

IFT6390  
Fondements de l'apprentissage machine

Learning terminology

Professor: Ioannis Mitliagkas

# Learning from examples

# Learning from examples



# Learning from examples



“horse”

# Learning from examples



“horse”

# Learning from examples



“horse”



“horse”

# Learning from examples



“horse”



“horse”



# Learning from examples



“horse”



“horse”



“horse”

# Learning from examples



“horse”



“horse”



“horse”

Principle much more general than to write by hand, starting from scratch, an algorithm to recognize a horse..

# Categories of problems (tasks) in machine learning

## Supervised learning

- o Classification
- o Regression

## Non-supervised learning

- o Density estimation (learn the distribution of examples)
- o Clustering
- o Dimensionality reduction
- o Feature extraction

## Reinforcement learning

(not scheduled for this class)

# Example of classification problem

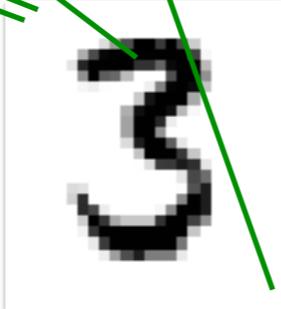
entry  $x_i$

(vector representation)

$$x \in \mathbb{R}^d$$

$$x = (0, 0, \dots, 54, 120, \dots, 0, 0)$$

Test point:  
(new  $x$ )



# Example of classification problem

*training set*

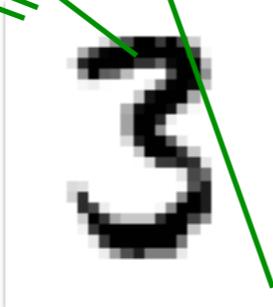
entry  $x_i$

(vector representation)

$$x \in \mathbb{R}^d$$

$$x = (0, 0, \dots, 54, 120, \dots, 0, 0)$$

Test point:  
(new  $x$ )



# Example of classification problem

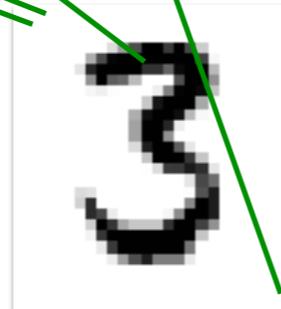
entry  $x_i$

(vector representation)

$$x \in \mathbb{R}^d$$

$$x = (0, 0, \dots, 54, 120, \dots, 0, 0)$$

Test point:  
(new  $x$ )



2	2
2	2
2	2
2	2
2	2
3	3
3	3
3	3
3	3
3	3

training set

# Example of classification problem

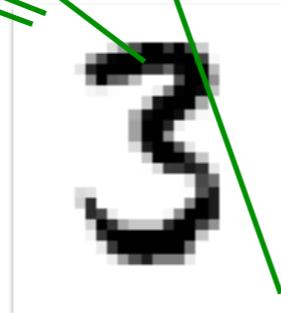
entry  $x_i$

(vector representation)

$$x \in \mathbb{R}^d$$

$$x = (0, 0, \dots, 54, 120, \dots, 0, 0)$$

Test point:  
(new  $x$ )



2	2
2	2
2	2
2	2
2	2
3	3
3	3
3	3
3	3
3	3

training set

label  $y_i$

# Example of classification problem

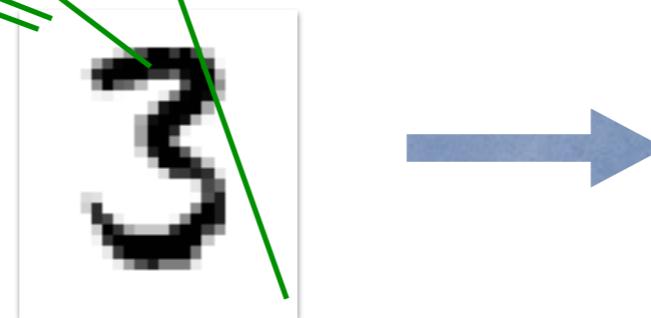
entry  $x_i$

(vector representation)

$$x \in \mathbb{R}^d$$

$$x = (0, 0, \dots, 54, 120, \dots, 0, 0)$$

Test point:  
(new  $x$ )



2	2
2	2
2	2
2	2
2	2
3	3
3	3
3	3
3	3
3	3

training set

label  $y_i$

# Example of classification problem

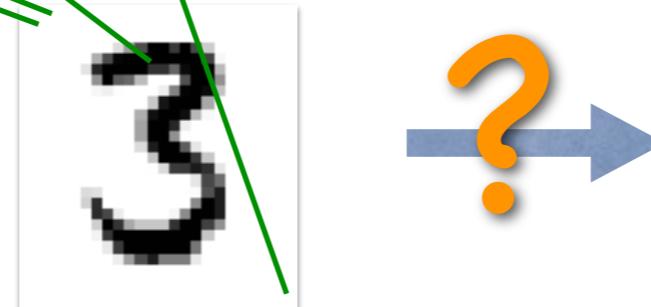
entry  $x_i$

(vector representation)

$$x \in \mathbb{R}^d$$

$$x = (0, 0, \dots, 54, 120, \dots, 0, 0)$$

Test point:  
(new  $x$ )



2	2
2	2
2	2
2	2
2	2
3	3
3	3
3	3
3	3
3	3

training set

label  $y_i$

# Example of classification problem

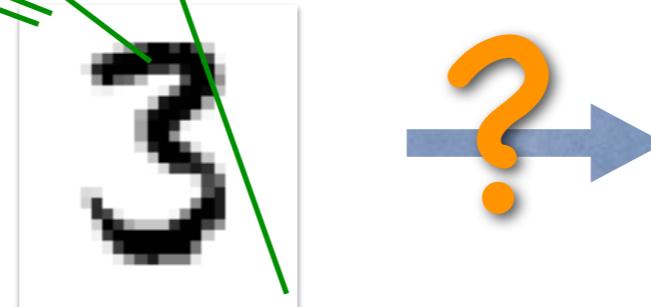
entry  $x_i$

(vector representation)

$$x \in \mathbb{R}^d$$

$$x = (0, 0, \dots, 54, 120, \dots, 0, 0)$$

Test point:  
(new  $x$ )



2 or 3?

2	2
2	2
2	2
2	2
2	2
3	3
3	3
3	3
3	3
3	3

training set

label  $y_i$

# Example of classification problem

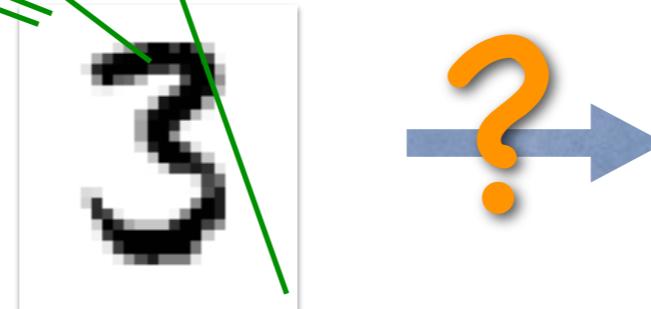
entry  $x_i$

(vector representation)

$$x \in \mathbb{R}^d$$

$$x = (0, 0, \dots, 54, 120, \dots, 0, 0)$$

Test point:  
(new  $x$ )



2	2
2	2
2	2
2	2
2	2
3	3
3	3
3	3
3	3
3	3

training set

label  $y_i$

Learning is not  
memorizing

# Example of classification problem

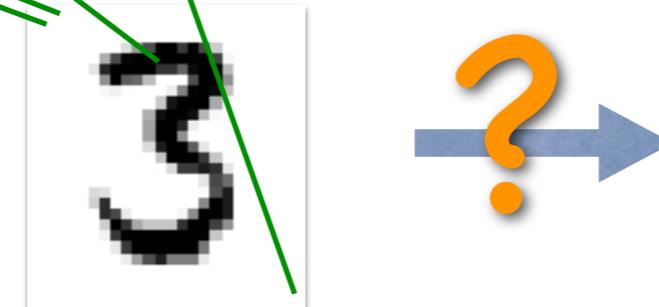
entry  $x_i$

(vector representation)

$$x \in \mathbb{R}^d$$

$$x = (0, 0, \dots, 54, 120, \dots, 0, 0)$$

Test point:  
(new  $x$ )



2 or 3?

2	2
2	2
2	2
2	2
2	2
3	3
3	3
3	3
3	3
3	3

training set

label  $y_i$

*Learning is not  
memorizing*

*It is the ability to  
generalize to unseen  
cases!*

# Data representation

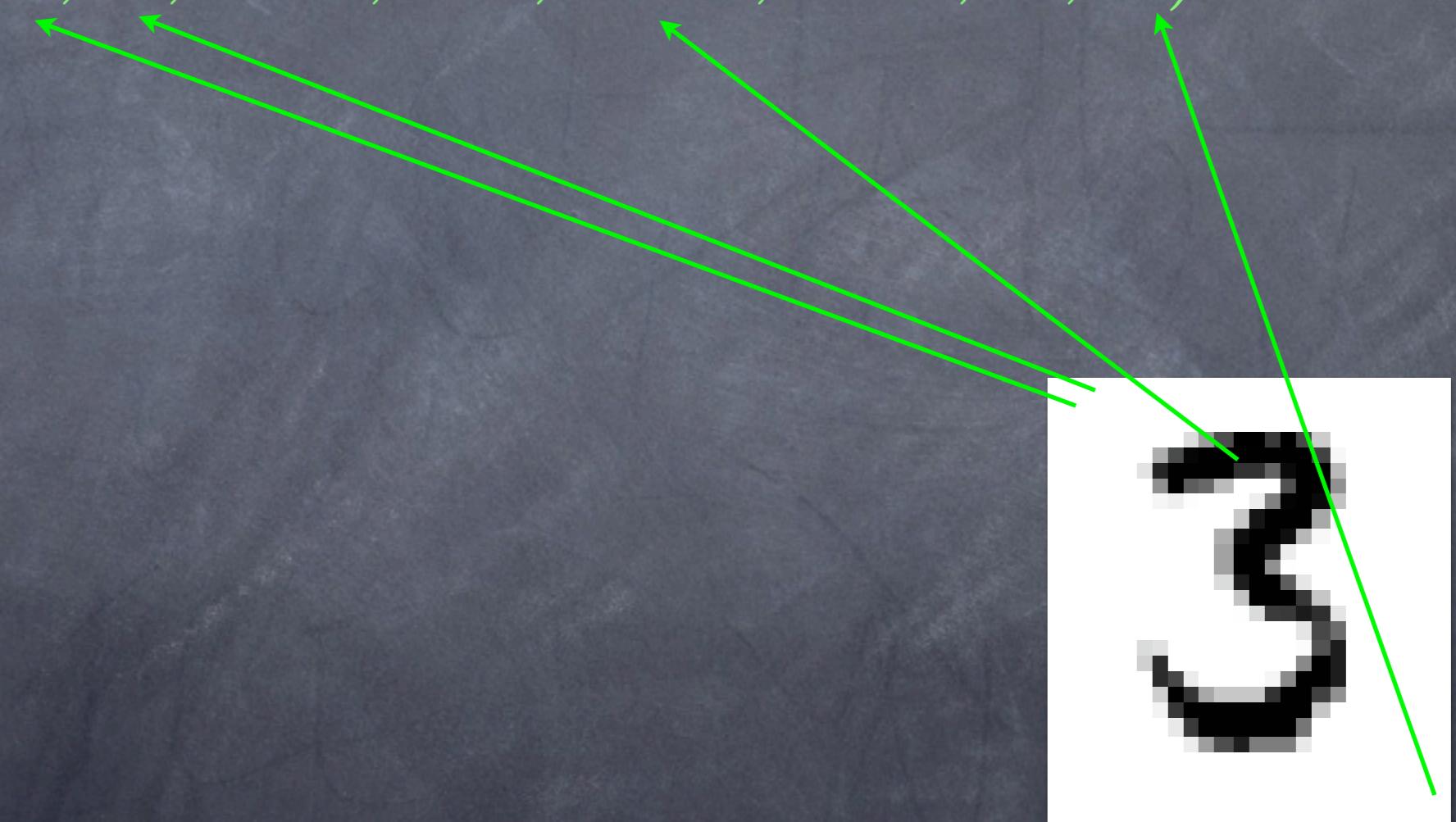
- Most learning algorithms require a **vector representation** of the examples
- Numerical vectors of fixed dimension **d**

$$x \in \mathbb{R}^d$$

ex:  $x = (0, 0, \dots, 54, 120, \dots, 0, 0)$

Example: raw vector representation of  
an image (bitmap) in grayscale

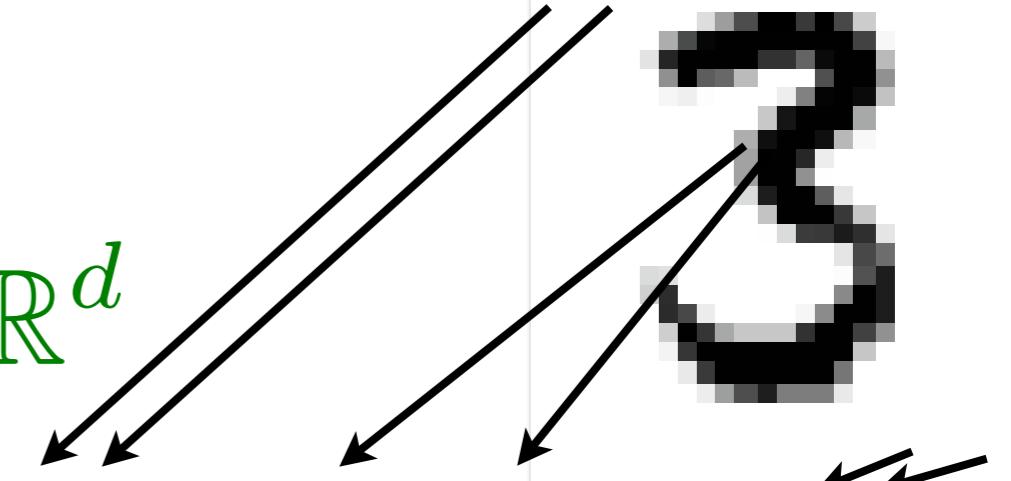
$$x = (0, 0, \dots, 54, 120, \dots, 0, 0)$$



Transform an example into a vector representation  $\mathbf{x} \in \mathbb{R}^d$

Raw representation:

$$\mathbf{x} = (0, 0, \dots, 54, 120, \dots, 0, 0)$$



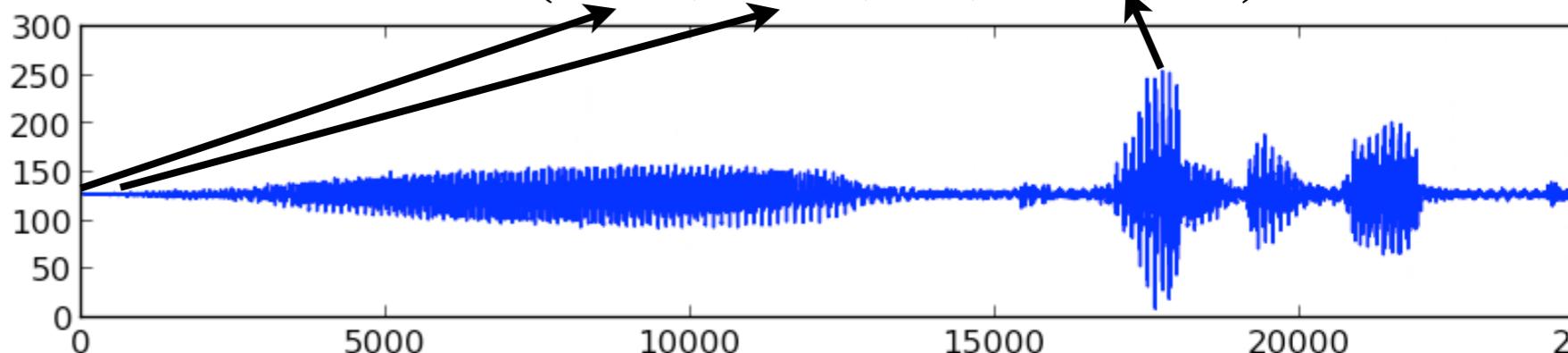
Transform an example into a vector representation  $\mathbf{x} \in \mathbb{R}^d$

Raw representation:

$$\mathbf{x} \in \mathbb{R}^d$$

$$\mathbf{x} = (0, 0, \dots, 54, 120, \dots, 0, 0)$$

$$\mathbf{x} = (125, 125, \dots, 250, \dots)$$



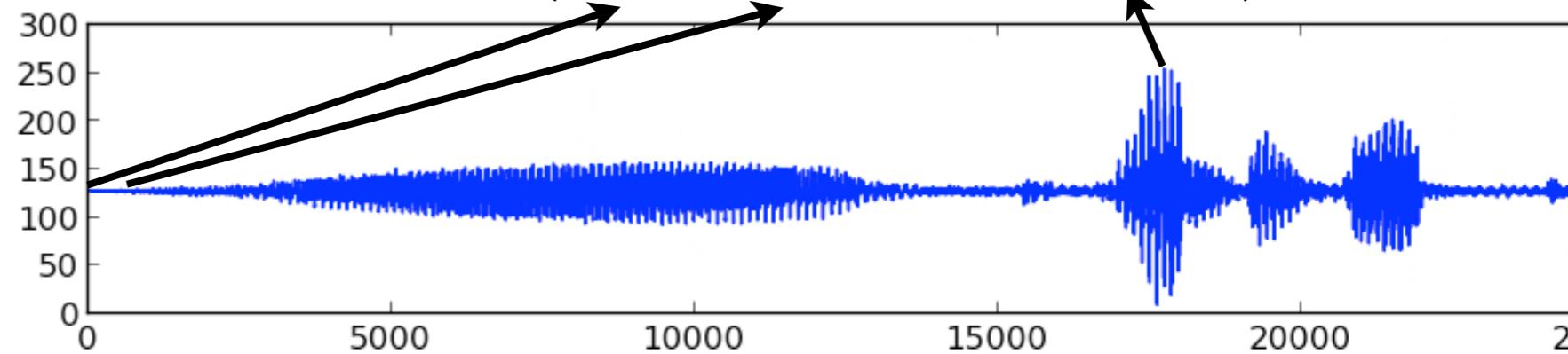
Transform an example into a vector representation  $\mathbf{x} \in \mathbb{R}^d$

Raw representation:

$$\mathbf{x} \in \mathbb{R}^d$$

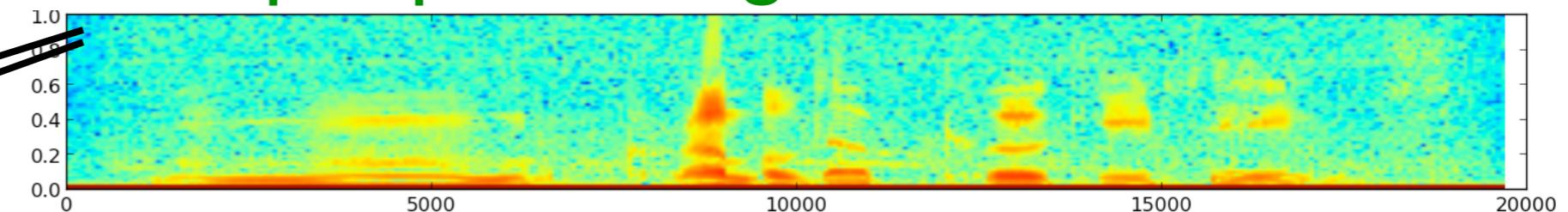
$$\mathbf{x} = (0, 0, \dots, 54, 120, \dots, 0, 0)$$

$$\mathbf{x} = (125, 125, \dots, 250, \dots)$$



Or feature extraction from pre-processing:

$$\mathbf{x} = (, \leftarrow, \leftarrow, , \dots)$$



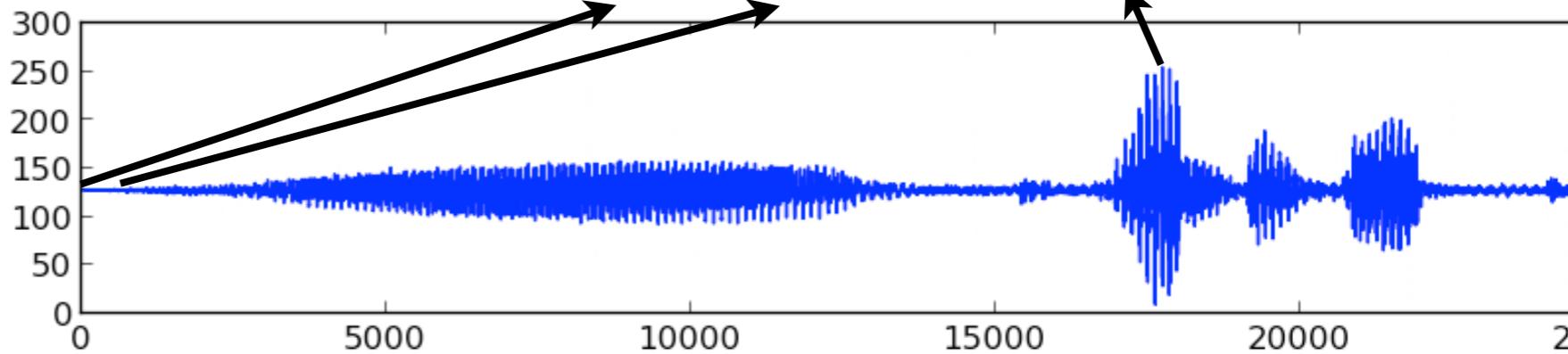
Transform an example into a vector representation  $\mathbf{x} \in \mathbb{R}^d$

Raw representation:

$$\mathbf{x} \in \mathbb{R}^d$$

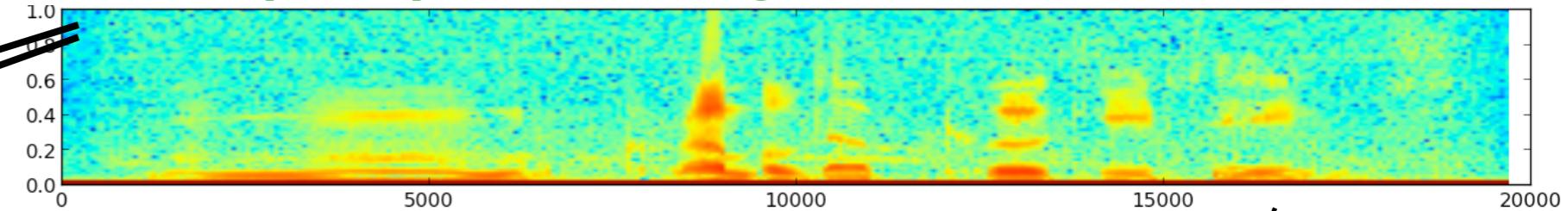
$$\mathbf{x} = (0, 0, \dots, 54, 120, \dots, 0, 0)$$

$$\mathbf{x} = (125, 125, \dots, 250, \dots)$$



Or feature extraction from pre-processing:

$$\mathbf{x} = (, , , , \dots)$$



Bag of words for «The cat jumped»:  $\mathbf{x} = (\dots 0 \dots, 0, 1, \dots 0 \dots, 1, 0, 0, \dots, 0, 0, 1, 0, \dots 0 \dots)$

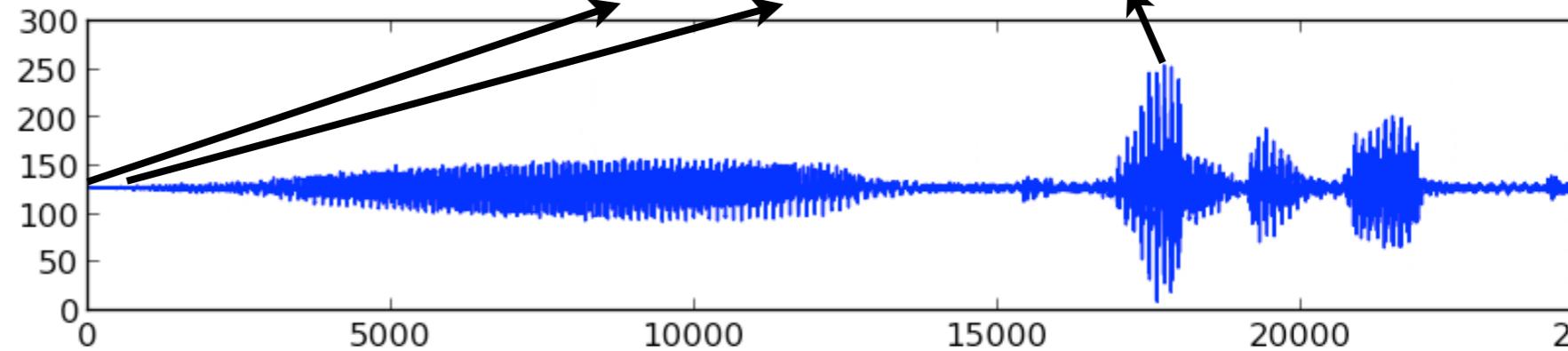
Transform an example into a vector representation  $\mathbf{x} \in \mathbb{R}^d$

Raw representation:

$$\mathbf{x} \in \mathbb{R}^d$$

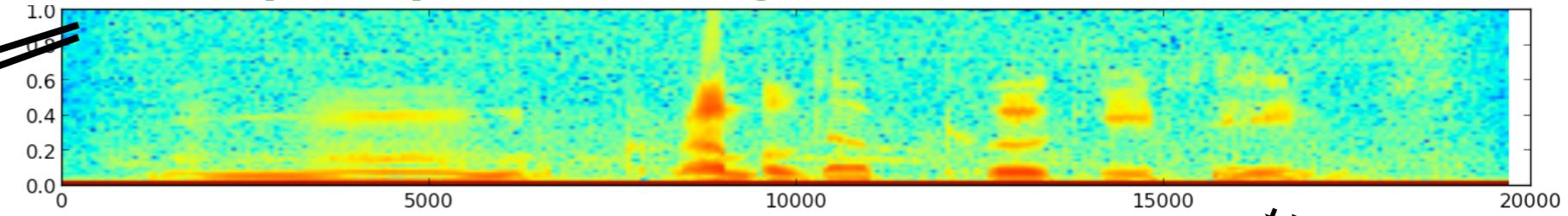
$$\mathbf{x} = (0, 0, \dots, 54, 120, \dots, 0, 0)$$

$$\mathbf{x} = (125, 125, \dots, 250, \dots)$$



Or feature extraction from pre-processing:

$$\mathbf{x} = (, , , , \dots)$$



we the jumped jumping run elephant dog cat horse

Bag of words for «The cat jumped»:  $\mathbf{x} = (\dots 0 \dots, 0, 1, \dots 0 \dots, 1, 0, 0, \dots, 0, 0, 1, 0, \dots 0 \dots)$

OR vector of handmade features:

ex: Histograms of Oriented Gradients

$$\mathbf{x} = (\text{feature } 1, \dots, \text{feature } d)$$

Another possibility:  
**Characteristic traits**  
describe the inputs  
(features, descriptors)

$\mathbf{x} = (\text{feature 1}, \text{feature 2}, \dots, \text{feature d})$

Another possibility:  
**Characteristic traits**  
describe the inputs  
(features, descriptors)

$$\mathbf{x} = (\text{feature 1}, \text{feature 2}, \dots, \text{feature d})$$

- ⦿ Often designed manually to facilitate the task of learning (feature engineering).
- ⦿ eg. to represent a natural image (vision): SIFT descriptors.
- ⦿ eg. to represent a document: frequency of occurrences of certain words (bag of words).

Another possibility:  
**Characteristic traits**  
describe the inputs  
(features, descriptors)

$$\mathbf{x} = (\text{feature 1}, \text{feature 2}, \dots, \text{feature d})$$

- ⦿ Often designed manually to facilitate the task of learning (feature engineering).
- ⦿ eg. to represent a natural image (vision): SIFT descriptors.
- ⦿ eg. to represent a document: frequency of occurrences of certain words (bag of words).  
eg. feature k = how many times the word "plane" appears in the document.





# Training set

input:



targets:

"horse"



"cat"

etc...



"horse"

# Training set

number  
of  
examples  
 $n$

input:



targets:

"horse"



"cat"

etc...



"horse"

# Training set

number  
of  
examples  
 $n$

input:



targets:

"horse"



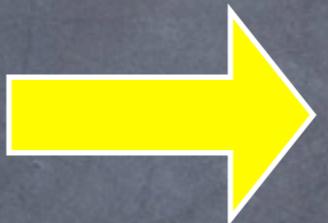
"cat"

etc...



"horse"

preprocessing,  
feature  
extraction



## Training set

number  
of  
examples  
 $n$

input:



targets:

"horse"



"cat"

etc...



"horse"

preprocessing,  
feature  
extraction

$X_1$

$X_n$

inputs:  $X$

(vector of features)

(3.5, -2, ..., 127, 0, ...)

+1

(-9.2, 32, ..., 24, 1, ...)

-1

etc...

$\underbrace{X_{n,2}}$

(6.8, 54, ..., 17, -3, ...)

+1

$Y_1$

$Y_n$

## Training set

number  
of  
examples  
 $n$

input:



targets:

"horse"

input:



targets:

"cat"

etc...



targets:

"horse"

preprocessing,  
feature  
extraction

$X_1$

$X_n$

inputs:  $X$

(vector of features)

(3.5, -2, ..., 127, 0, ...)

+1

$Y_1$

(-9.2, 32, ..., 24, 1, ...)

-1

etc...

$X_{n,2}$

(6.8, 54, ..., 17, -3, ...)

+1

$Y_n$

Dimensionality  
of input

$d$

## Training set

number  
of  
examples  
 $n$

input:



targets:

"horse"

input:



targets:

"cat"

etc...



targets:

"horse"

Test point:



Dimensionality  
of input

$d$

inputs:  $X$

(vector of features)

$(3.5, -2, \dots, 127, 0, \dots)$

+1

$X_1$

$(-9.2, 32, \dots, 24, 1, \dots)$

-1

etc...

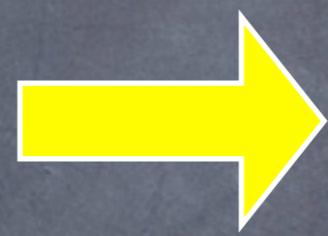
$X_n$

$X_{n,2}$

$(6.8, 54, \dots, 17, -3, \dots)$

+1

$Y_n$



preprocessing,  
feature  
extraction

## Training set

number  
of  
examples  
 $n$

input:



targets:

"horse"

input:



targets:

"cat"

etc...



targets:

"horse"

Test point:



Dimensionality  
of input

$d$

inputs:  $X$

(vector of features)

$X_1$   
(3.5, -2, ..., 127, 0, ...)

+1

$X_n$   
(-9.2, 32, ..., 24, 1, ...)

-1

etc...

$X_{n,2}$

$X_n$   
(6.8, 54, ..., 17, -3, ...)

+1

$Y_1$

preprocessing,  
feature  
extraction

## Training set

number  
of  
examples  
 $n$

input:



targets:

"horse"



"cat"

etc...



"horse"

Test point:



→ ?

Dimensionality  
of input

$d$

inputs:  $X$

(vector of features)

$X_1$   
(3.5, -2, ..., 127, 0, ...)

+1

$X_n$   
(-9.2, 32, ..., 24, 1, ...)

-1

etc...

$X_{n,2}$

$X_n$   
(6.8, 54, ..., 17, -3, ...)

+1

$Y_1$

$Y_n$

# Training set

number  
of  
examples  
 $n$

input:



targets:

"horse"



"cat"

etc...



"horse"



→ ?

Test point:

Dimensionality  
of input

$d$

inputs:  $X$

(vector of features)

$X_1$   
(3.5, -2, ..., 127, 0, ...)

+1

$X_n$   
(-9.2, 32, ..., 24, 1, ...)

-1

etc...

$X_{n,2}$

$X_n$   
(6.8, 54, ..., 17, -3, ...)

+1

$\mathcal{X} = (5.7, -27, \dots, 64, 0, \dots) \rightarrow$

$Y_1$

$Y_n$

preprocessing,  
feature  
extraction

## Training set

number  
of  
examples  
 $n$

input:



targets:

"horse"



"cat"

etc...



"horse"



→ ?

Test point:

Dimensionality  
of input

$d$

inputs:  $X$

(vector of features)

$X_1$   
(3.5, -2, ..., 127, 0, ...)

+1

$X_n$   
(-9.2, 32, ..., 24, 1, ...)

-1

etc...

$X_{n,2}$

$X_n$   
(6.8, 54, ..., 17, -3, ...)

+1

$Y_n$

preprocessing,  
feature  
extraction

$X_n$

$\mathcal{X} =$

(5.7, -27, ..., 64, 0, ...)

→ ?

# Importance of the problem dimension

⇒ Determines which learning algorithms are applicable or not (because of their algorithmic complexity).

- Number of examples: **n**  
(often in the millions of more)
- Dimension of input: **d**  
how many features we are using  
(often in the order of 100 to 1000, sometimes in the 10s of thousands)
- Dimensions of the target we want to predict eg., number of classes **m** (often low, sometimes thousands)

→ A dataset will often be organized  
in a matrix  $n \times (d+1)$  or  $n \times (d+m)$

# Terminology of supervised learning

One input  
generally  
represented by  
a vector of  
dimension d.

$$x \in \mathbb{R}^d$$

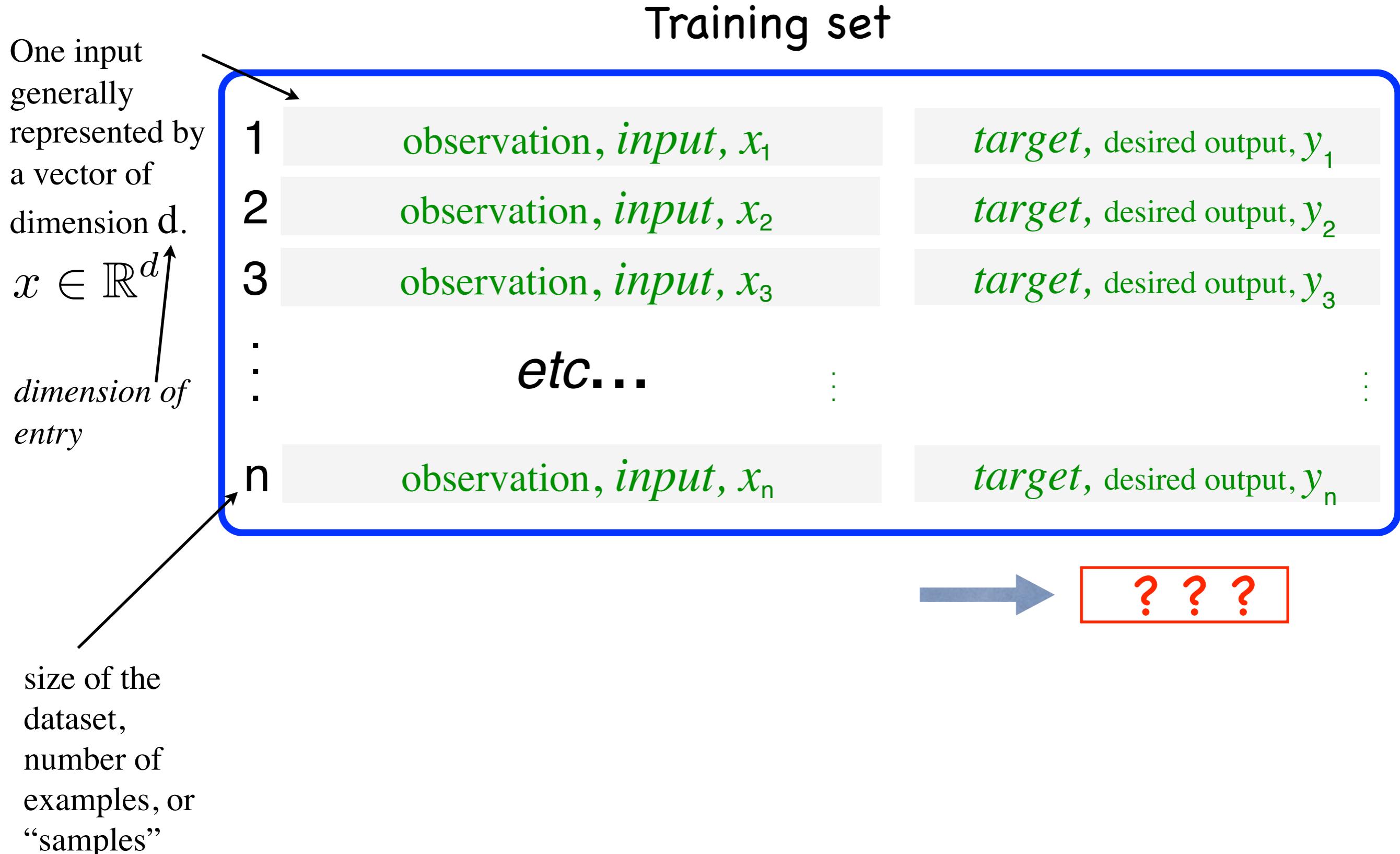
*dimension of  
entry*

1	observation, <i>input</i> , $x_1$	<i>target</i> , desired output, $y_1$
2	observation, <i>input</i> , $x_2$	<i>target</i> , desired output, $y_2$
3	observation, <i>input</i> , $x_3$	<i>target</i> , desired output, $y_3$
.	etc...	:
n	observation, <i>input</i> , $x_n$	<i>target</i> , desired output, $y_n$

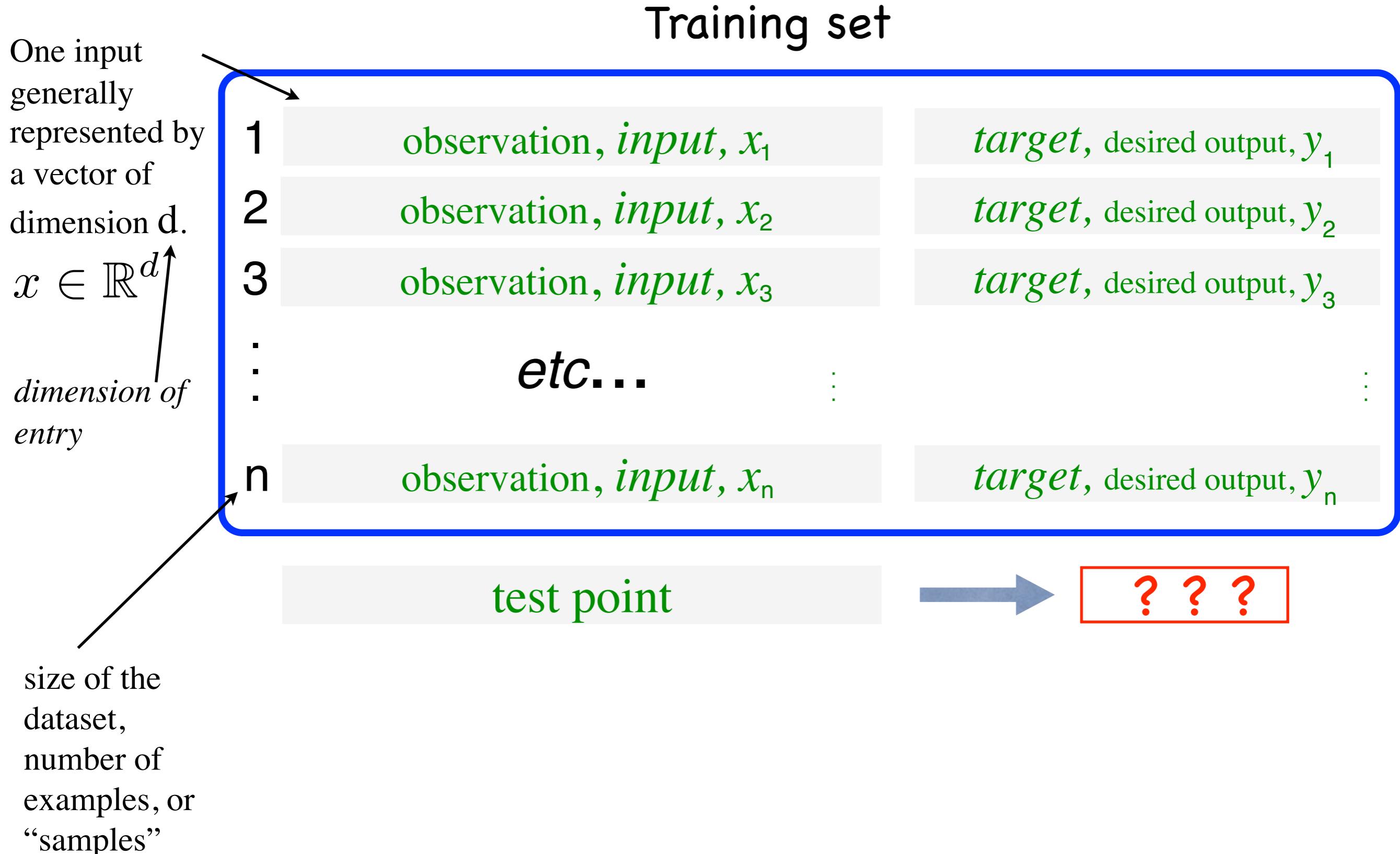


size of the  
dataset,  
number of  
examples, or  
“samples”

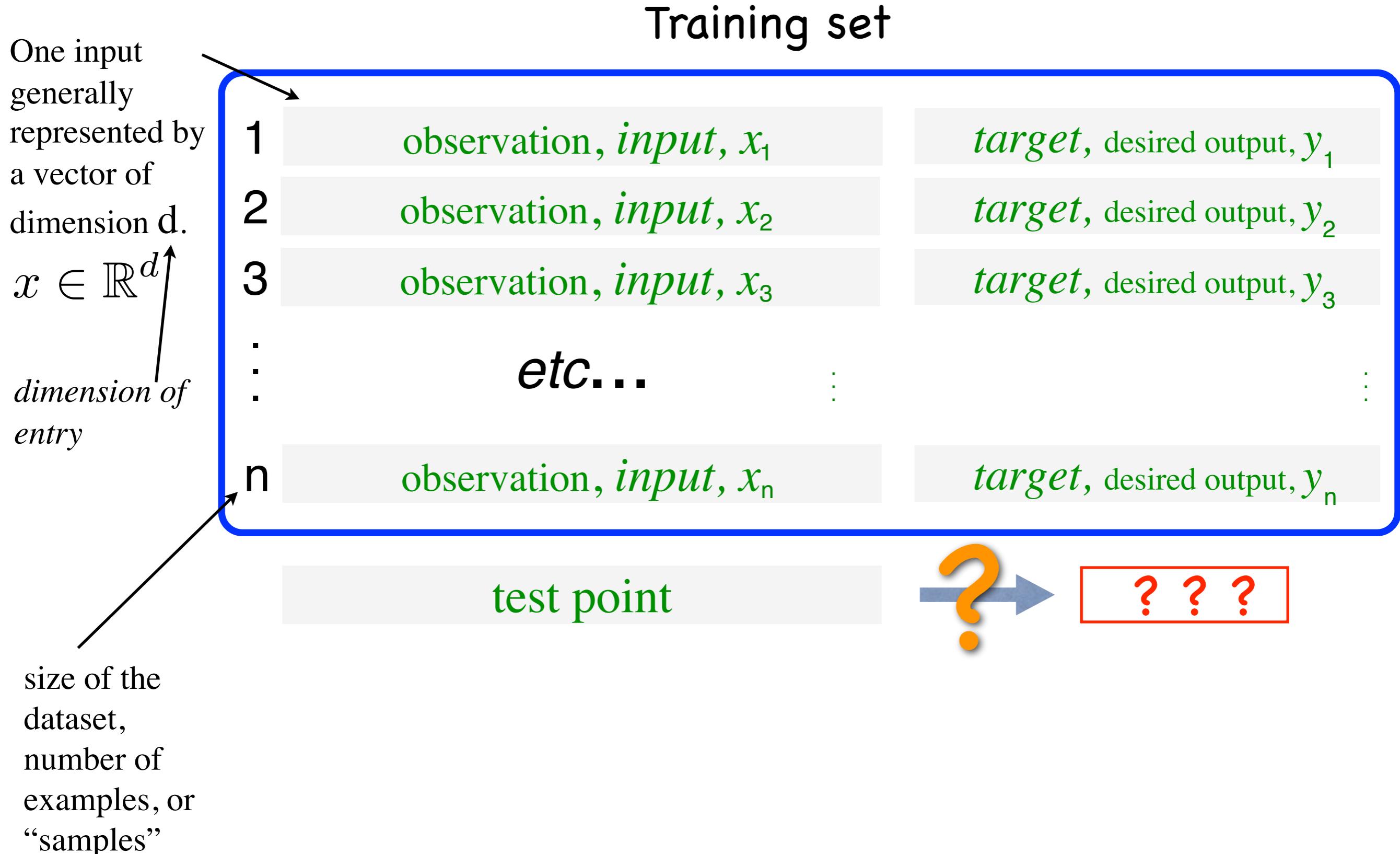
# Terminology of supervised learning



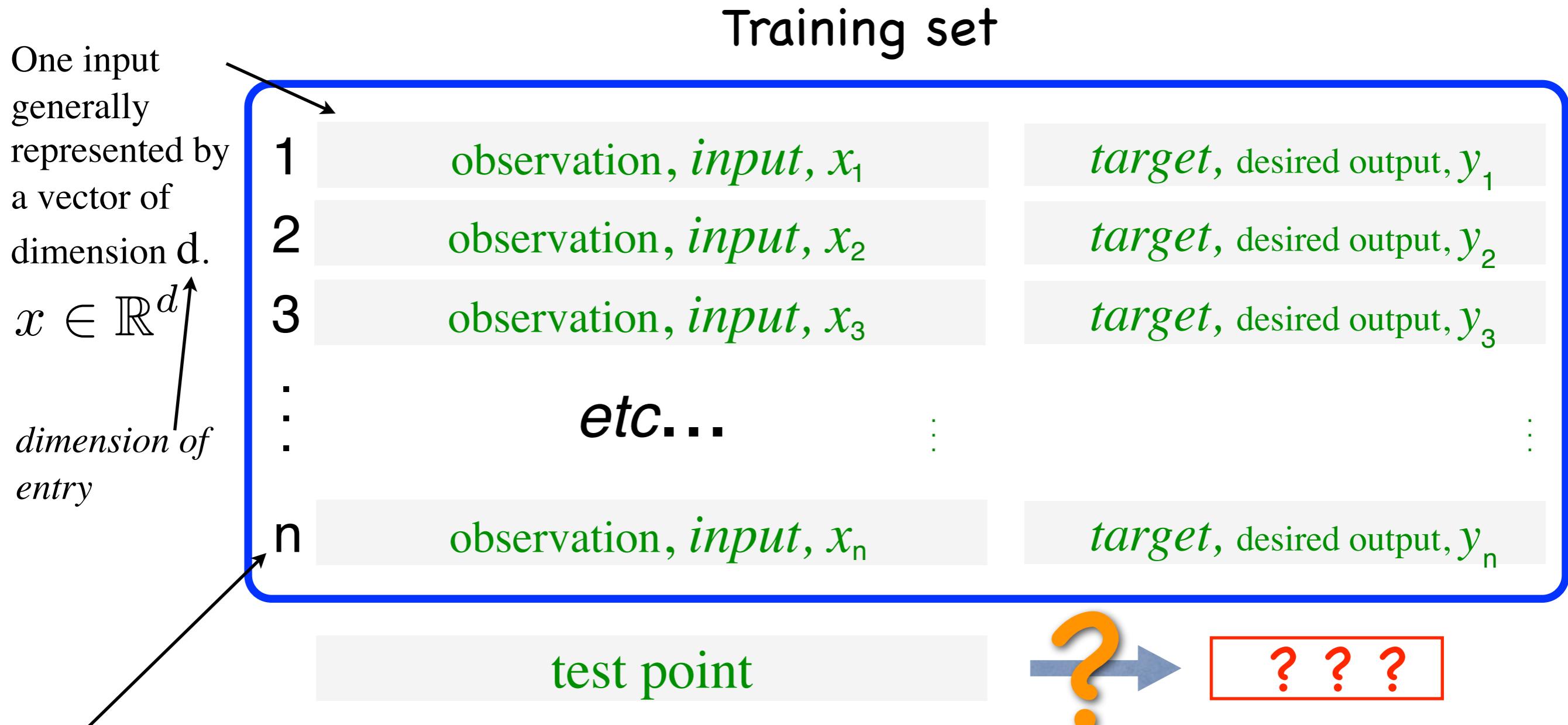
# Terminology of supervised learning



# Terminology of supervised learning



# Terminology of supervised learning



We are looking for an algorithm that produces an output that is a good prediction of the target.

This algorithm finds a good function  $x \rightarrow y$

# Terminology of supervised learning

- o When the target is a class label, a categorical variable (indicating which class or category the entry belongs to, among many) is said to be a **CLASSIFICATION** problem. (we often use an integer as a label).
- o When the target is one (or more) real values to predict, we talk about a **REGRESSION** problem.

When we do not have an explicit target, we are in the context of **unsupervised learning**.

# Learning tasks

Supervised learning = predict a target  $y$  from the input  $x$

# Learning tasks

**Supervised learning** = predict a target  $y$  from the input  $x$

**Unsupervised learning:** no explicit target (label)  $y$

# Learning tasks

Supervised learning = predict a target  $y$  from the input  $x$

- $y$  represents a category or “class”  
→ classification (binary or multiclass)
- $y$  is a real value  
→ regression

Unsupervised learning: no explicit target (label)  $y$

# Learning tasks

Supervised learning = predict a target  $y$  from the input  $x$

- $y$  represents a category or “class”  
    ⇒ classification (binary or multiclass)
- $y$  is a real value  
    ⇒ regression

Unsupervised learning: no explicit target (label)  $y$

- modeling the distribution of  $x$   
    ⇒ density estimation
- discover an underlying structure in  $x$   
    ⇒ clustering  
    ⇒ dimensionality reduction (eg. for visualisation)

# Learning tasks

Supervised learning = predict a target  $y$  from the input  $x$

- $y$  represents a category or “class”  
    ⇒ classification (binary or multiclass)
- $y$  is a real value  
    ⇒ regression

} predictive models

Unsupervised learning: no explicit target (label)  $y$

- modeling the distribution of  $x$   
    ⇒ density estimation
- discover an underlying structure in  $x$   
    ⇒ clustering  
    ⇒ dimensionality reduction (eg. for visualisation)

# Learning tasks

Supervised learning = predict a target  $y$  from the input  $x$

- $y$  represents a category or “class”  
    ⇒ classification (binary or multiclass)
- $y$  is a real value  
    ⇒ regression

} predictive models

Unsupervised learning: no explicit target (label)  $y$

- modeling the distribution of  $x$   
    ⇒ density estimation
- discover an underlying structure in  $x$   
    ⇒ clustering  
    ⇒ dimensionality reduction (eg. for visualisation)

} descriptive models

# Learning tasks

## Semi-supervised learning:

same goal as supervised learning, but using both labeled (with target) and unlabeled examples.

## Reinforcement learning

an artificial agent must learn to decide what actions to perform in a changing environment to maximize a total reward.

# Data set seen as a scatter plot in a high-dimensional vector space

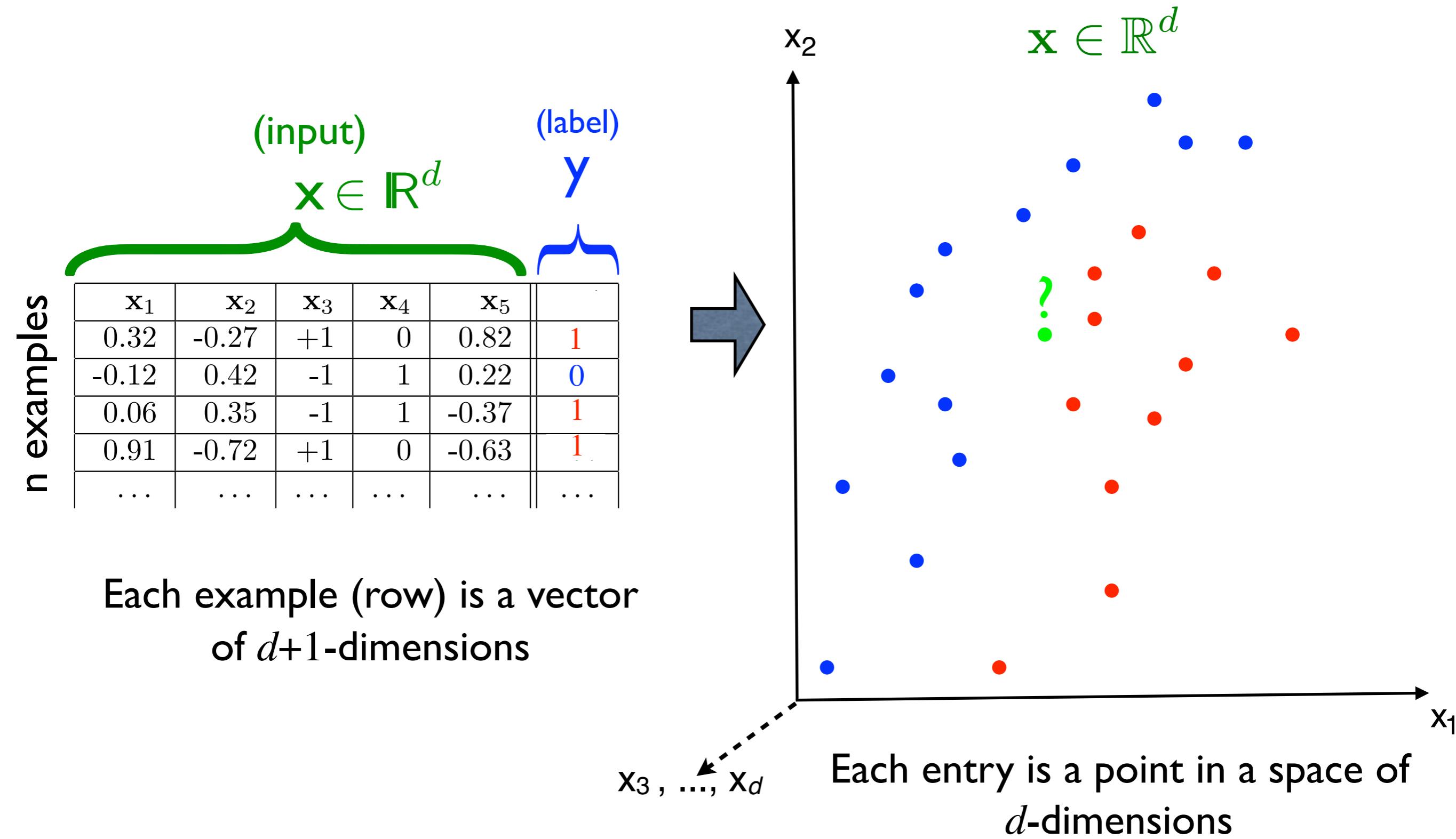
(input)  $x \in \mathbb{R}^d$  (label)  $y$

n examples

$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	
0.32	-0.27	+1	0	0.82	1
-0.12	0.42	-1	1	0.22	0
0.06	0.35	-1	1	-0.37	1
0.91	-0.72	+1	0	-0.63	1
...	...	...	...	...	...

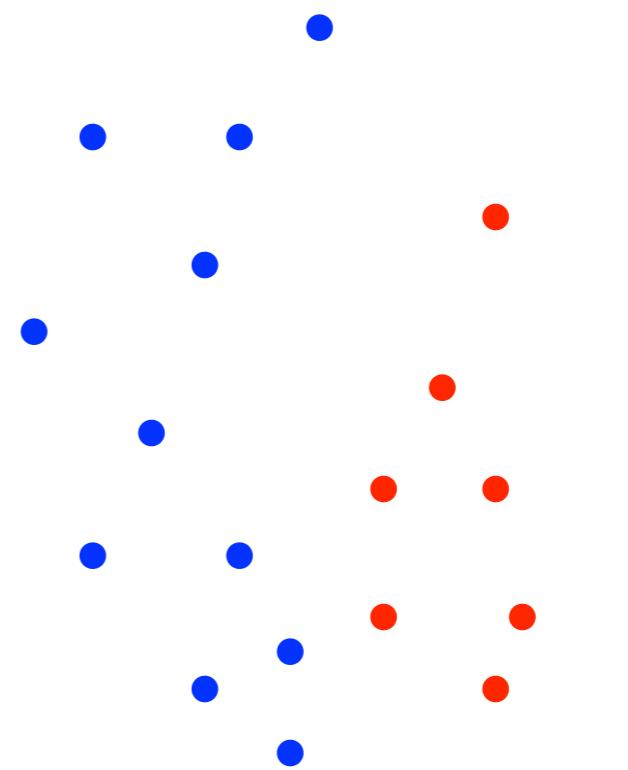
Each example (row) is a vector  
of  $d+1$ -dimensions

# Data set seen as a scatter plot in a high-dimensional vector space



# The nearest neighbor algorithm

---

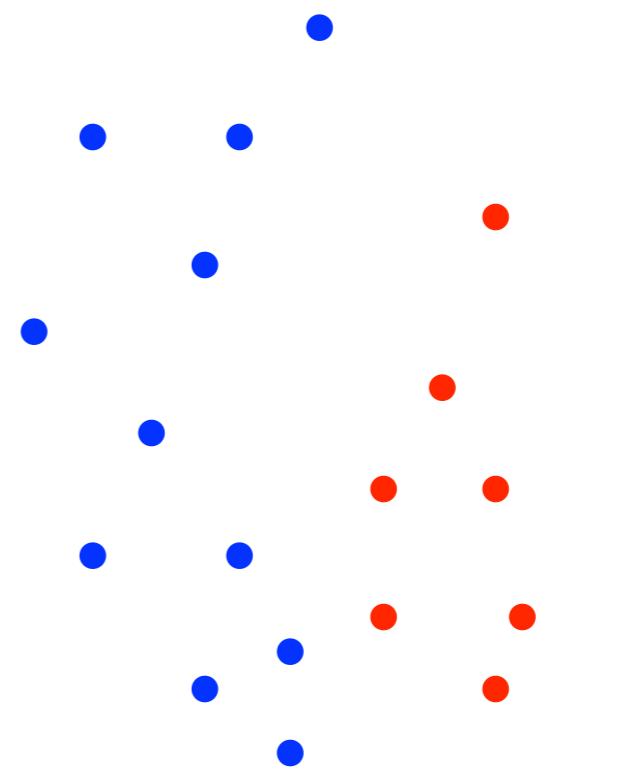


Training set

# The nearest neighbor algorithm

---

For a test point  $\mathbf{x}$ :

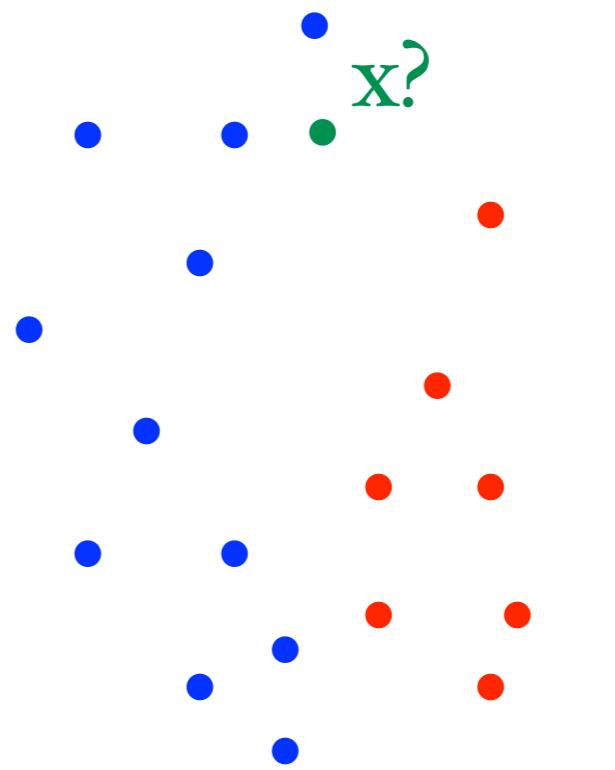


Training set

# The nearest neighbor algorithm

---

For a test point  $\mathbf{x}$ :



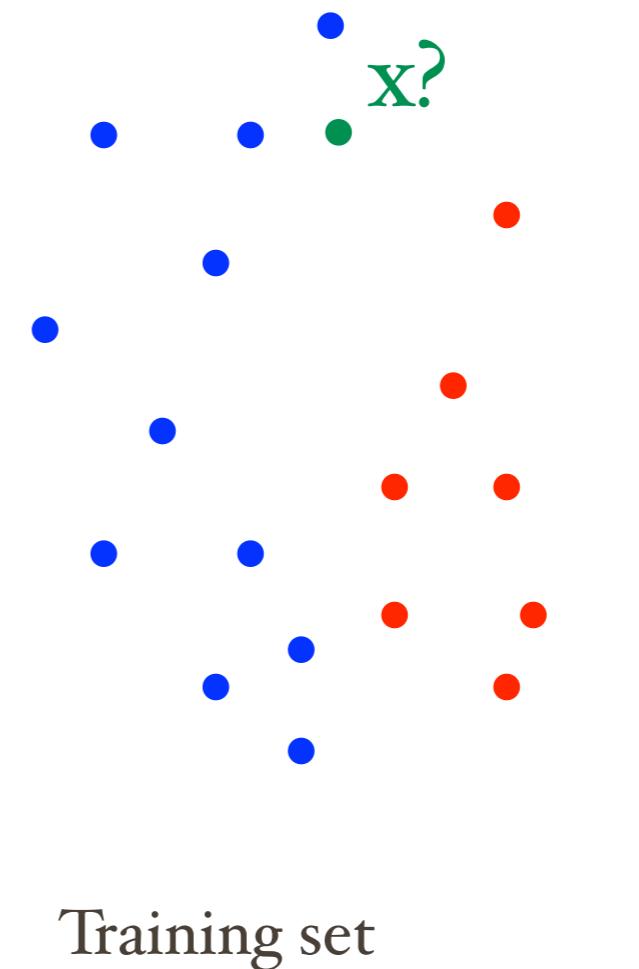
Training set

# The nearest neighbor algorithm

---

For a test point  $x$ :

- We find the **nearest neighbour** of  $x$  within the training set by some measure of distance (eg Euclidean distance).

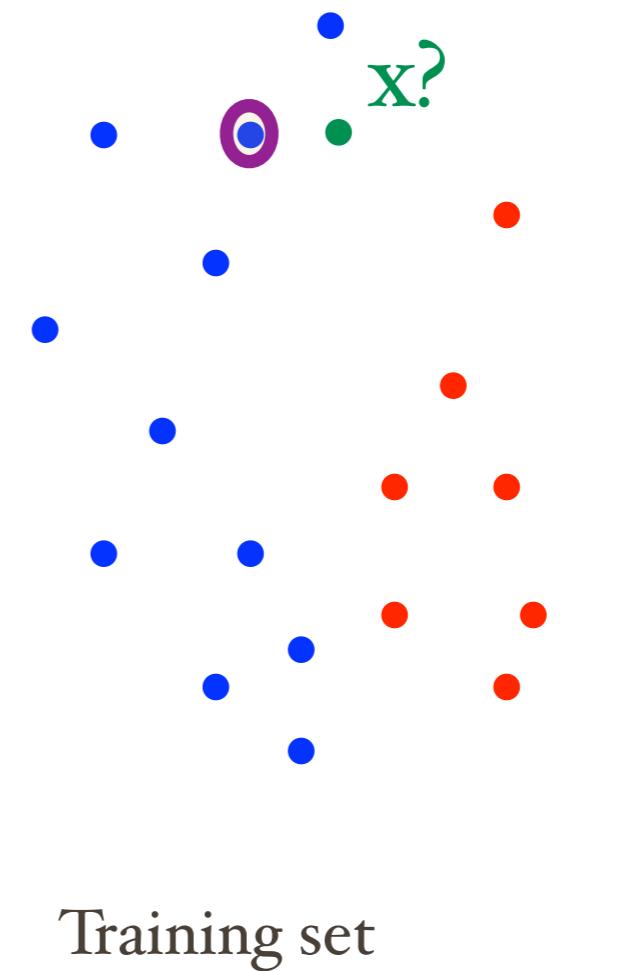


# The nearest neighbor algorithm

---

For a test point  $x$ :

- We find the **nearest neighbour** of  $x$  within the training set by some measure of distance (eg Euclidean distance).

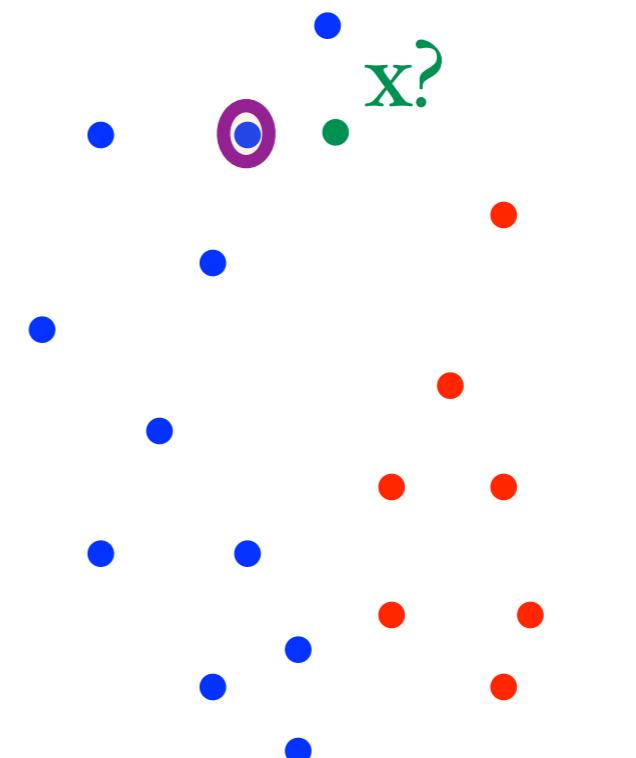


# The nearest neighbor algorithm

---

For a test point  $x$ :

- We find the **nearest neighbour** of  $x$  within the training set by some measure of distance (eg Euclidean distance).
- We associate  $x$  with the class of this nearest neighbor.



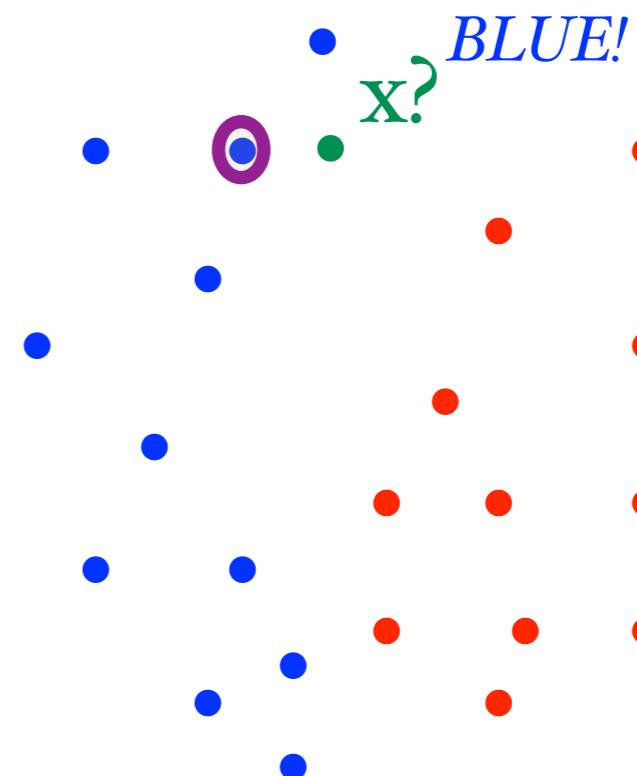
Training set

# The nearest neighbor algorithm

---

For a test point  $x$ :

- We find the **nearest neighbour** of  $x$  within the training set by some measure of distance (eg Euclidean distance).
- We associate  $x$  with the class of this nearest neighbor.

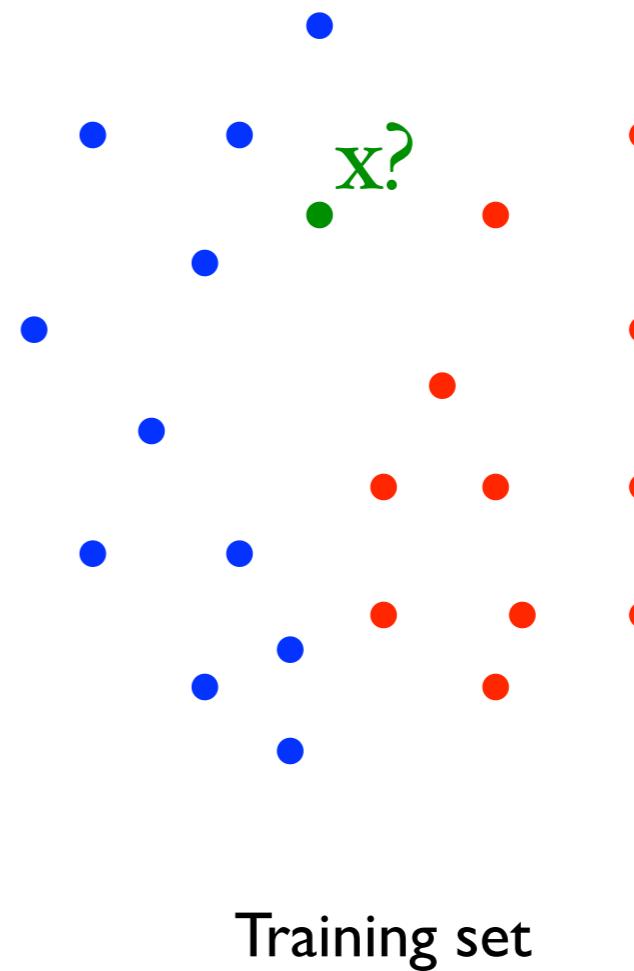


# The k-nearest neighbors algorithm (k-NN)

For a test point  $x$ :

$k=3$

- We find the  **$k$  nearest neighbors** of  $x$  within the training set apprentissage by some measure of distance (eg Euclidean distance).
- We associate  $x$  with the class of the **majority** of neighbors.

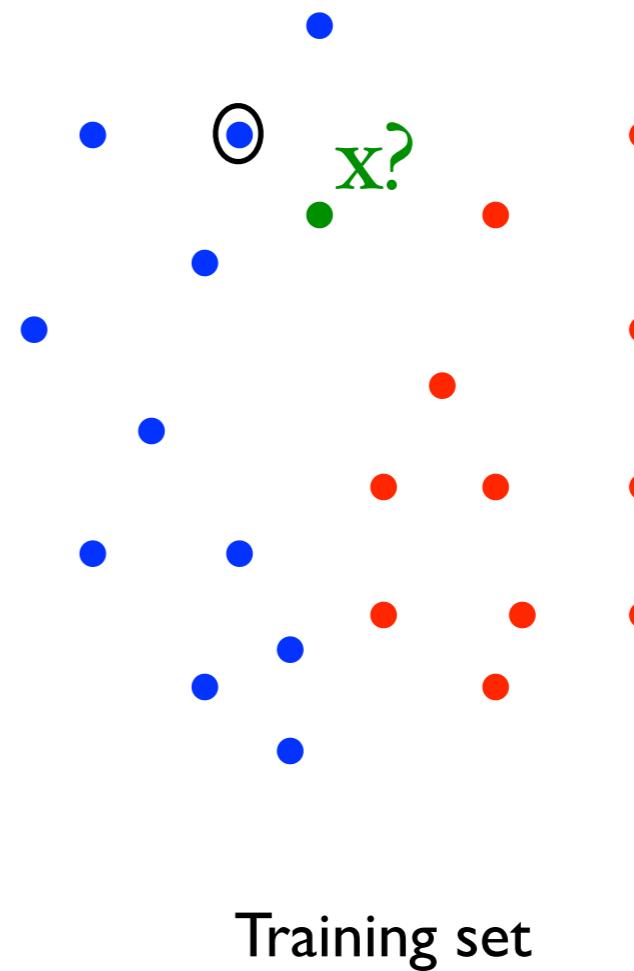


# The k-nearest neighbors algorithm (k-NN)

For a test point  $x$ :

$k=3$

- We find the  **$k$  nearest neighbors** of  $x$  within the training set apprentissage by some measure of distance (eg Euclidean distance).
- We associate  $x$  with the class of the **majority** of neighbors.

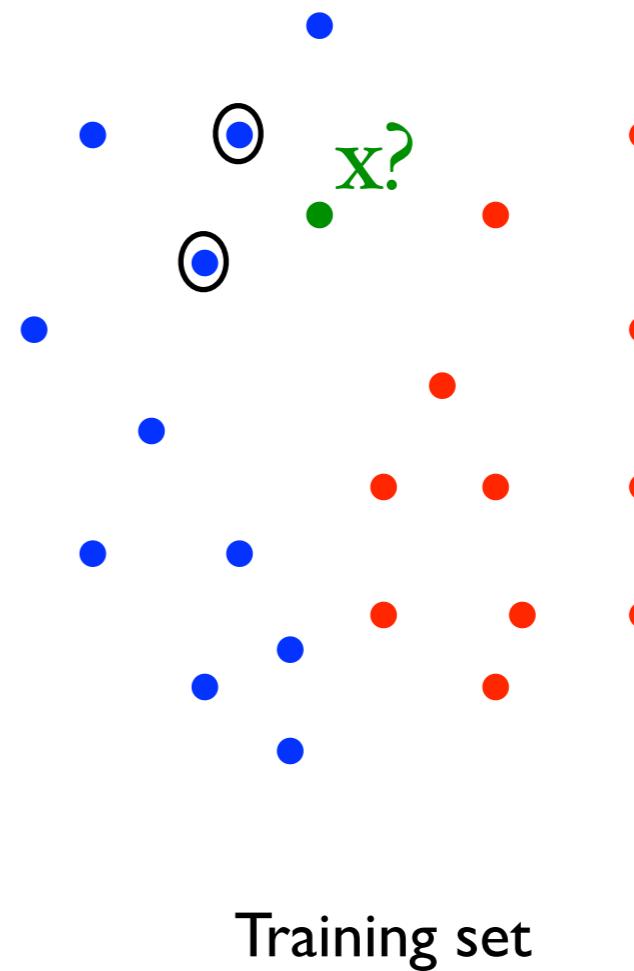


# The k-nearest neighbors algorithm (k-NN)

For a test point  $x$ :

$k=3$

- We find the  **$k$  nearest neighbors** of  $x$  within the training set apprentissage by some measure of distance (eg Euclidean distance).
- We associate  $x$  with the class of the **majority** of neighbors.

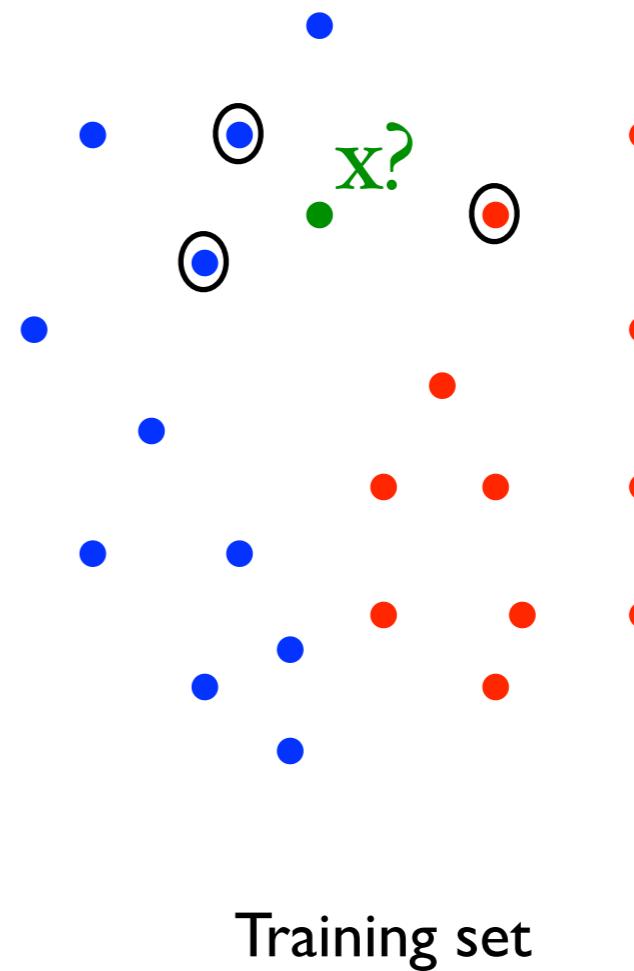


# The k-nearest neighbors algorithm (k-NN)

For a test point  $x$ :

$k=3$

- We find the  **$k$  nearest neighbors** of  $x$  within the training set apprentissage by some measure of distance (eg Euclidean distance).
- We associate  $x$  with the class of the **majority** of neighbors.

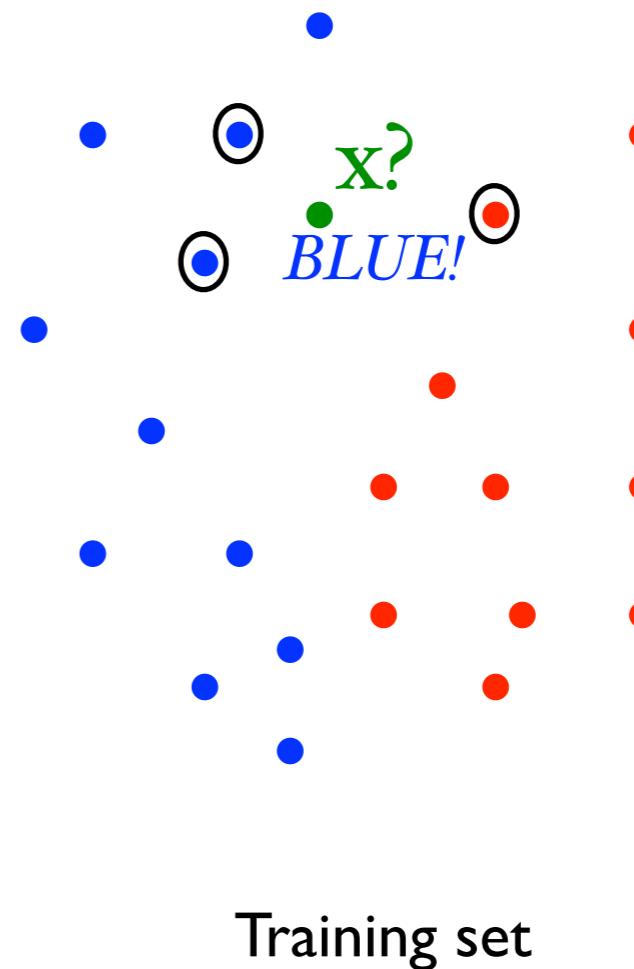


# The k-nearest neighbors algorithm (k-NN)

For a test point  $x$ :

$k=3$

- We find the  **$k$  nearest neighbors** of  $x$  within the training set apprentissage by some measure of distance (eg Euclidean distance).
- We associate  $x$  with the class of the **majority** of neighbors.

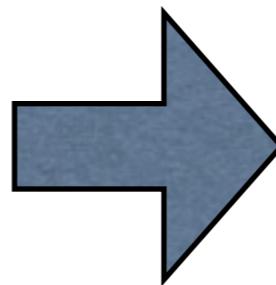


# Short review of vectors

(of dimension d)

- what is a vector? mathematical, computer, graphic representations  
...
- Euclidean distance
- norms
- scalar product
- efficient calculation of the distance from a point to a set of points

# The notion of level of representation



very high level representation:

CAT      JUMPING

... etc ...

slightly higher level representation

raw input vector representation:

$$\mathcal{X} = \begin{bmatrix} 23 & 19 & 20 & \cdots & 18 \end{bmatrix}$$

$x_1 \quad x_2 \quad x_3 \quad \vdots \quad x_n$

