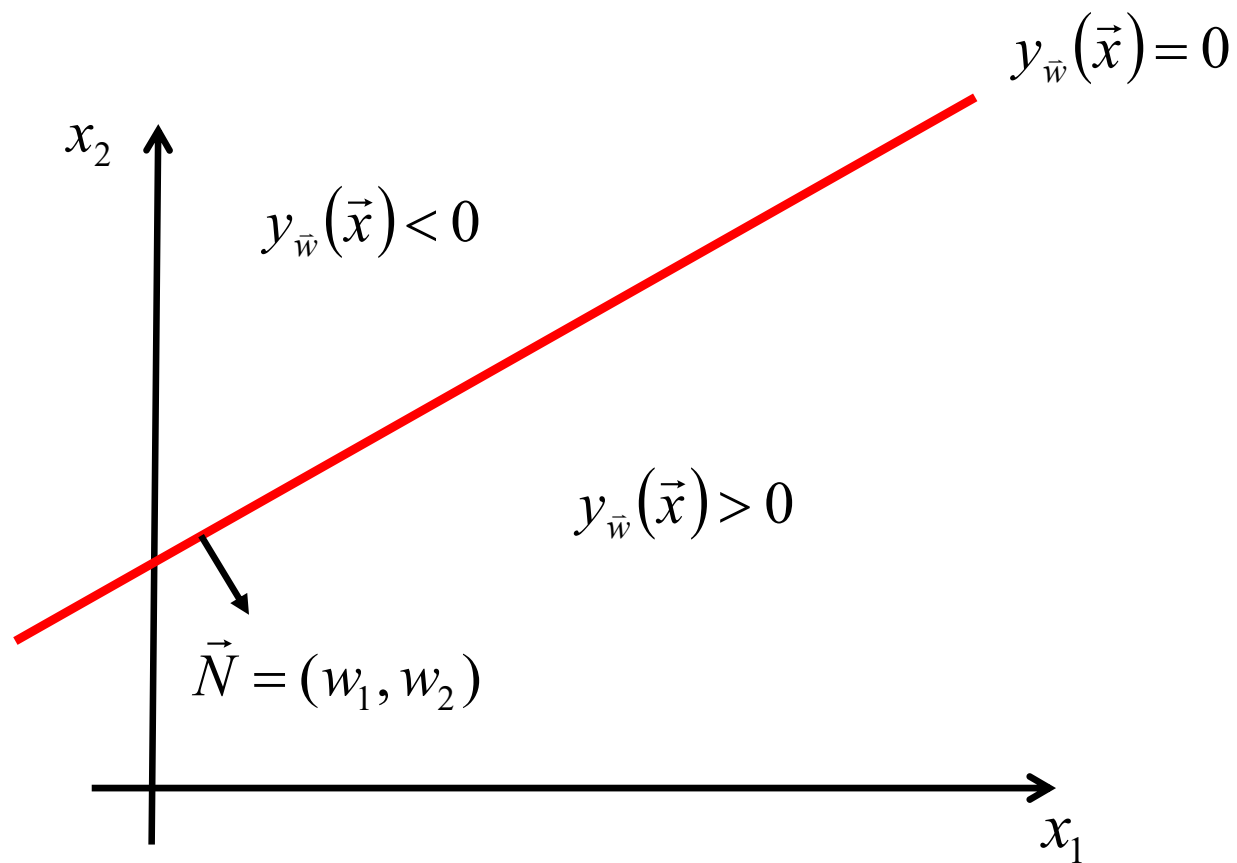
$$w_2 = -2.0$$

$$w_0 = 4.0$$



Classification function

$$y_{\vec{w}}(\vec{x}) = w_1 x_1 + w_2 x_2 + w_0$$



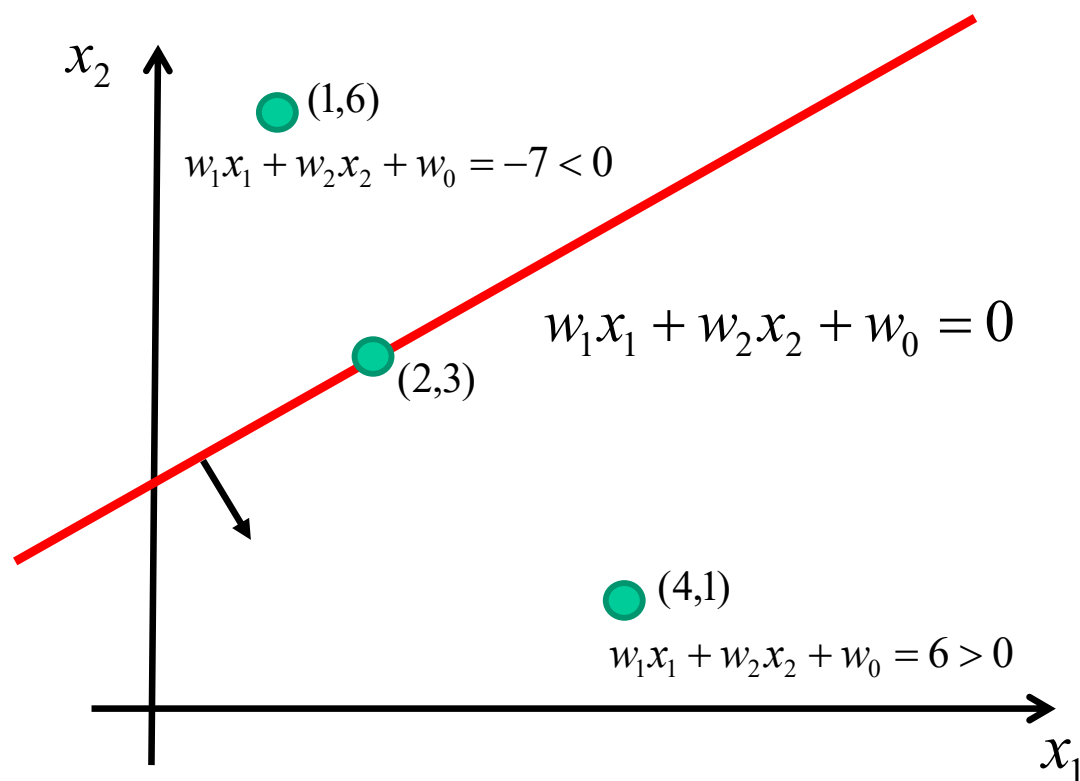
Classification function

$$y_{\vec{w}}(\vec{x}) = w_1 x_1 + w_2 x_2 + w_0$$

$$w_1 = 1.0$$

$$w_2 = -2.0$$

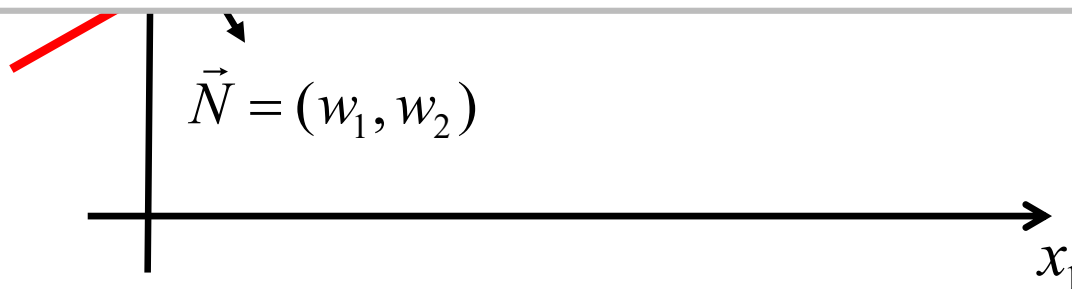
$$w_0 = 4.0$$



Implicit function

$$y_{\vec{w}}(\vec{x}) = w_0 + w_1 x_1 + w_2 x_2$$
$$= \underbrace{(w_0, w_1, w_2)}_{\vec{w}} \cdot \underbrace{(1, x_1, x_2)}_{\vec{x}'}$$

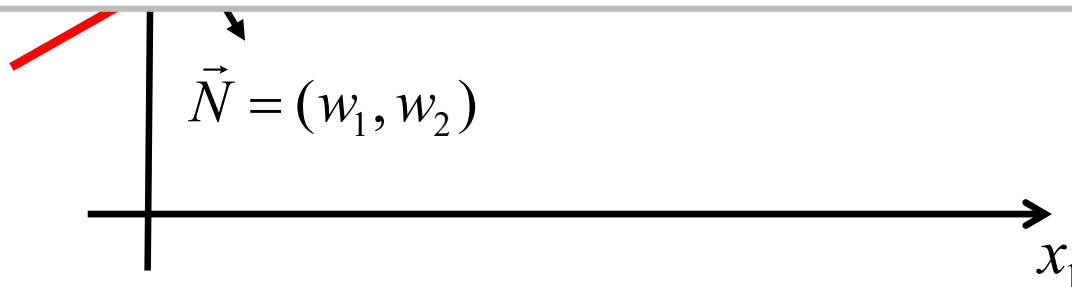
**DOT
product**



Implicit function

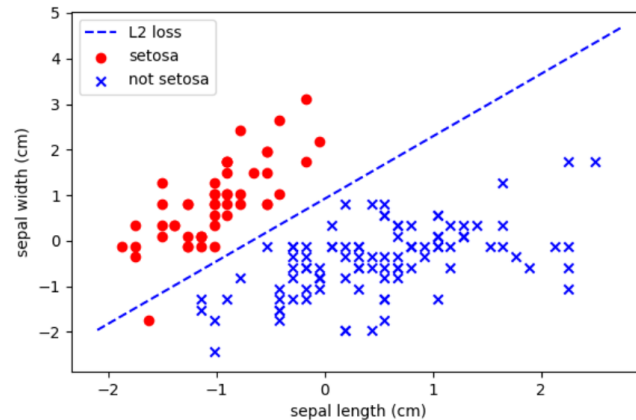
$$\begin{aligned}y_{\vec{w}}(\vec{x}) &= w_1 x_1 + w_2 x_2 + w_0 \\ &= (w_0, w_1, w_2) \cdot (1, x_1, x_2) \\ &= \vec{w}^T \vec{x}'\end{aligned}$$

**DOT
product**

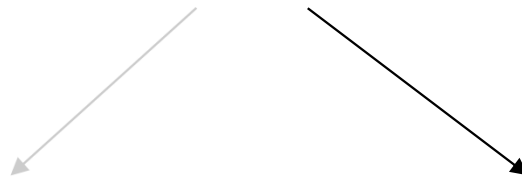


Linear classifier = dot product with bias included

$$y_{\vec{w}}(\vec{x}) = \vec{w}^T \vec{x} = \begin{cases} > 0 & \text{if in front} \\ < 0 & \text{otherwise} \end{cases}$$



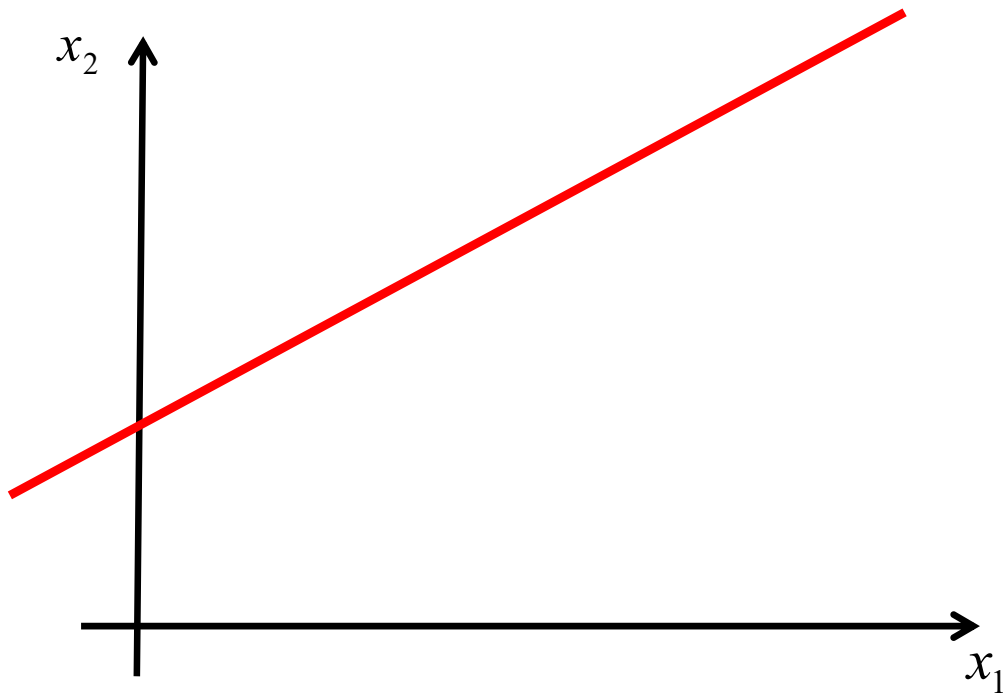
Linear models



Classification

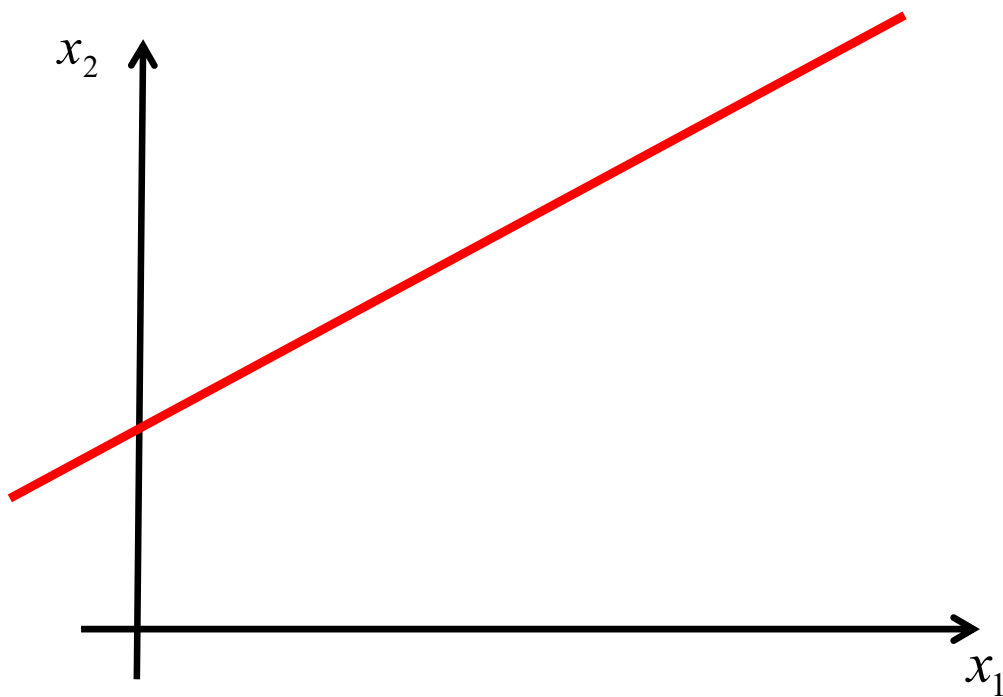
Regression

Classification : implicit function



$$0 = w_1 x_1 + w_2 x_2 + w_0$$

Explicit function



$$0 = w_1 x_1 + w_2 x_2 + w_0$$

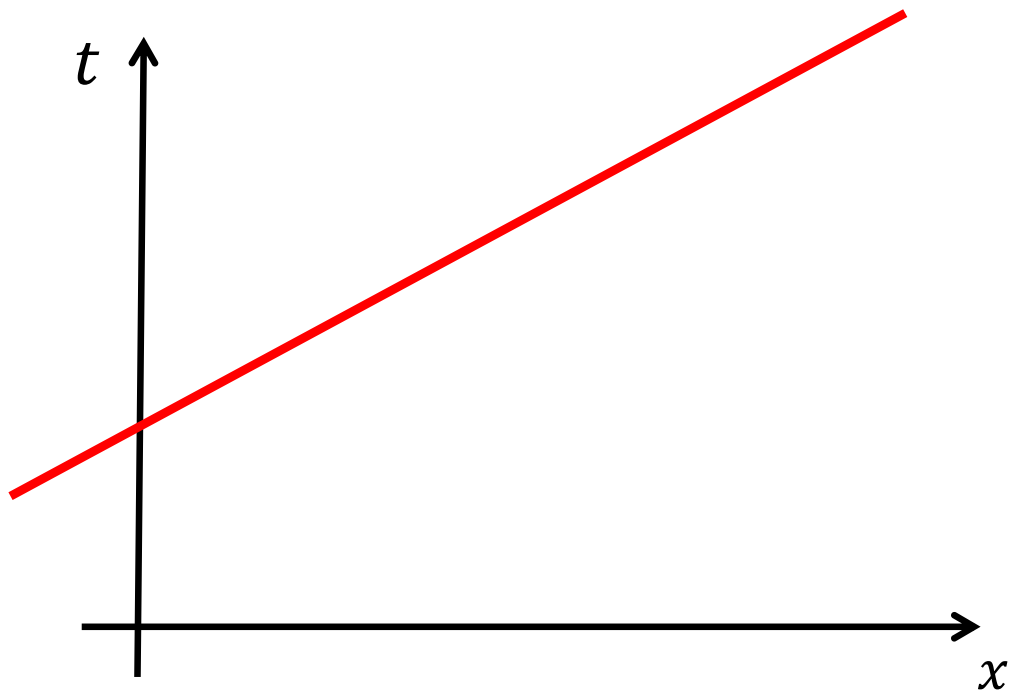
$$x_2 = \frac{1}{-w_2} (w_1 x_1 + w_0)$$

$$x_2 = w'_1 x_1 + w'_0$$

Rename, t for target

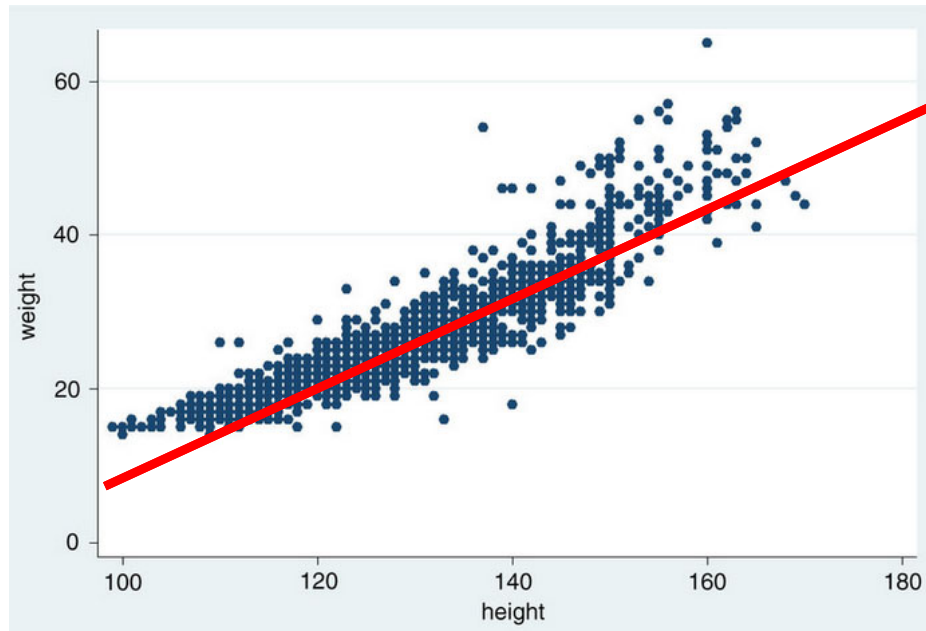
$$x_2 \rightarrow t$$

$$t = w_1 x + w_0$$



Goal : predict the target t given x

1,694 children surveyed in Tanzania.



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

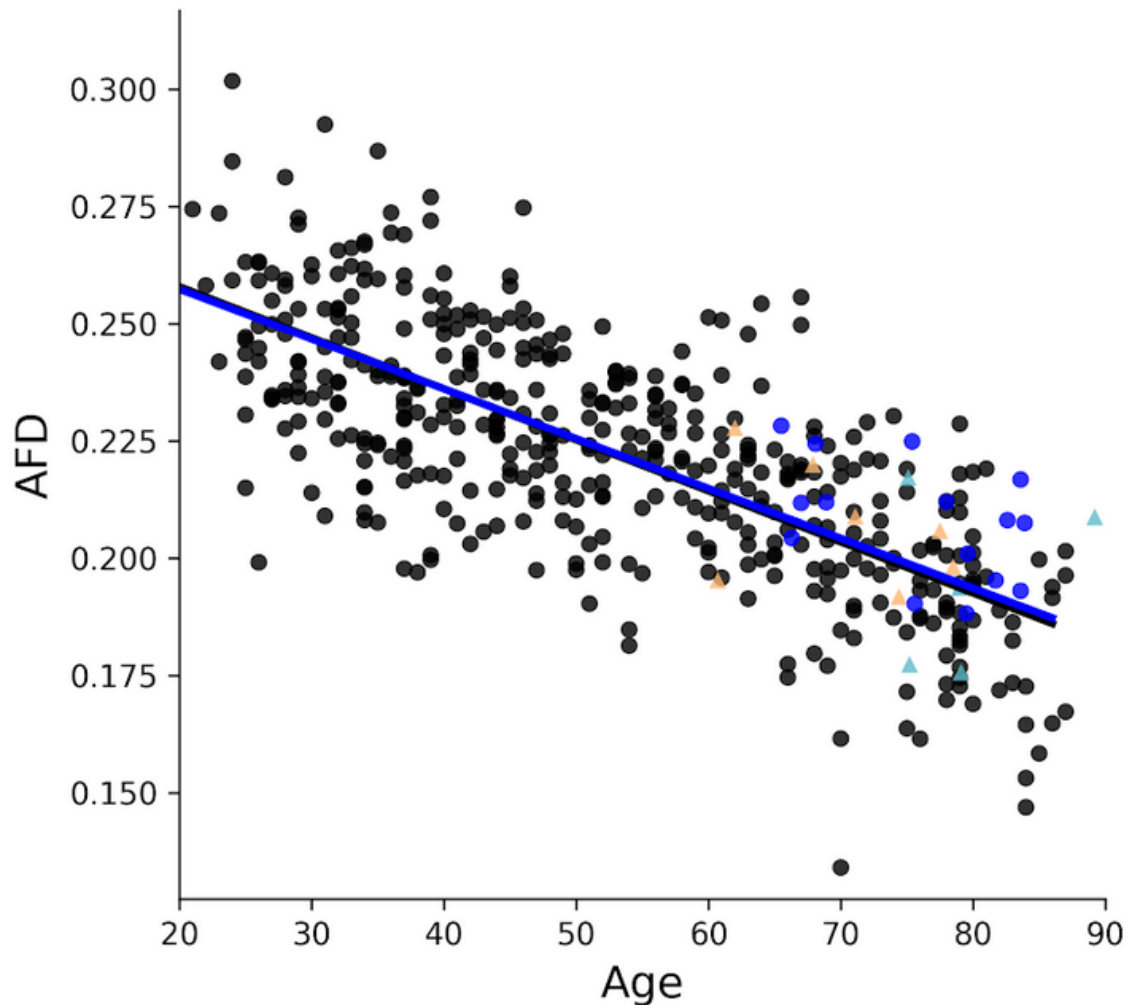
$$y_{\vec{w}}(x_i) = t_i \quad \forall i$$

Nordin P, Poggensee G, Mtweve S, Krantz I. **From a weighing scale to a pole: a comparison of two different dosage strategies in mass treatment of Schistosomiasis haematobium.** Glob Health Action. 2014

A line : 1D regression

Example

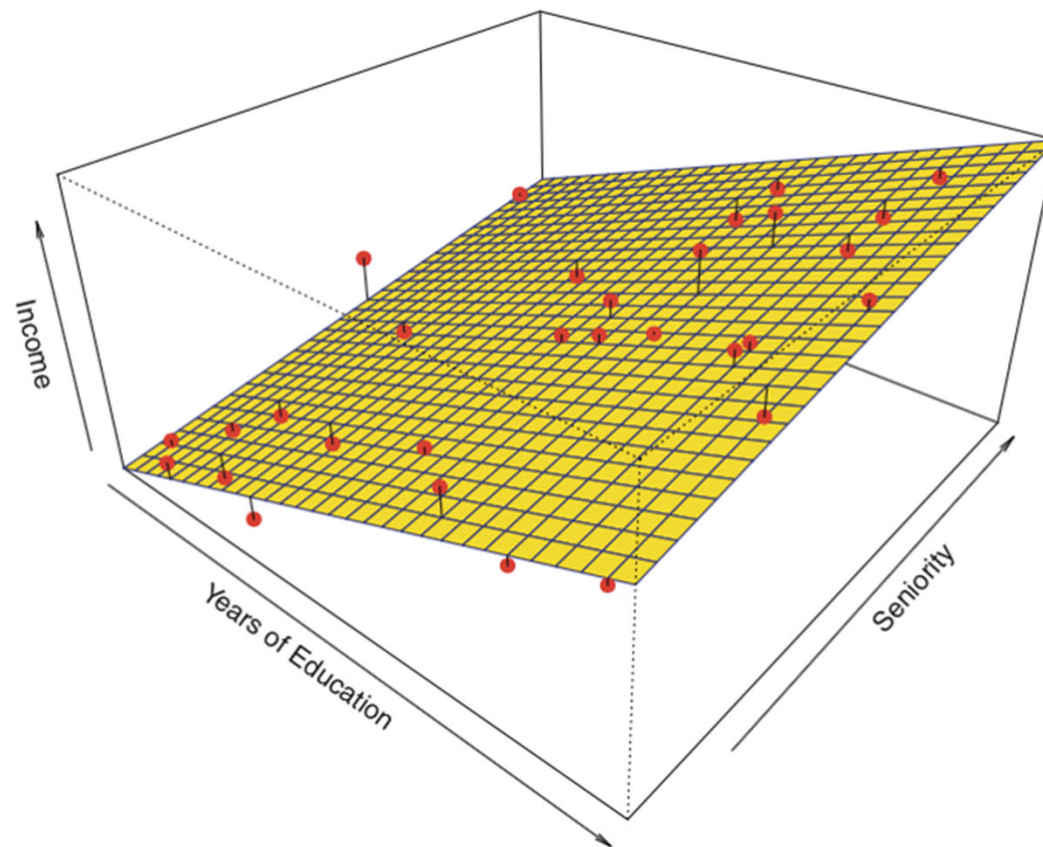
$$y_{\vec{w}}(x) = w_0 + w_1 x$$



A plane : 2D regression

Example

$$y_{\vec{w}}(\vec{x}) = w_0 + w_1x_1 + w_2x_2$$



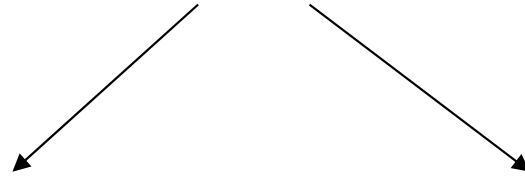
A hyper plane : dD regression

$$y_{\vec{w}}(\vec{x}) = w_0 + w_1x_1 + w_2x_2 + \cdots + w_dx_d$$

$$= \vec{w}^T \vec{x}$$

Dot product

Linear models



Classification

$$y_{\vec{w}}(\vec{x}) = \vec{w}^T \vec{x} = \begin{cases} > 0 & \text{if in front} \\ < 0 & \text{otherwise} \end{cases}$$

Regression

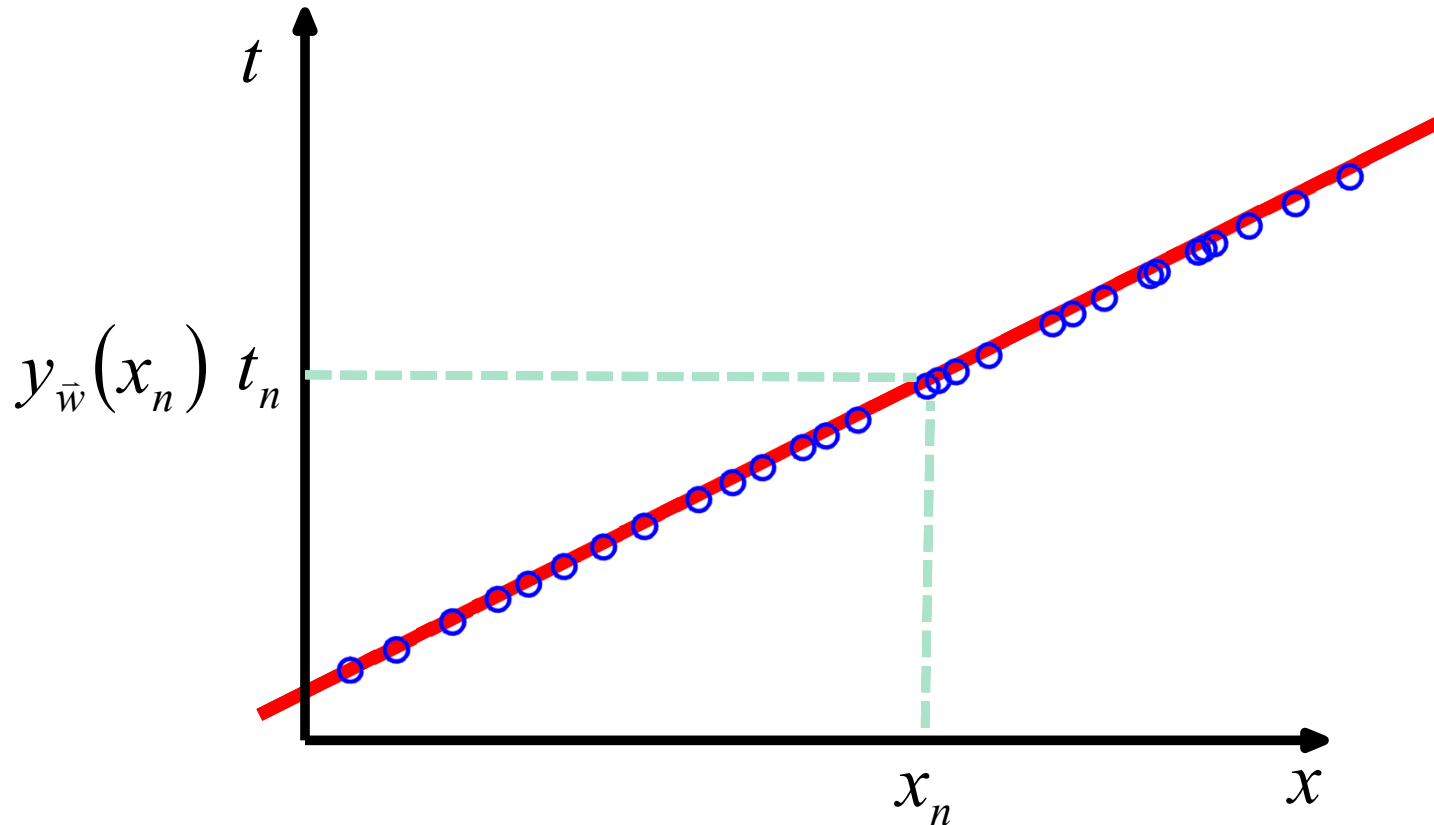
$$y_{\vec{w}}(\vec{x}) = \vec{w}^T \vec{x} = t$$

Problem to solve

Given a training example

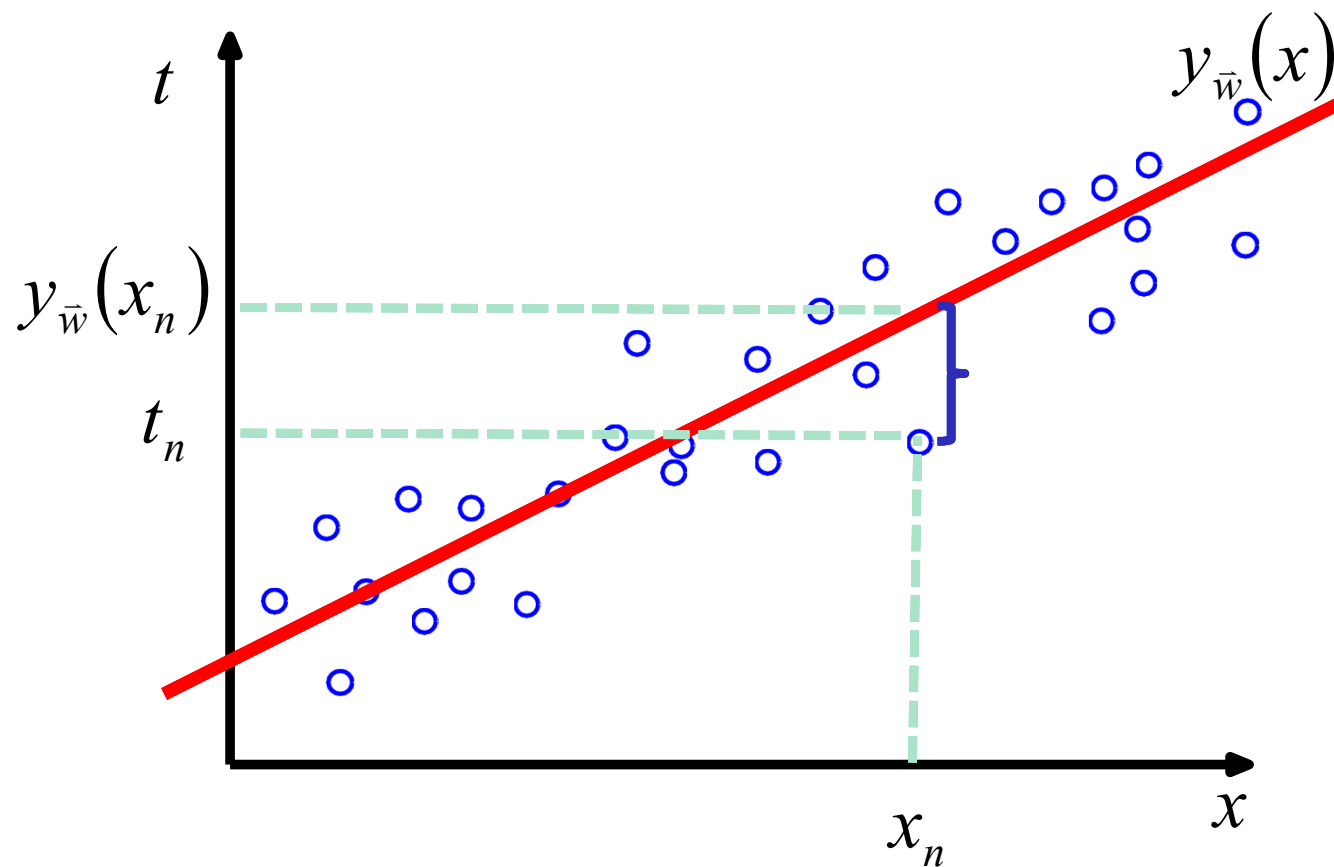
$$D = \{(x_1, t_1), (x_2, t_2), \dots, (x_N, t_N)\}$$

Ideally, we wish $y_{\bar{w}}(x_i) = t_i$



Problem to solve

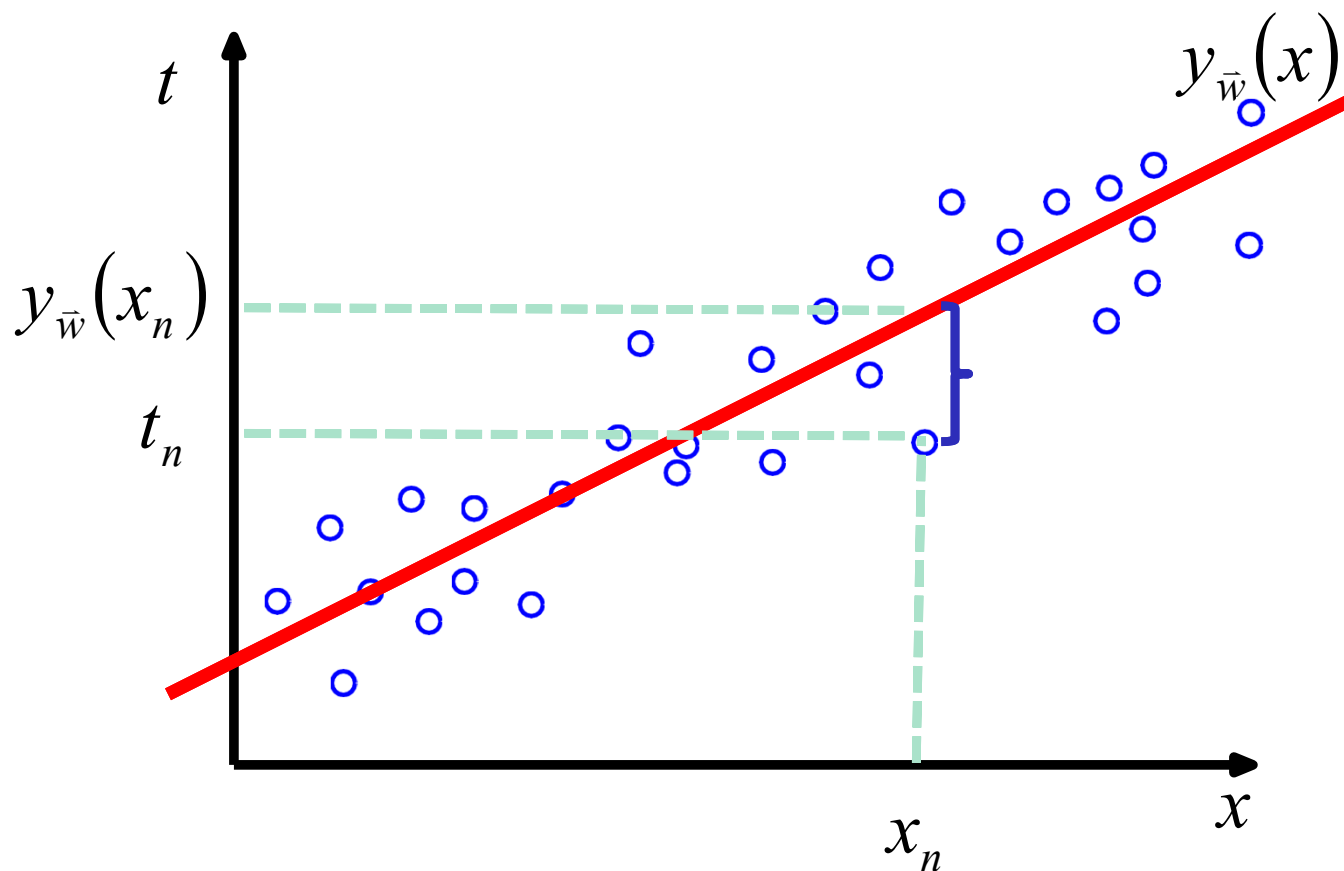
Unfortunately, real data are **noisy**



Here the goal is to make **small mistakes**.

Problem to solve

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



Problem to solve

$$\vec{w} = \arg \min_{\vec{w}} \sum_{n=1}^N (\vec{w}^T \vec{x}_n - t_n)^2$$

If the data is **linear** + the noise is **Gaussian**,
the best possible weights are those
minimizing this function

Problem to solve

$$\vec{w} = \arg \min_{\vec{w}} \underbrace{\sum_{n=1}^N (\vec{w}^T \vec{x}_n - t_n)^2}_{\mathcal{L}_D(\vec{w})}$$

the « best » \vec{w} is the one for which the gradient is **zero**

$$\nabla_{\vec{w}} \mathcal{L}_D(\vec{w}) = \sum_{n=1}^N 2(\vec{w}^T \vec{x}_n - t_n) \vec{x}_n^T = 0$$

$$\vec{w}^T \sum_{n=1}^N \vec{x}_n \vec{x}_n^T - \sum_{n=1}^N t_n \vec{x}_n^T = 0$$

Problem to solve

$$\vec{w}^T \sum_{n=1}^N \vec{x}_n \vec{x}_n^T - \sum_{n=1}^N t_n \vec{x}_n^T = 0$$

By isolating \vec{w} , we get

$$\vec{w} = (X^T X)^{-1} X^T T$$

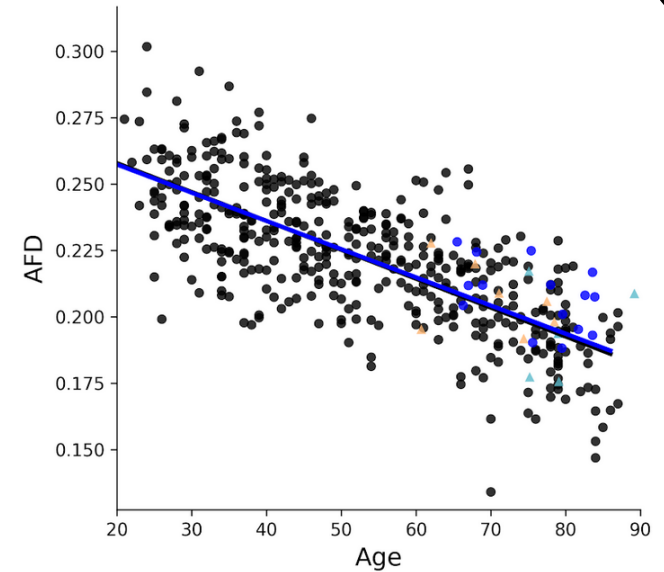
where

$$X = \begin{pmatrix} 1 & x_{1,1} & \cdots & x_{1,d} \\ 1 & x_{2,1} & \cdots & x_{2,d} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N,1} & \cdots & x_{N,d} \end{pmatrix} \quad T = \begin{pmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{pmatrix}$$

For a 1D regression

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

$$\vec{w} = (X^T X)^{-1} X^T T$$

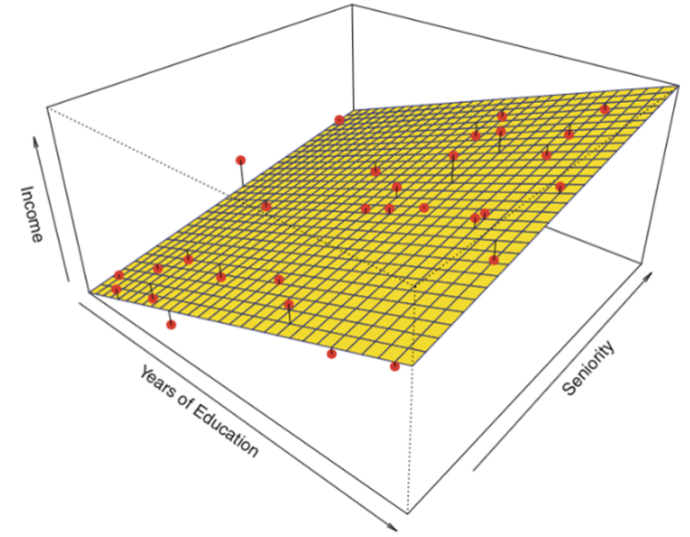


where

$$X = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_N \end{pmatrix} \quad T = \begin{pmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{pmatrix}$$

For a 2D regression

$$y_{\vec{w}}(\vec{x}) = w_0 + w_1 x_1 + w_2 x_2$$



$$\vec{w} = (X^T X)^{-1} X^T T$$

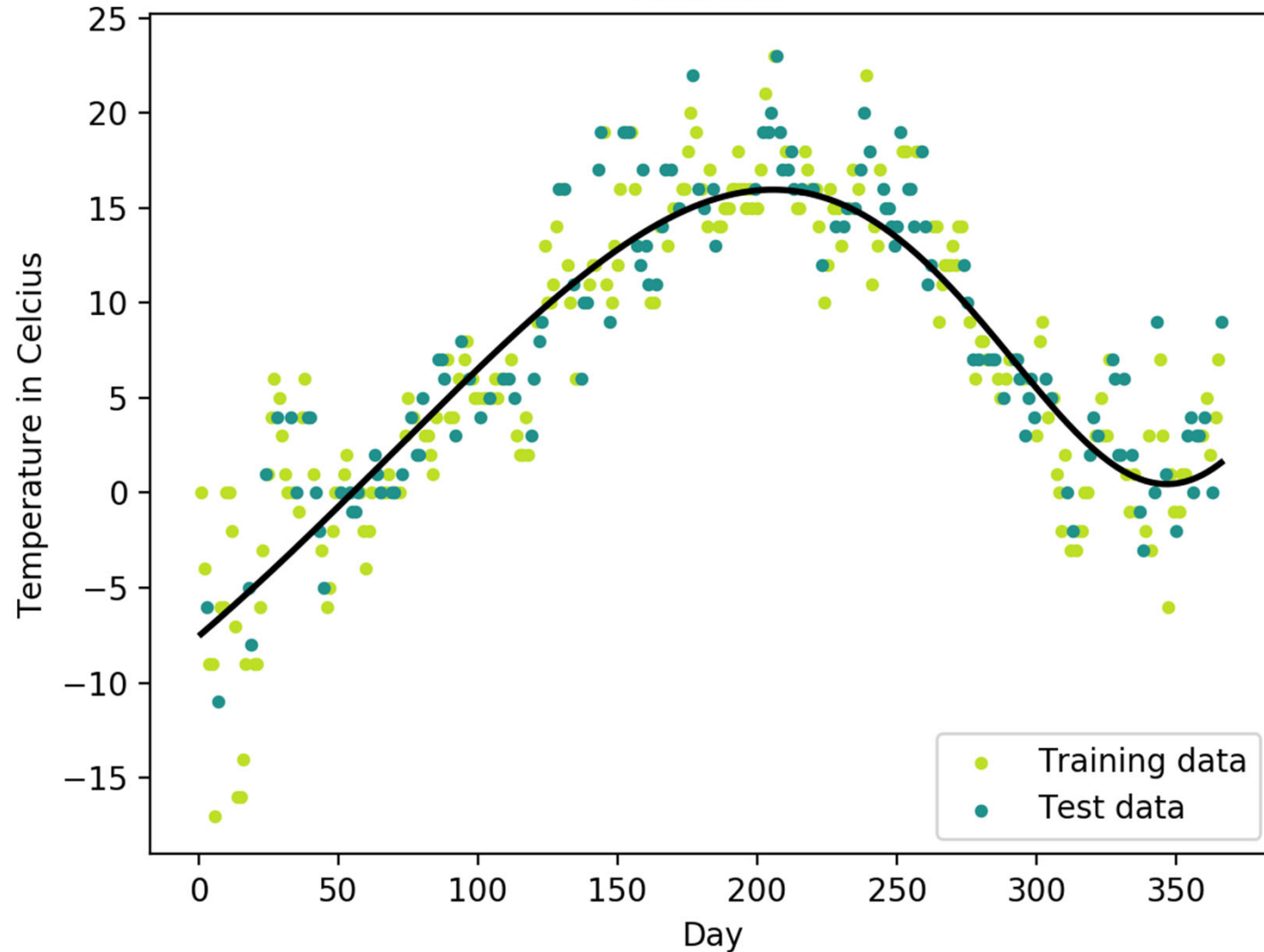
where

$$X = \begin{pmatrix} 1 & x_{1,1} & x_{1,2} \\ 1 & x_{2,1} & x_{2,2} \\ \vdots & \vdots & \vdots \\ 1 & x_{N,1} & x_{N,2} \end{pmatrix} \quad T = \begin{pmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{pmatrix}$$

What about non-linear data?

Polynomial Ridge Regression

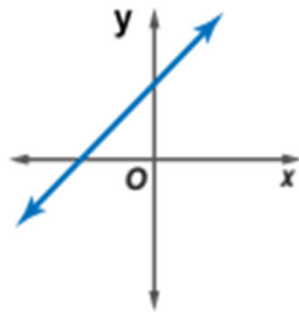
MSE: 9.58



What about non-linear data?

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

Linear function
Degree 1



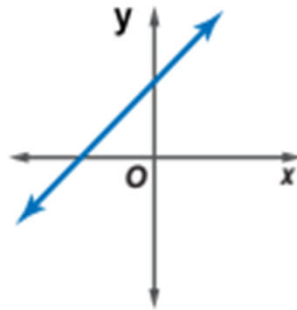
<http://www.math.glencoe.com/>

What about non-linear data?

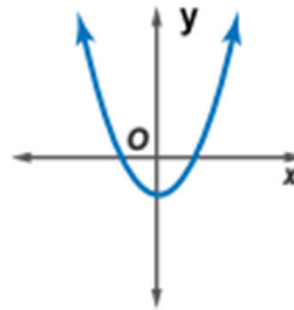
$$y_{\vec{w}}(\vec{x}) = w_0 + w_1x$$

$$y_{\vec{w}}(\vec{x}) = w_0 + w_1x + w_2x^2$$

Linear function
Degree 1



Quadratic function
Degree 2

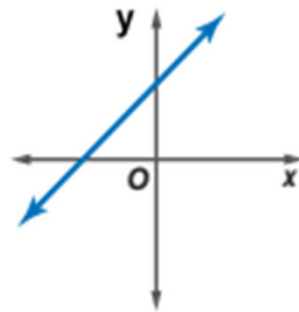


What about non-linear data?

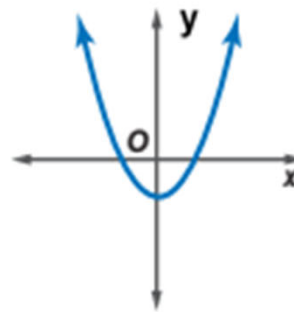
$$y_{\vec{w}}(\vec{x}) = w_0 + w_1x$$

$$y_{\vec{w}}(\vec{x}) = w_0 + w_1x + w_1x^2$$

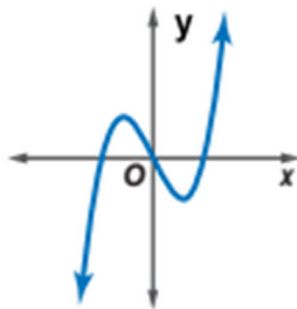
Linear function
Degree 1



Quadratic function
Degree 2



Cubic function
Degree 3



$$y_{\vec{w}}(\vec{x}) = w_0 + w_1x + w_1x^2 + w_1x^3$$

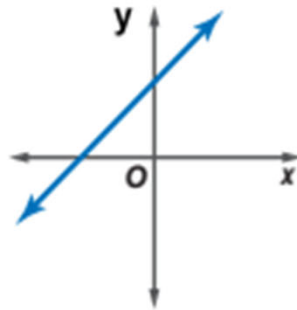
<http://www.math.glencoe.com/>

What about non-linear data?

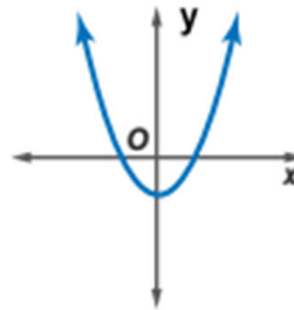
$$y_{\vec{w}}(\vec{x}) = w_0 + w_1x$$

$$y_{\vec{w}}(\vec{x}) = w_0 + w_1x + w_1x^2$$

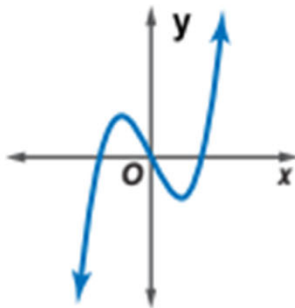
Linear function
Degree 1



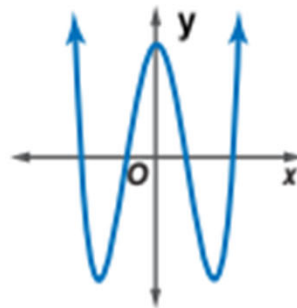
Quadratic function
Degree 2



Cubic function
Degree 3



Quartic function
Degree 4



$$y_{\vec{w}}(\vec{x}) = w_0 + w_1x + w_1x^2 + w_1x^3 + w_1x^4$$

$$y_{\vec{w}}(\vec{x}) = w_0 + w_1x + w_1x^2 + w_1x^3$$

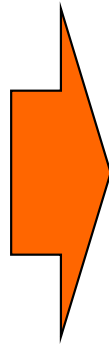
<http://www.math.glencoe.com/>

Basis function

Example: Instead of a **1D regression**, lets do a **4D regression**

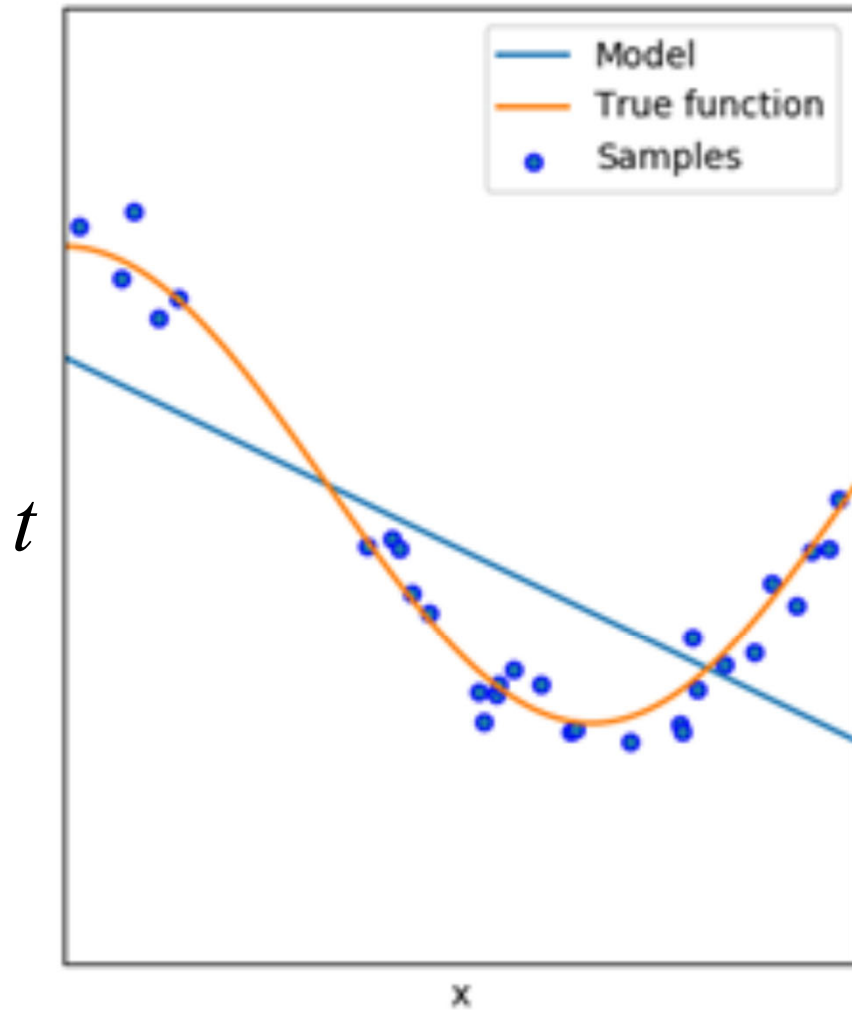
$$\varphi(x) \rightarrow (x, x^2, x^3, x^4)$$

$$y_{\vec{w}}(x) = w_0 + w_1x$$



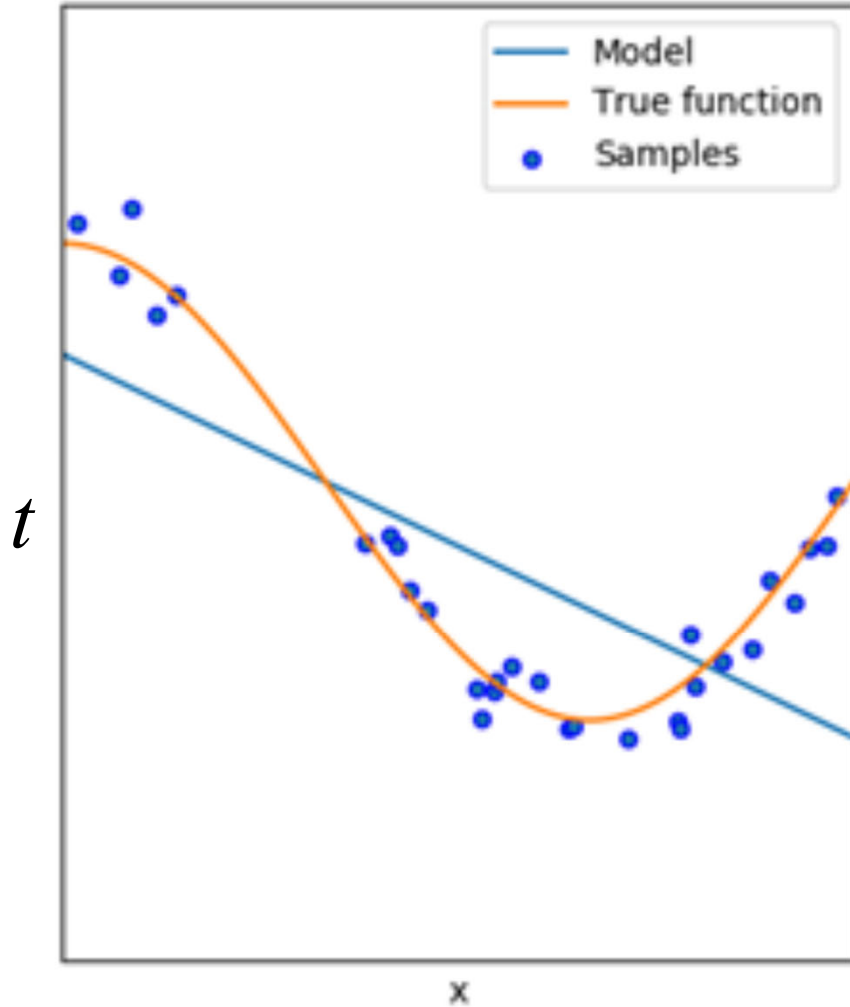
$$\begin{aligned} y_{\vec{w}}(x) &= w_0 + w_1x + w_2x^2 + w_3x^3 + w_4x^4 \\ &= w_0 + \sum_{i=1}^4 w_i\varphi_i(x) \end{aligned}$$

Degree 1
MSE = 4.08e-01(+/- 4.25e-01)



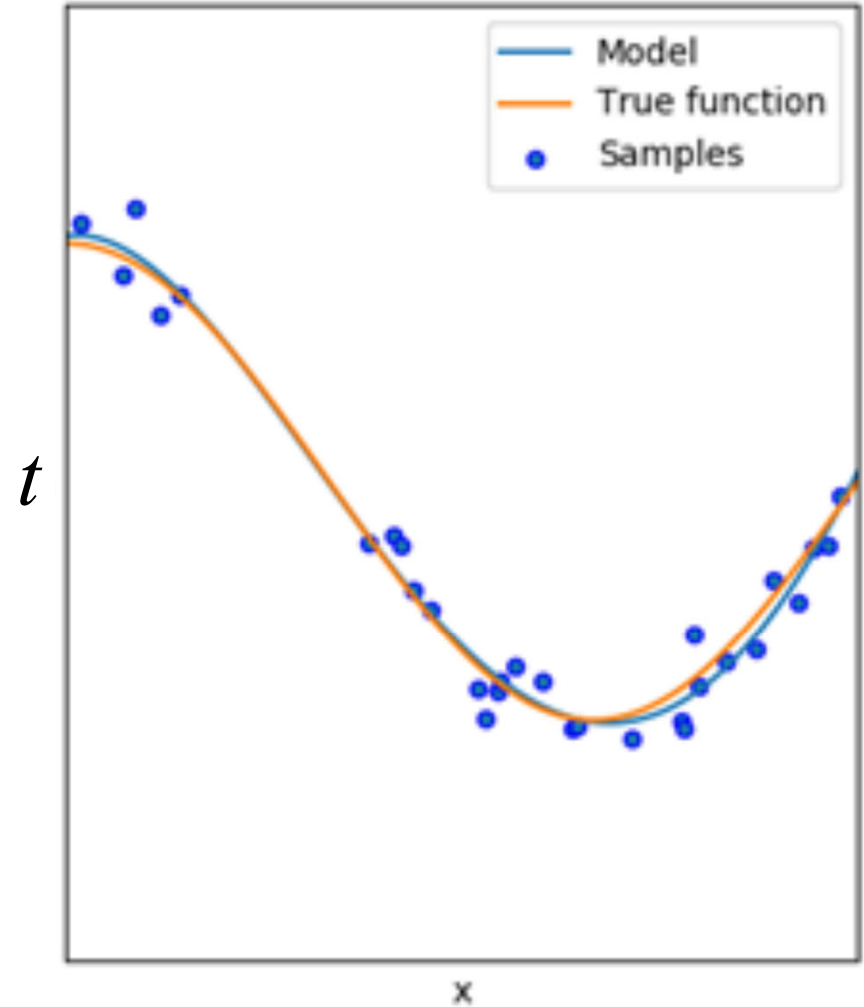
$$y_{\vec{w}}(x) = w_0 + w_1 x$$

Degree 1
MSE = 4.08e-01(+/- 4.25e-01)



$$y_{\vec{w}}(x) = w_0 + w_1x$$


Degree 4
MSE = 4.32e-02(+/- 7.08e-02)




$$y_{\vec{w}}(x) = w_0 + \sum_{i=1}^4 w_i \phi_i(x)$$

Basis function

$$y_{\vec{w}}(\vec{x}) = \sum_{i=0}^{d-1} w_i \varphi_i(\vec{x})$$

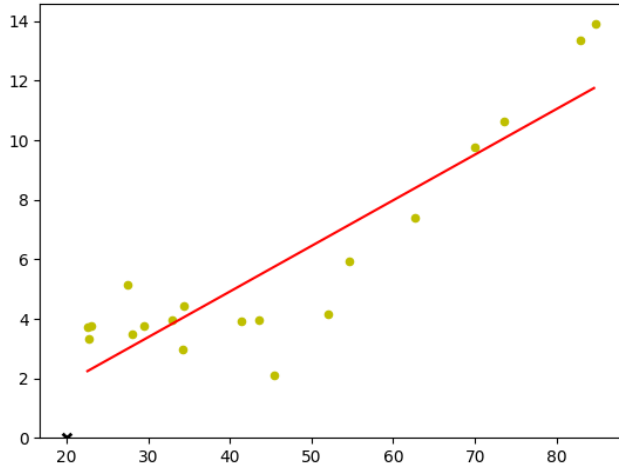
parameter 

Basis function 

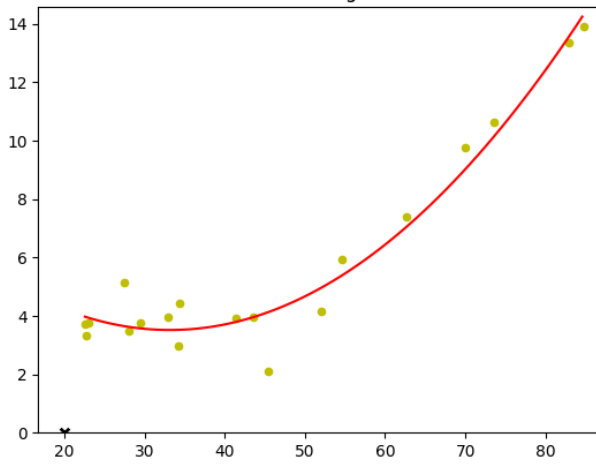
where $\varphi_0(\vec{x}) = 1$

Regression

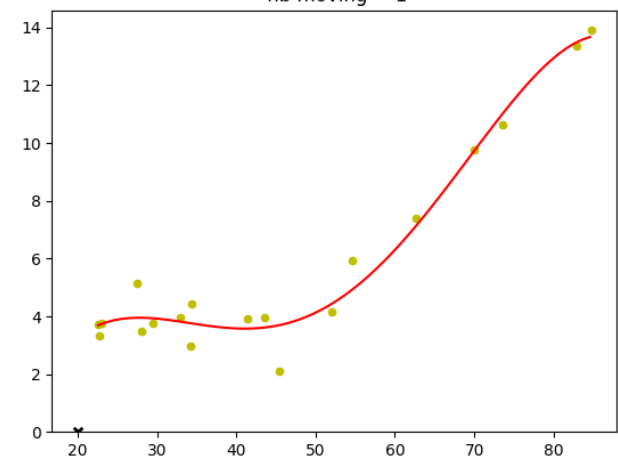
Polynom degree : 1
nb target = 19
nb moving = 1



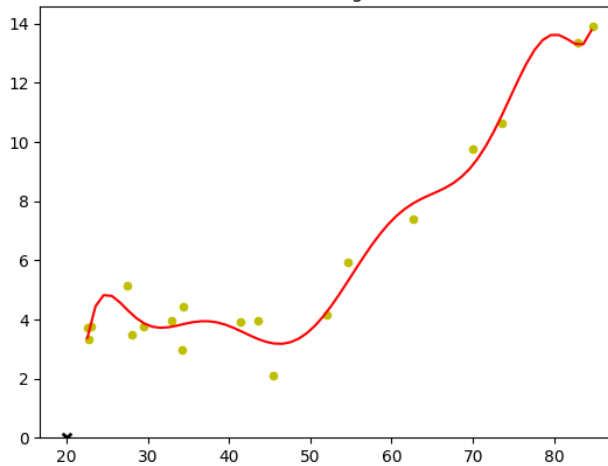
Polynom degree : 2
nb target = 19
nb moving = 1



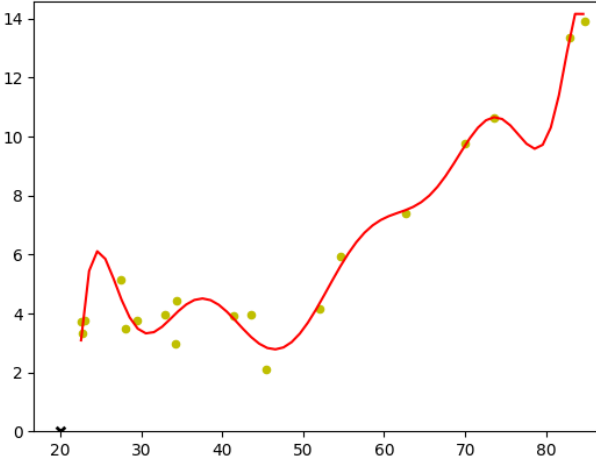
Polynom degree : 4
nb target = 19
nb moving = 1



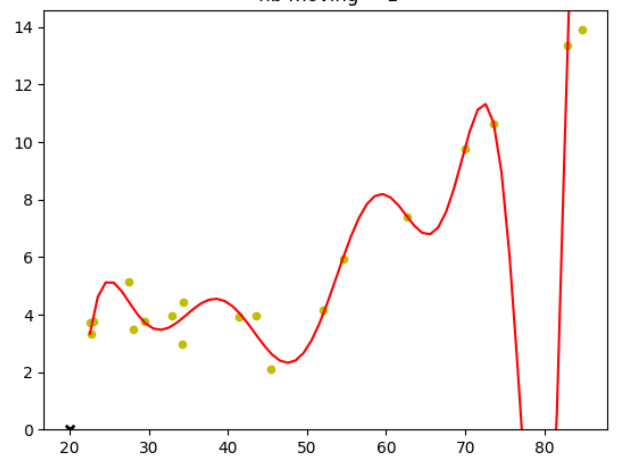
Polynom degree : 10
nb target = 19
nb moving = 1



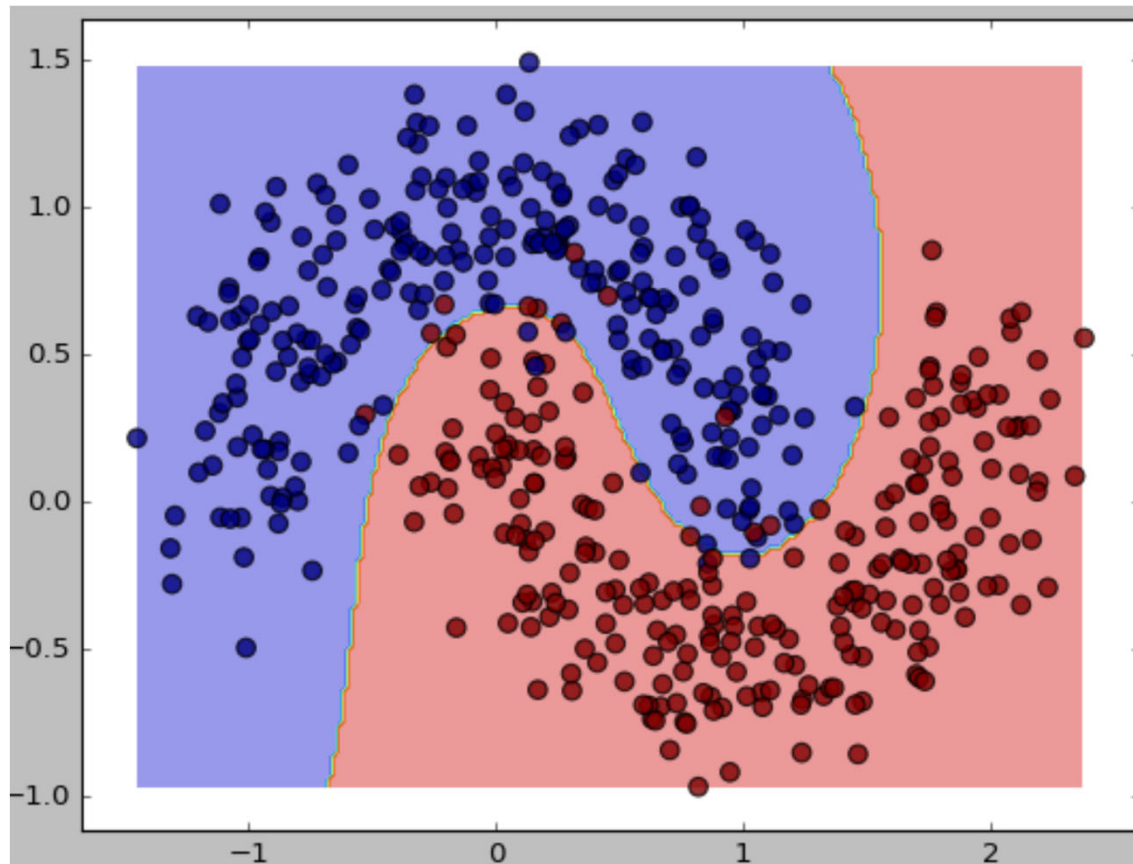
Polynom degree : 15
nb target = 19
nb moving = 1



Polynom degree : 20
nb target = 19
nb moving = 1

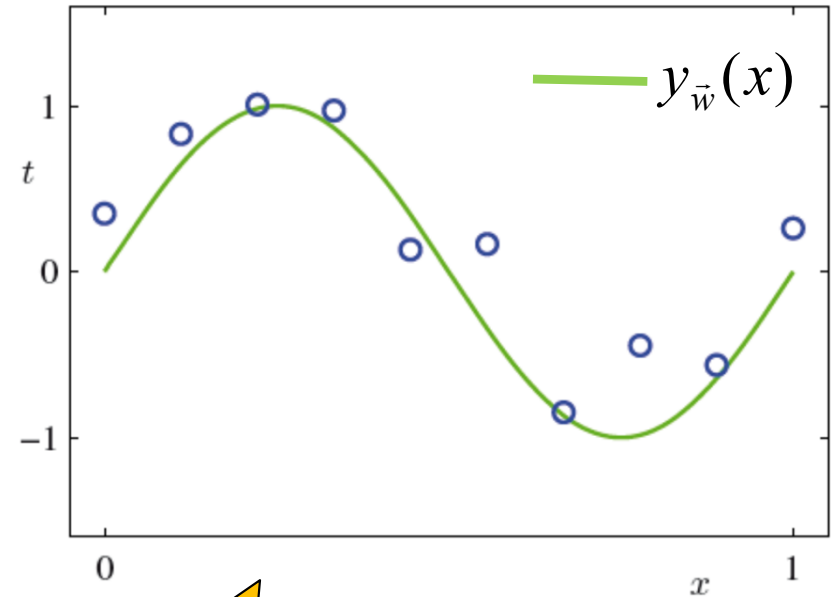


Similar approach for classification



Unknowns

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



Two unknowns

$$\vec{w} \in \mathbb{R}^d$$

Parameters

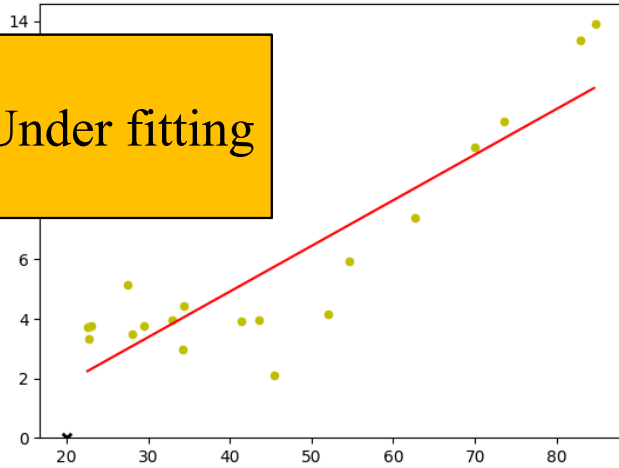
$$d \in \mathbb{N}^{\geq 0}$$

Hyper-parameters

Hyperparameter d

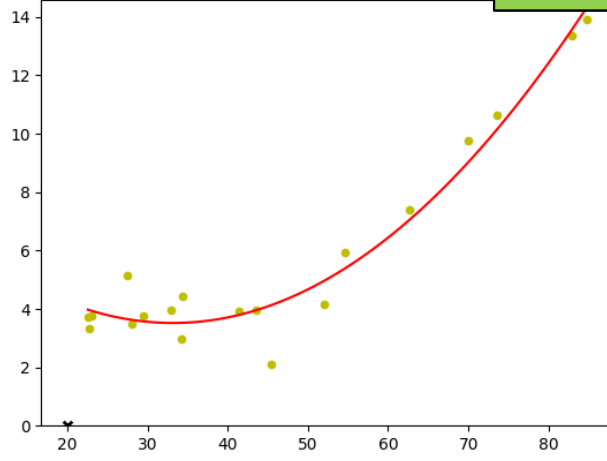
Polynom degree : 1
nb target = 19
nb moving = 1

Under fitting

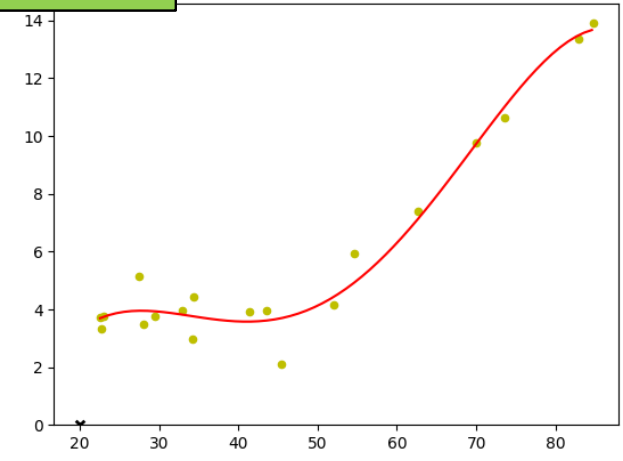


Polynom degree : 2
nb target = 19
nb moving = 1

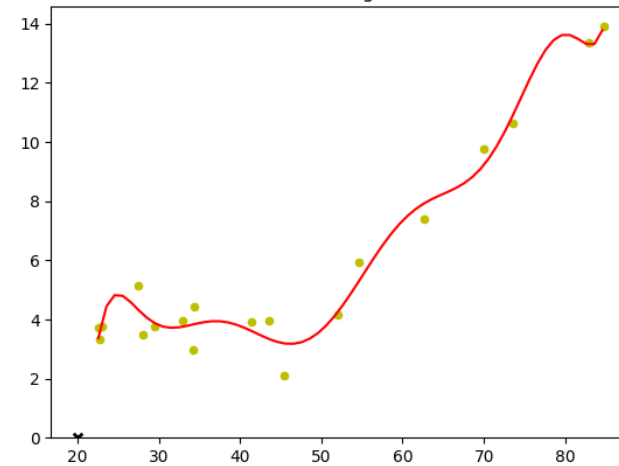
Best fit



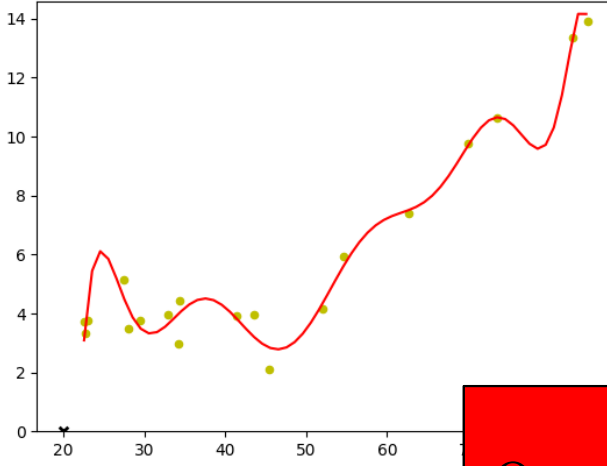
Polynom degree : 4
nb target = 19
nb moving = 1



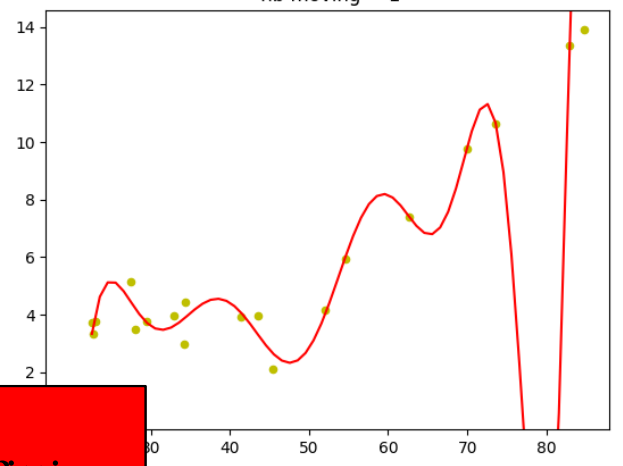
Polynom degree : 10
nb target = 19
nb moving = 1



Polynom degree : 15
nb target = 19
nb moving = 1



Polynom degree : 20
nb target = 19
nb moving = 1

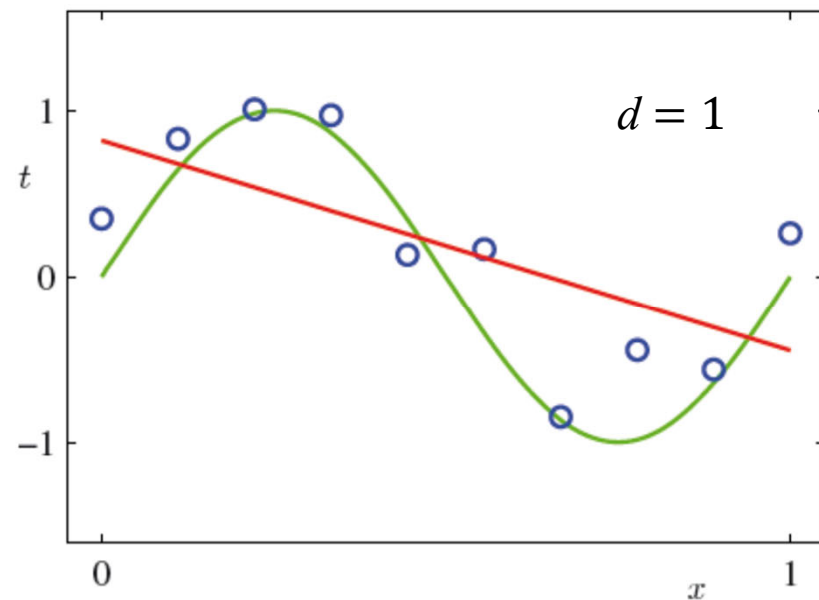
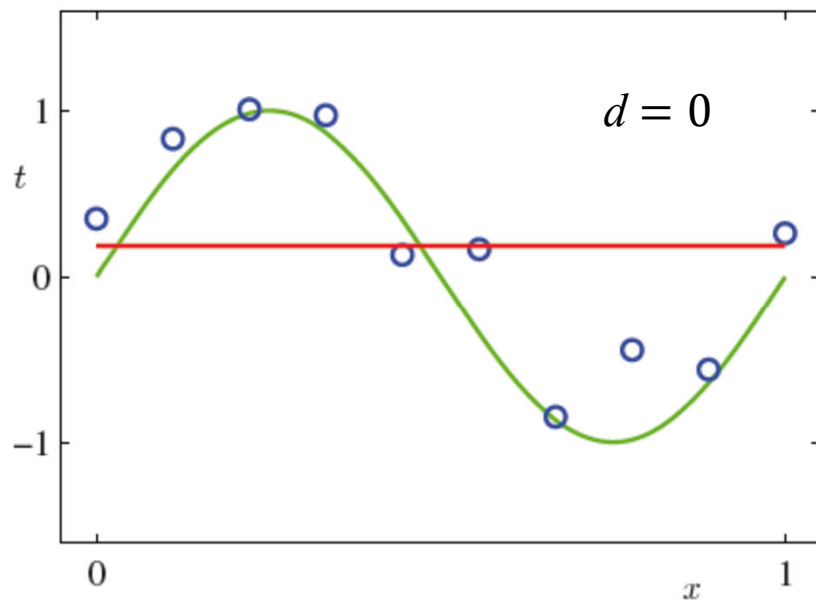


Over fitting

Underfitting

$$\begin{aligned}d = 0 & \Rightarrow y_{\vec{w}}(x) = w_0 \\d = 1 & \Rightarrow y_{\vec{w}}(x) = w_0 + w_1 x\end{aligned}$$

A small d gives a simplistic model that **underfits** the data.

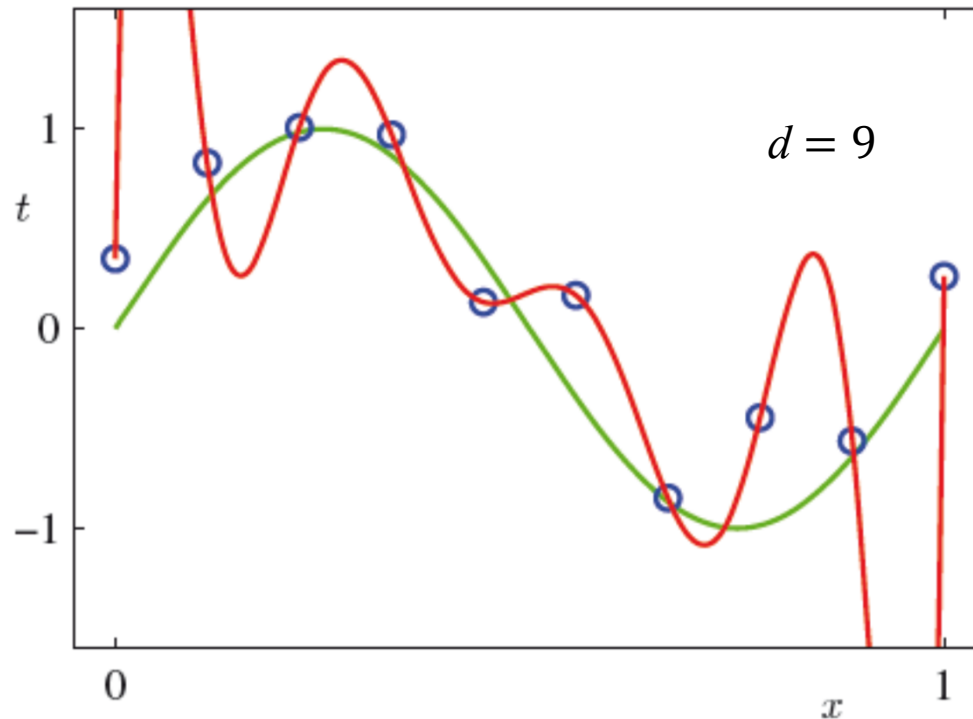


$$\mathcal{L}_D(\vec{w}) \Rightarrow \text{large}$$

$$\mathcal{L}_{D_{test}}(\vec{w}) \Rightarrow \text{large}$$

Overfitting

A large d gives a model that « **learn by heart** » and thus **overfits** training data

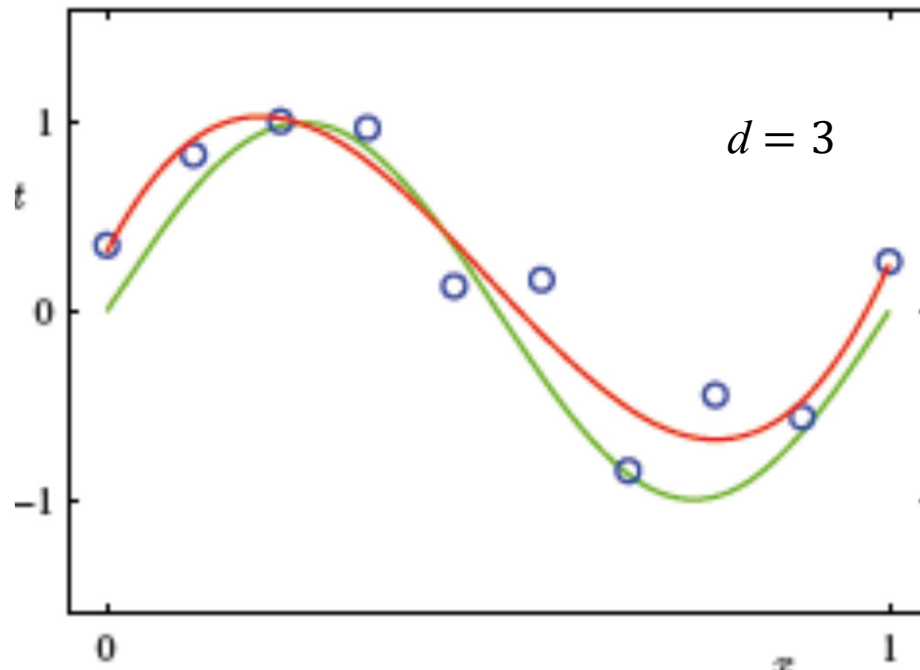


$$\mathcal{L}_D(\vec{w}) \Rightarrow \text{VERY low}$$

$$\mathcal{L}_{D_{test}}(\vec{w}) \Rightarrow \text{large}$$

Over- and underfitting

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



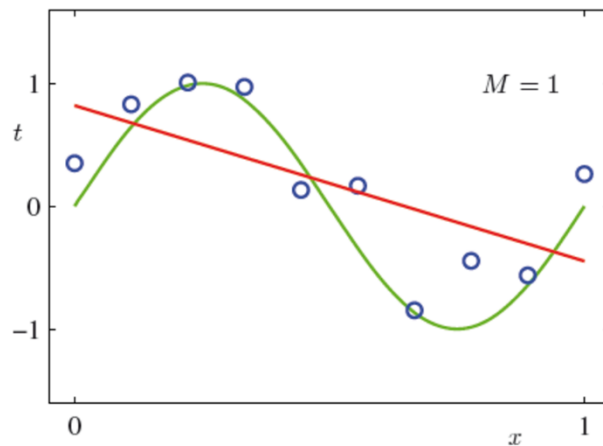
$$\mathcal{L}_D(\vec{w}) \Rightarrow \text{low}$$

$$\mathcal{L}_{D_{test}}(\vec{w}) \Rightarrow \text{low}$$

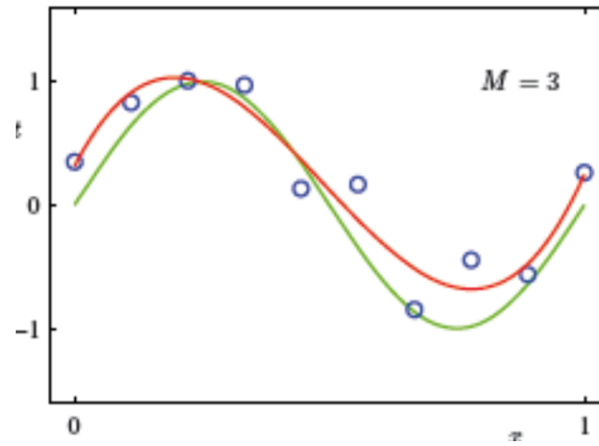
Hyperparameters often control the **capacity** of a model

Capacity: ability of a model to fit the training data

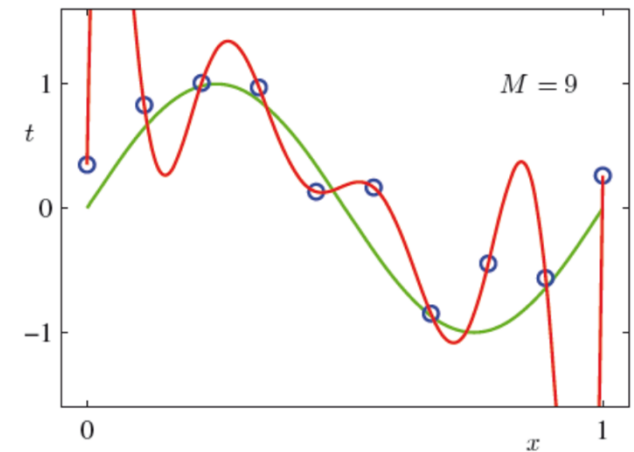
Low capacity



Medium capacity



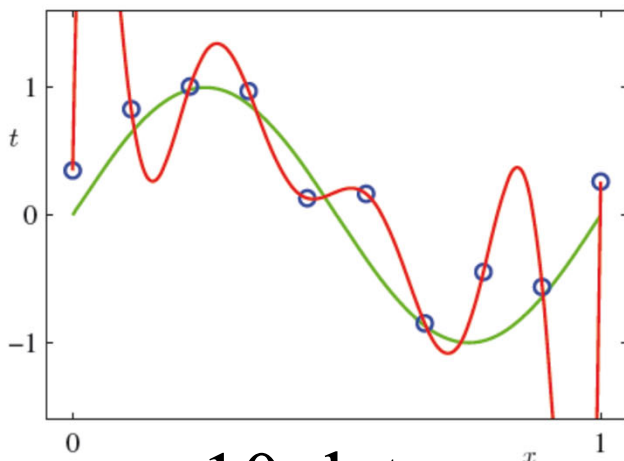
Large capacity



Generalization

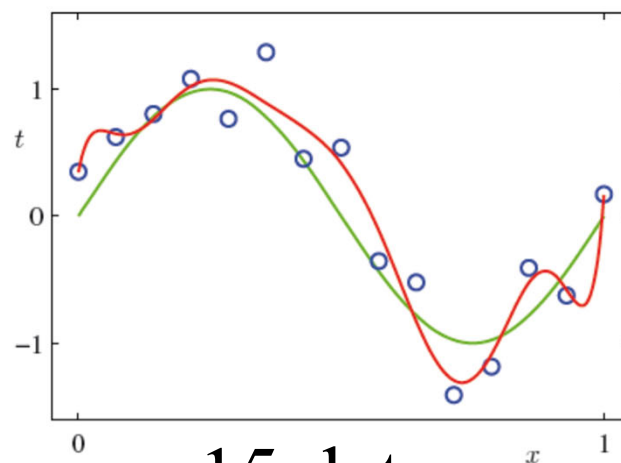
The more data you have, the better a high capacity model will generalize.

$d = 9$



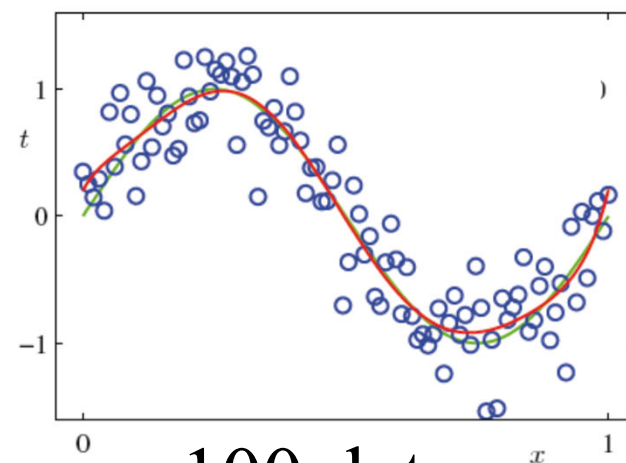
10 data

$d = 9$



15 data

$d = 9$



100 data

How do we prevent our model from under- and overfitting?



Regularization

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

	$d = 0$	$d = 1$	$d = 3$	$d = 9$
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 « d »
2. Reduce capacity by **regularization**

Exemple : penalise the **L2 norm**

Constant that controls regularization

$$E_D(\vec{w}) = \frac{1}{N} \sum_{n=1}^N (t_n - y_{\vec{w}}(\vec{x}))^2 + \lambda \|\vec{w}\|^2$$
$$\|\vec{w}\|^2 = \vec{w}^T \vec{w} = w_0^2 + w_1^2 + \dots + w_d^2$$

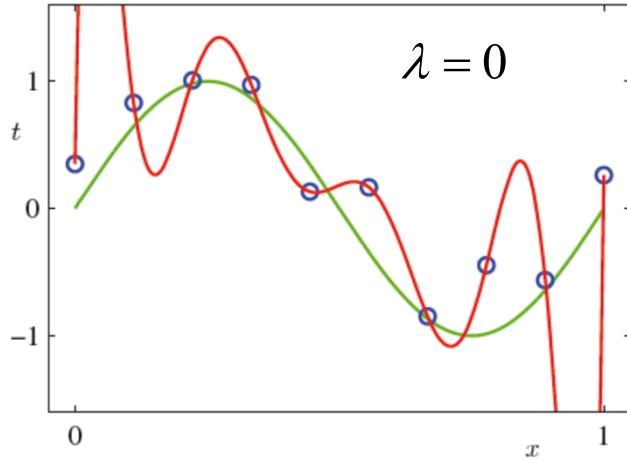


Ridge model

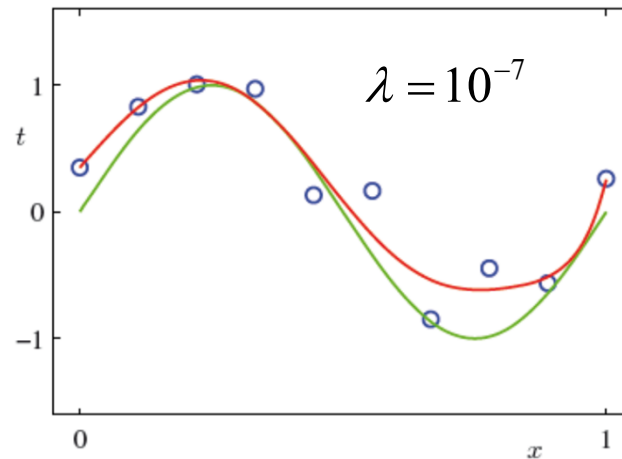
Regularization

Strong regularization = less capacity

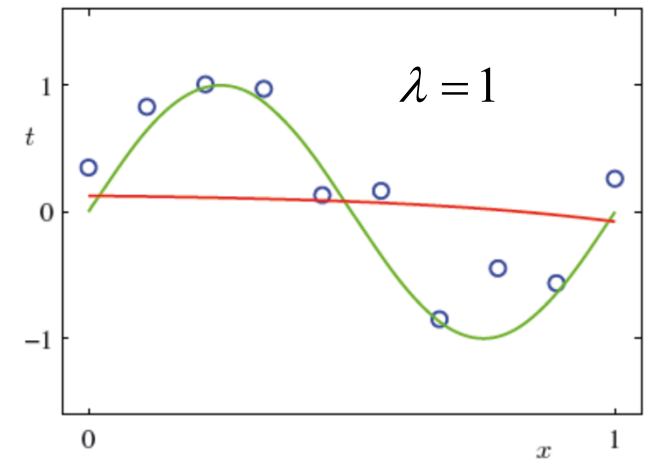
$d = 9$



$d = 9$

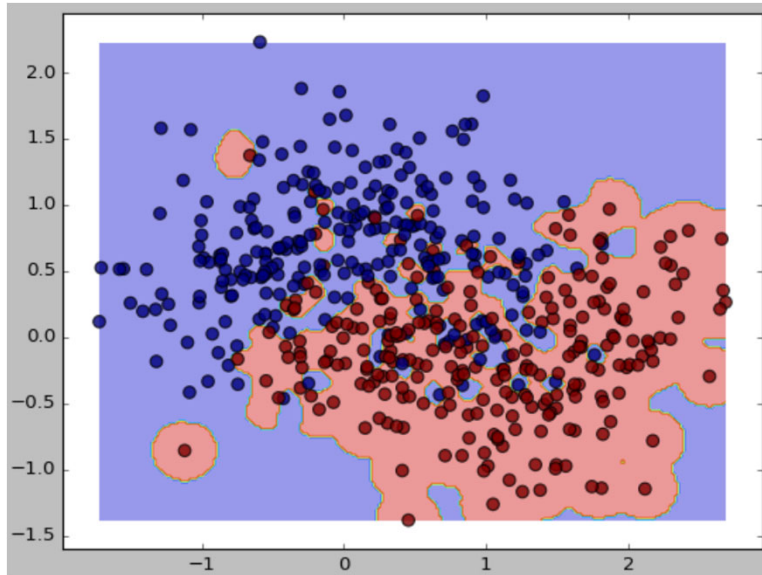


$d = 9$

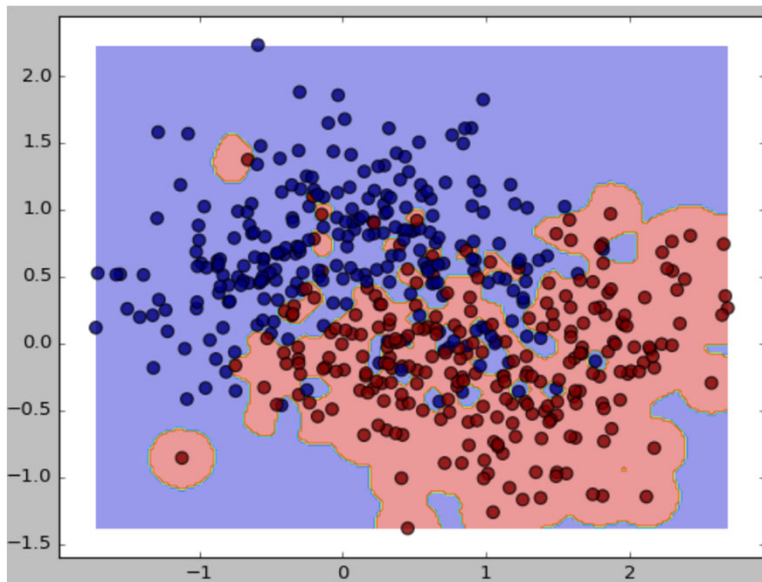


Over- and under-fitting
also influence classification

Overfitting (Classification)

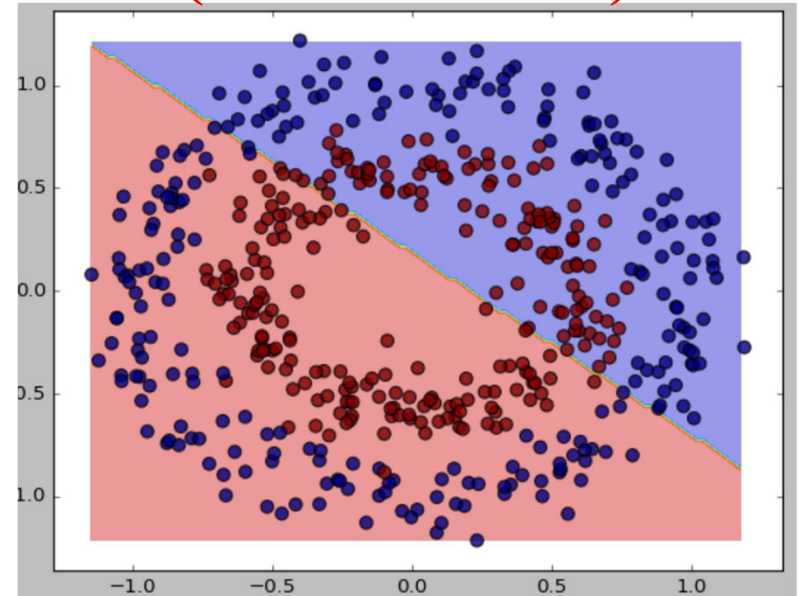


Training accuracy = 99.6%

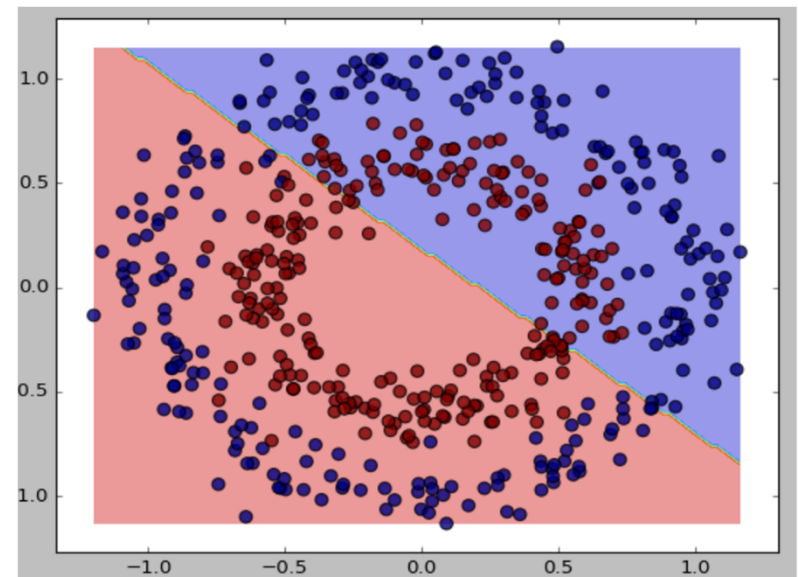


Testing accuracy = 78%

Underfitting (Classification)

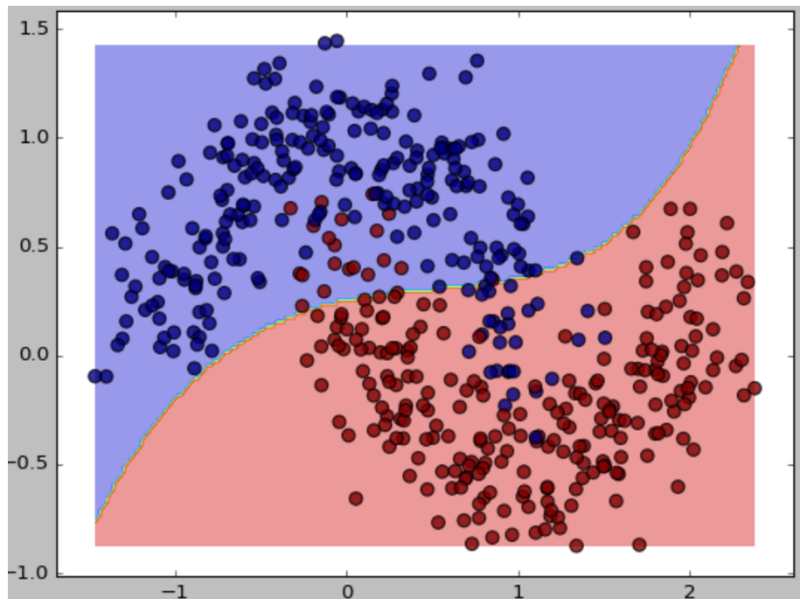


Training accuracy = 52.2%

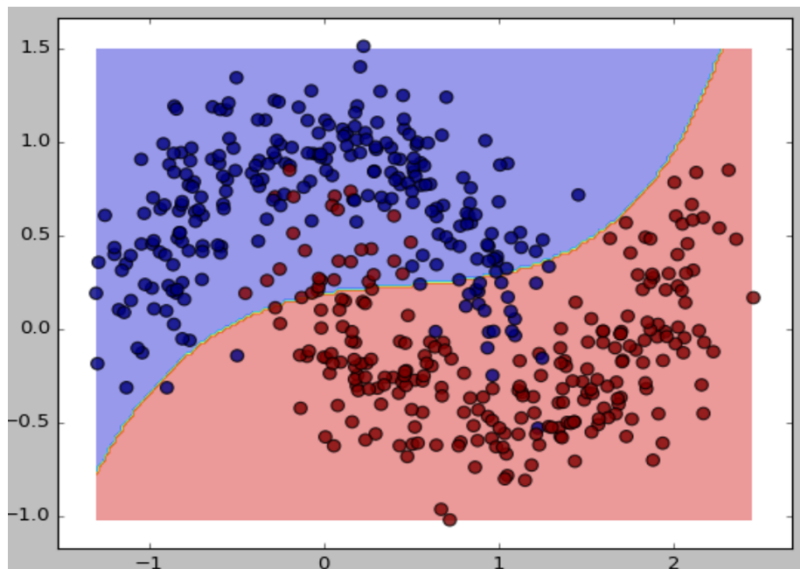


Testing accuracy = 51.2%

Could be better...

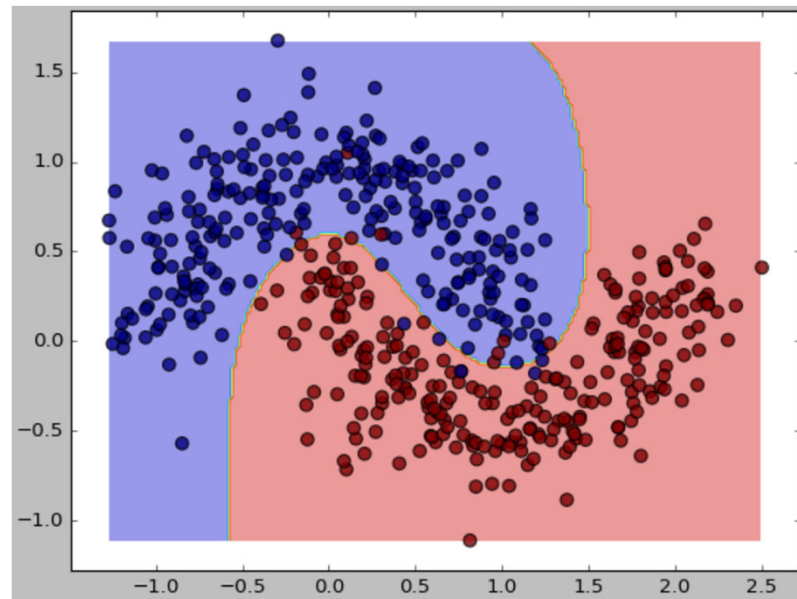


Training accuracy = 82%

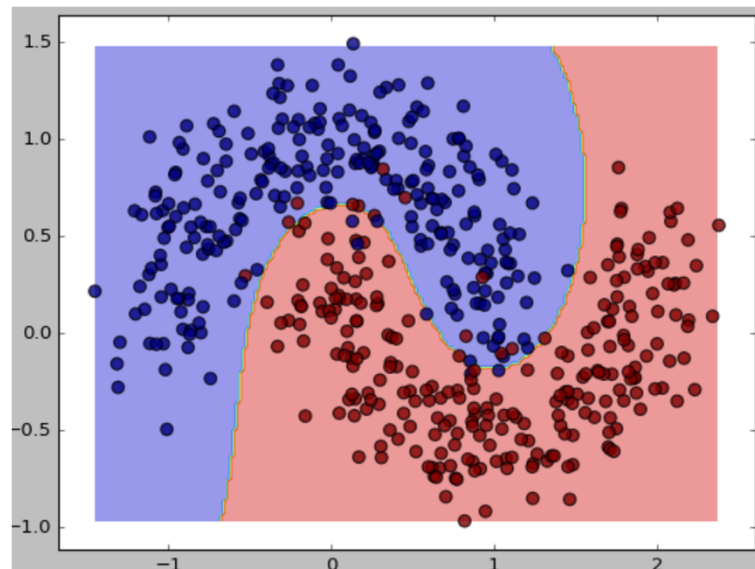


Testing accuracy = 80%

Wonderful !!!



Training accuracy = 97.8%



Testing accuracy = 96.2%

$$E_D(\vec{w}) = \frac{1}{N} \sum_{n=1}^N (y_{\vec{w}}(x_n) - t_n)^2 + \lambda \|\vec{w}\|^2$$
$$\|\vec{w}\|^2 = \vec{w}^T \vec{w} = w_0^2 + w_1^2 + \dots + w_d^2$$

Model selection

How to find the right hyper-parameters?

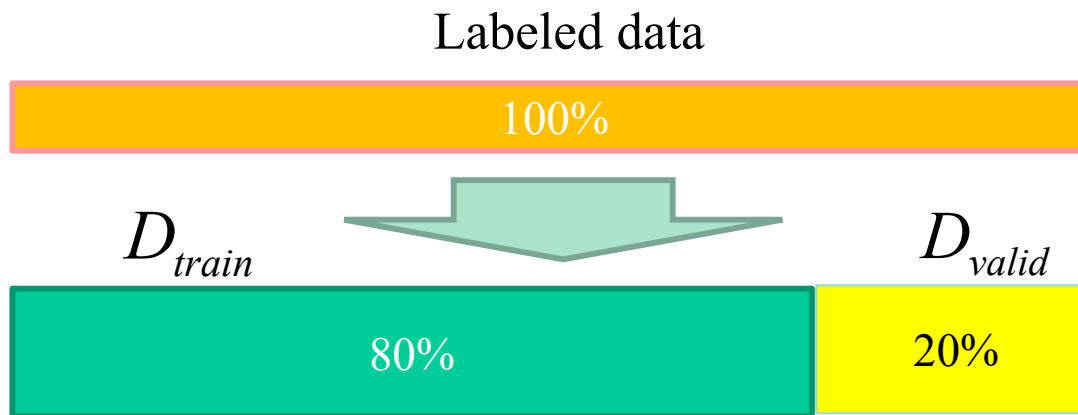
d and λ

How to find the right d and the right λ ?

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

Cross-validation

1- Randomly divide data in 2 groups



2- FOR M from M_{\min} to M_{\max}
FPR λ from λ_{\min} to λ_{\max}

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

3- Keep (M, λ) with the lowest **validation error**

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

Mean validation error

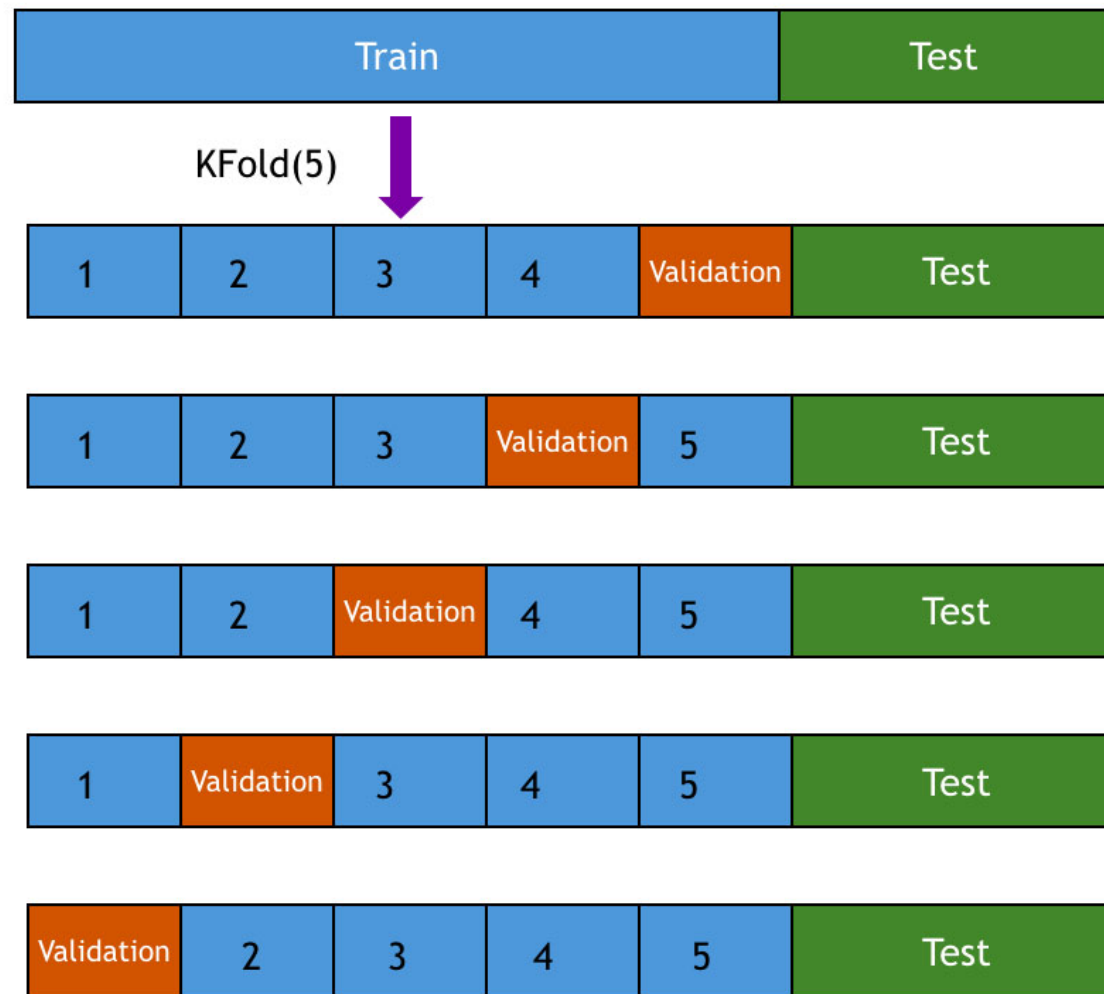
↓ ↙ STD

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

BEST!

$d=5,$
 $\lambda = 1$

k-fold cross-validation



In short

- ✓ The goal is to **train** a **model** on a **training dataset** with good **generalization** capabilities
- ✓ Training = minimization of a **loss function**
- ✓ 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)

Thank you!