

CIFAR
CANADIAN
INSTITUTE
FOR
ADVANCED
RESEARCH



Deep Learning Summer School 2015

Introduction to Machine Learning

by **Pascal Vincent**



Montreal Institute for Learning Algorithms

August 4, 2015

Université
de Montréal

Département d'informatique et de recherche
opérationnelle

What is machine learning ?

Historical perspective

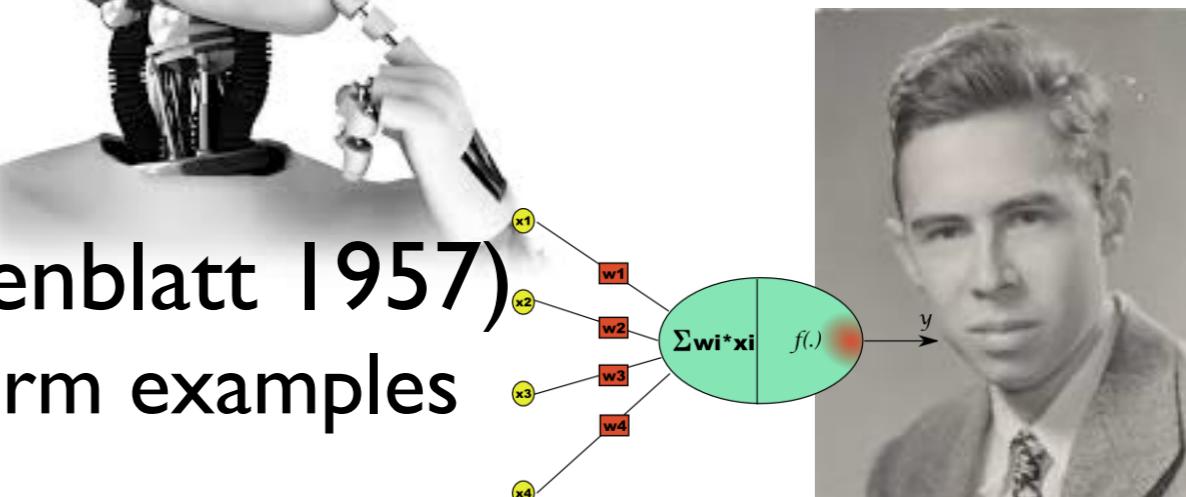
- Born from the ambitious goal of **Artificial Intelligence**



- Founding project:

The **Perceptron** (Frank Rosenblatt 1957)

First artificial neuron **learning** form examples



- Two historically opposed approaches to AI:

Neuroscience inspired:

→ neural nets **learning from examples** for artificial perception

Classical symbolic AI:

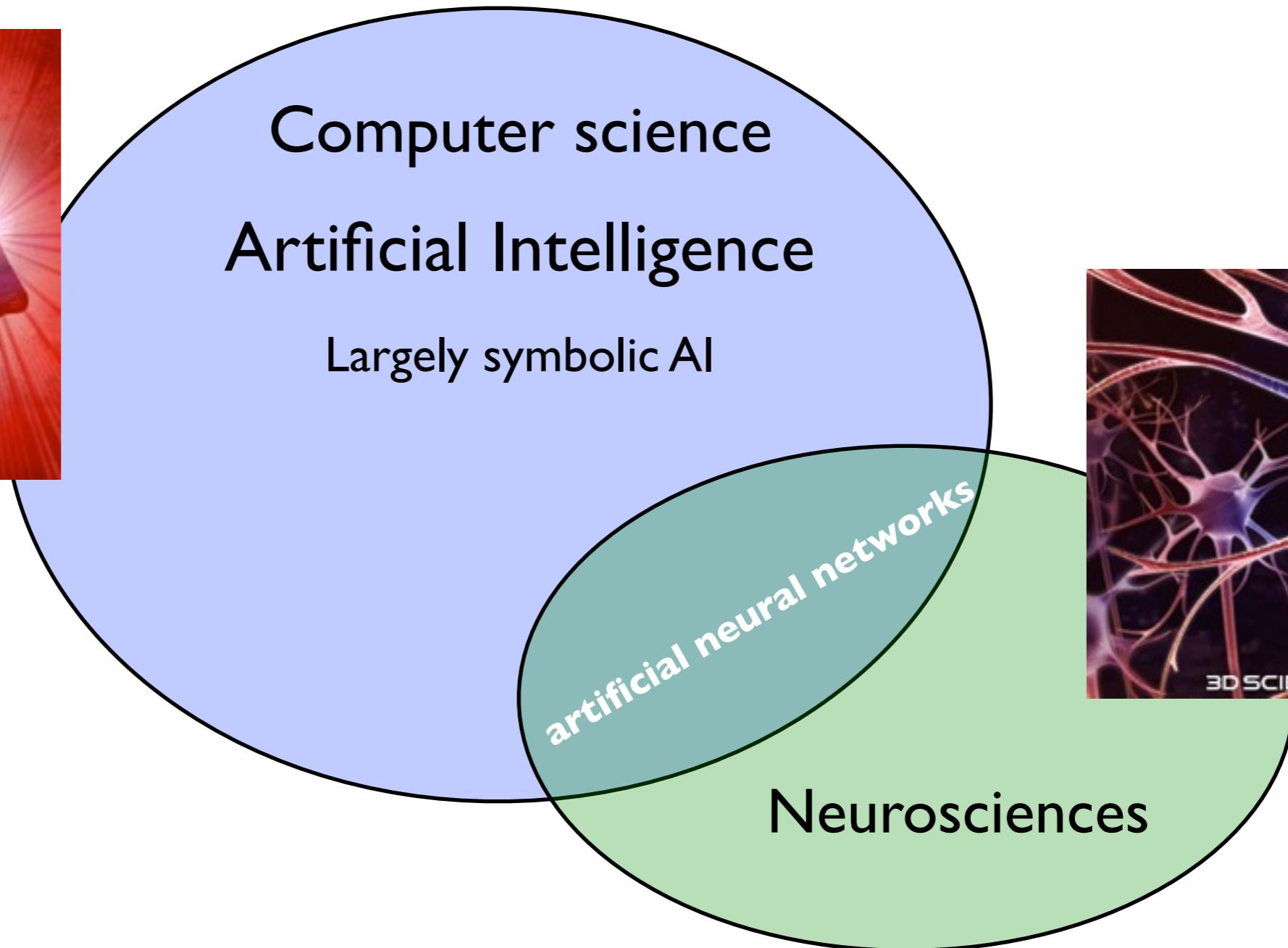
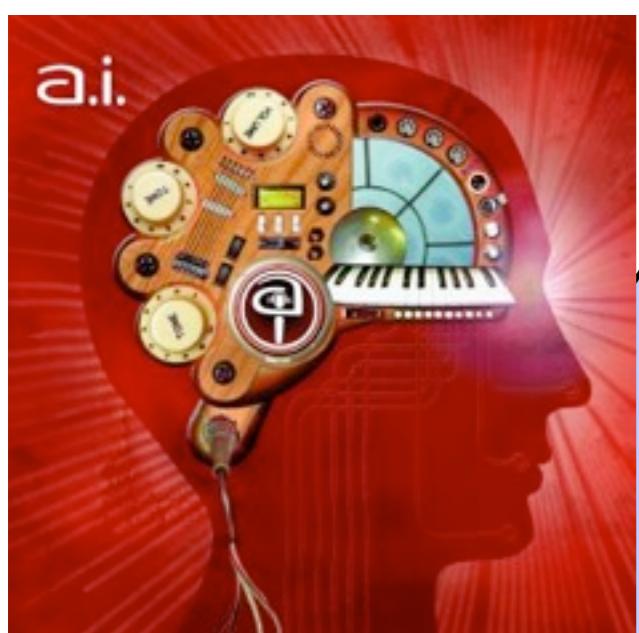
Primacy of **logical reasoning** capabilities

→ No learning (humans coding rules)
→ poor handling of uncertainty

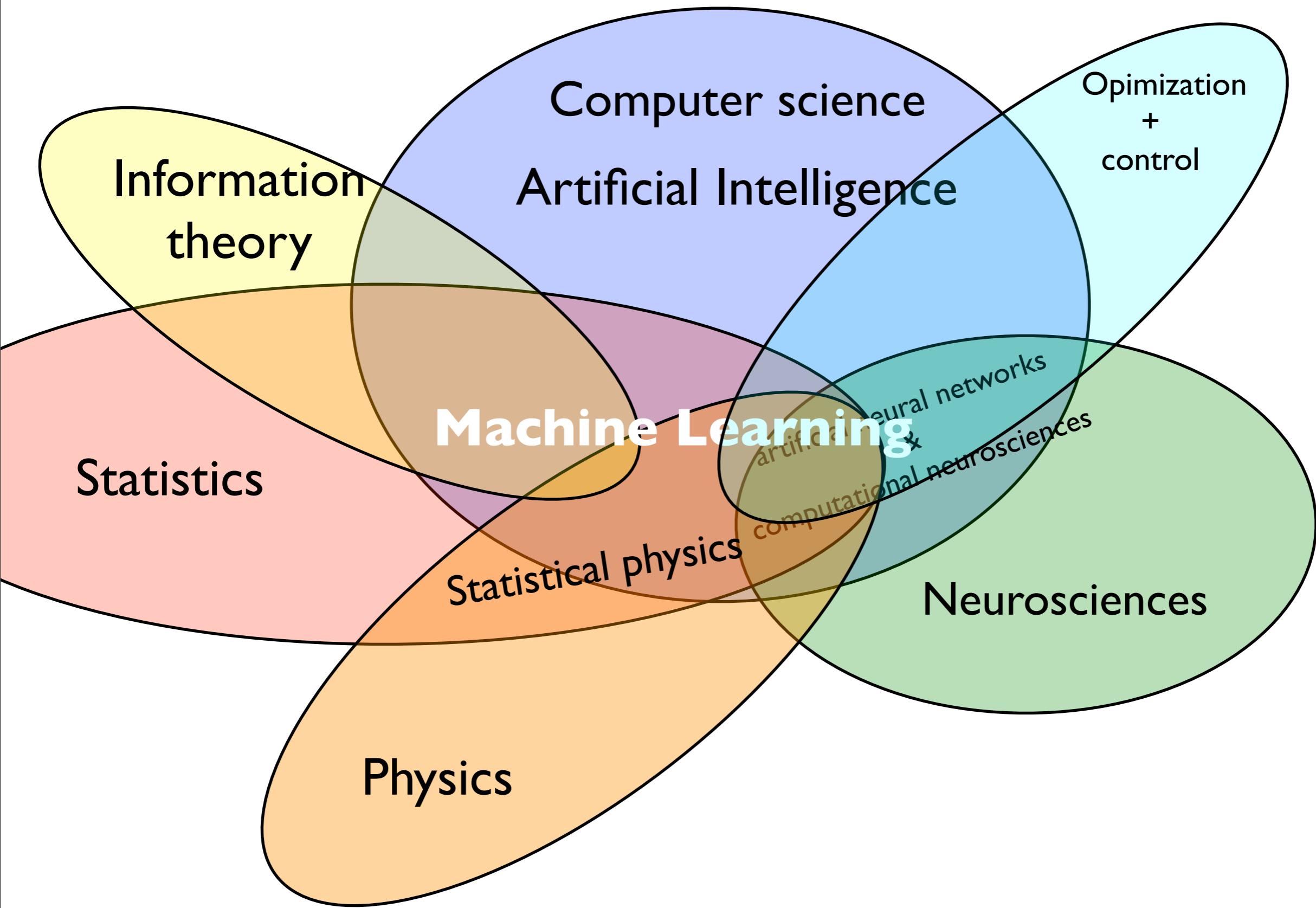
Got eventually fixed (Bayes Nets...)

Learning and probabilistic models largely won → machine learning

Artificial Intelligence in the 60s



Current view of ML founding disciplines



What is machine-learning?



A (hypnotized) user's perspective

A scientific (witchcraft) field that

- researches fundamental principles (potions)
- and **develops** magical algorithms (spells to invoke)
- capable of leveraging collected data to (automagically)
produce accurate predictive functions
applicable to similar data (in the future!)

(may also yield informative descriptive functions of data)

The key ingredient of machine learning is...



Data!

- Collected from nature... or industrial processes.
- Comes stored in many forms (and formats...), structured, unstructured, occasionally clean, **usually messy**, ...
- In ML we like to view data as a **list of examples**
 (or we'll turn it into one)
 - ➡ ideally **many examples** of the **same nature**.
 - ➡ preferably with each example a **vector of numbers**
(or we'll first turn it into one!)

D_n

Training data set (training set)



New test point:



→ ?

Input dimensionality:

d

inputs: X (input feature vector)	targets: Y (label)
X_1 (3.5, -2, ..., 127, 0, ...)	+1
X_n (-9.2, 32, ..., 24, 1, ...)	-1
$X_{n,2}$ (6.8, 54, ..., 17, -3, ...)	+1

$$X = (5.7, -27, \dots, 64, 0, \dots) \xrightarrow{f_\theta} +1$$

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

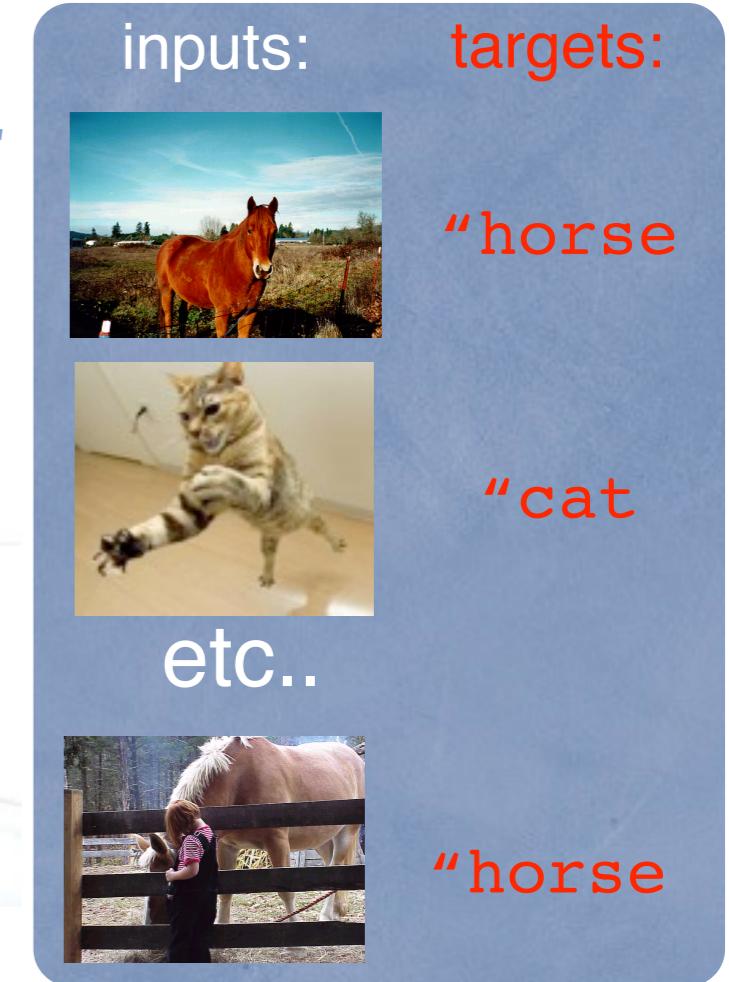
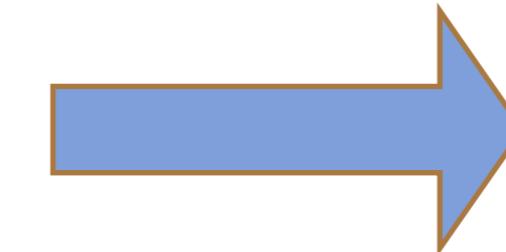
Importance of the Problem dimensions

- ⇒ Détermines which learning algorithms will be practically applicable (based on their algorithmic complexity and memory requirements).
 - Number of examples: **n**
(sometimes several millions)
 - Input dimensionality: **d**
number of input features characterizing each example
(often 100 to 1000, sometimes 10000 or much more)
 - Target dimensionality ex. number of classes **m**
(often small, sometimes huge)
- Data suitable for ML will often be organized as a matrix: **$n \times (d+1)$** ou **$n \times (d+m)$**

Turning ~~messy~~ data into a nice list of examples



data-plumbing

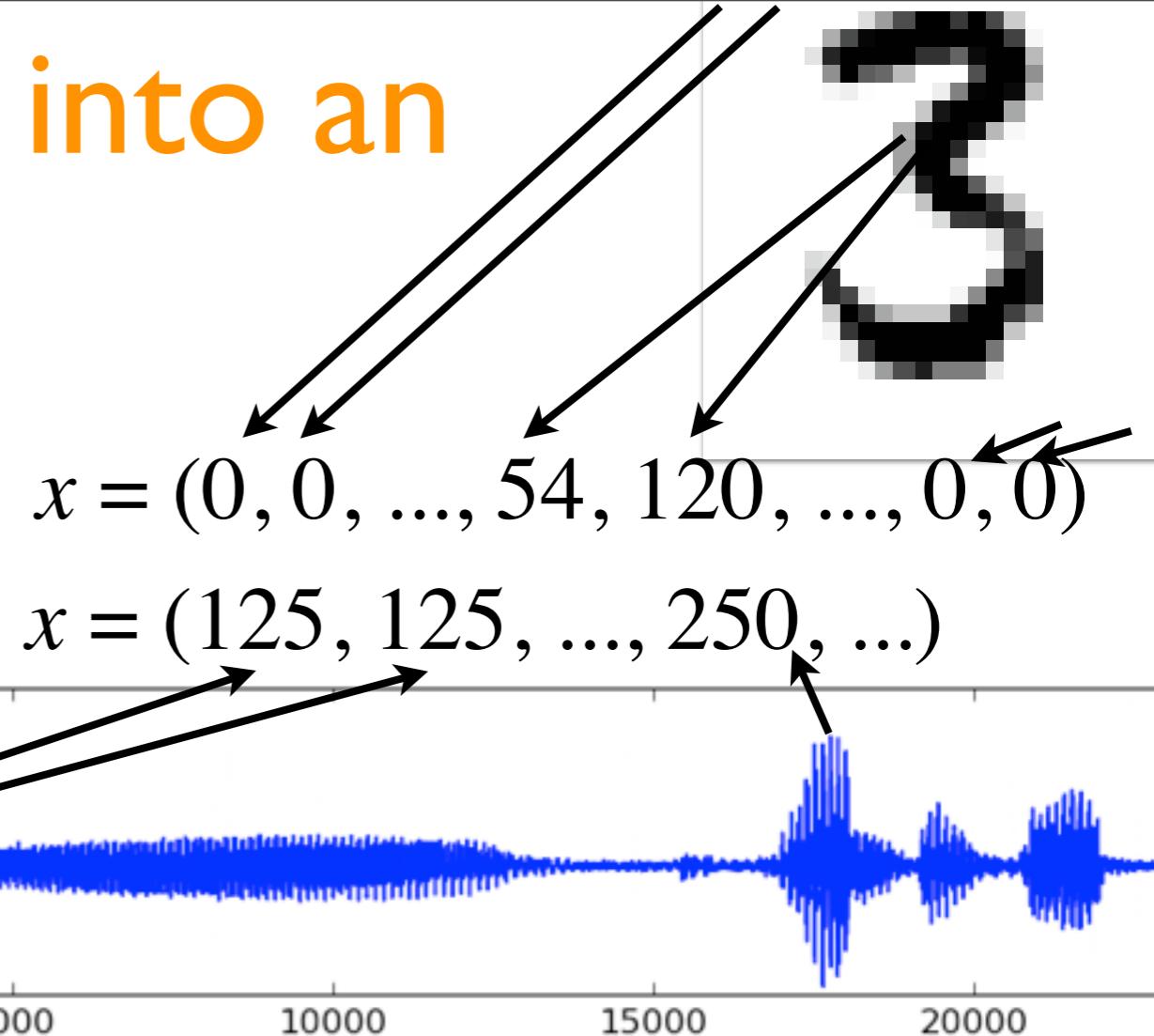


Key questions to decide what «examples» should be:

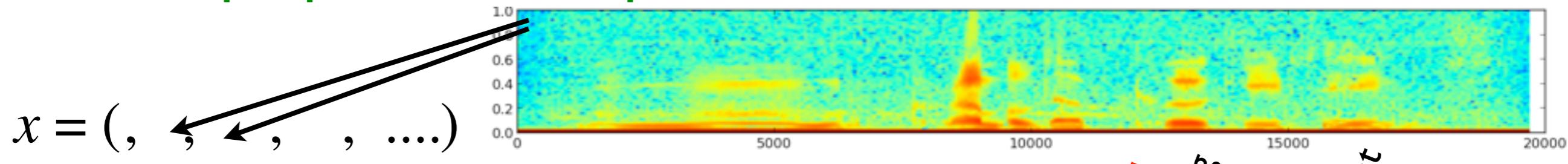
- **input:** What is all the (potentially relevant) **information** I will have at my disposal about a case when I will have to make a prediction about it? (at test time)
- **target:** what I want to predict: Can I get my hands on many such examples that are actually labeled with prediction targets?

Turning an example into an input vector $\mathbf{x} \in \mathbb{R}^d$

Raw input representation:



OR some preprocessed representation:

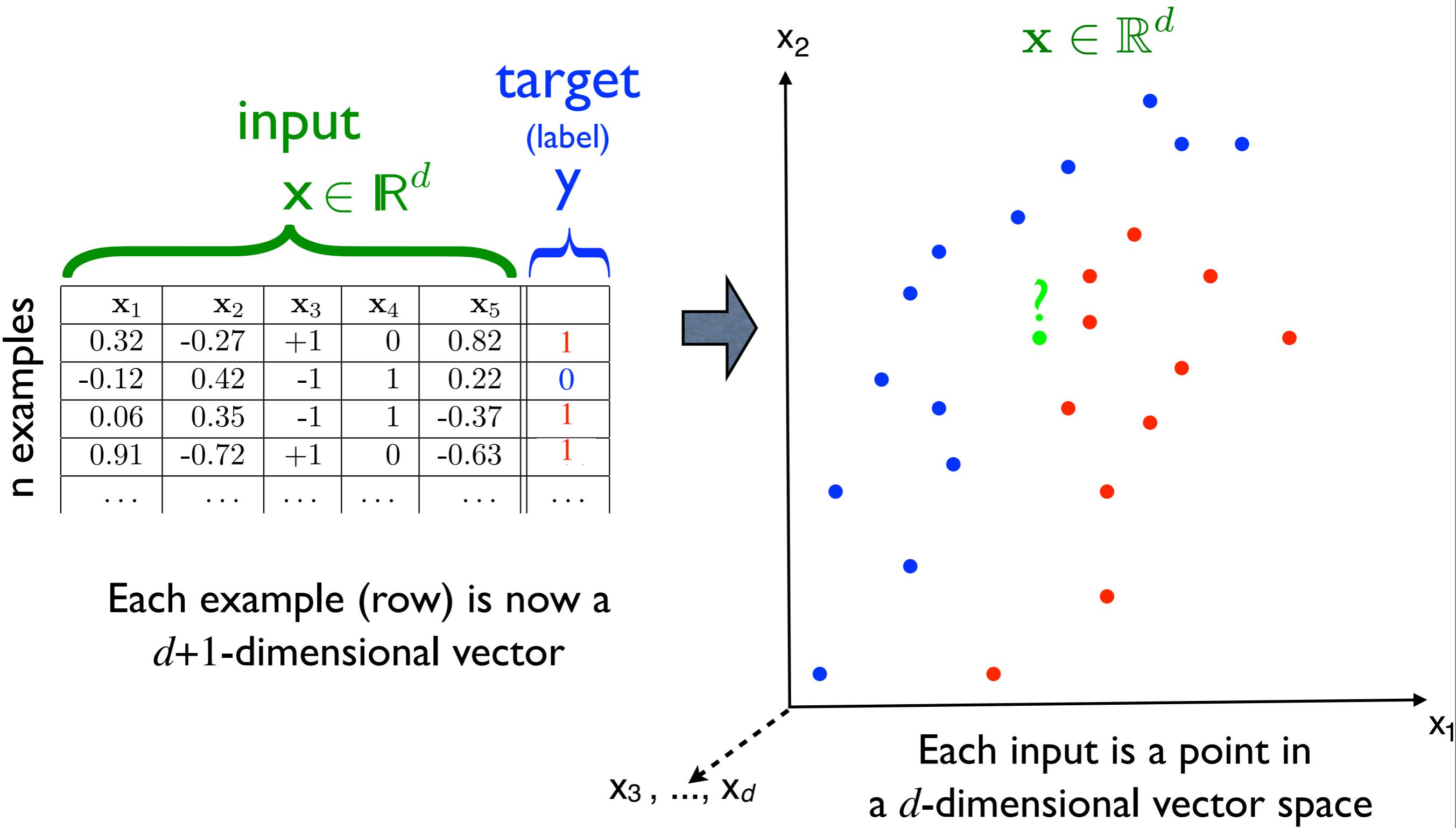


Bag of words for «The cat jumped»: $x = (\dots 0 \dots, 0, 1, \dots 0 \dots, 1, 0, 0, \dots, 0, 0, 1, 0, \dots 0 \dots)$

OR vector of hand-engineered features:
ex: Histograms of Oriented Gradients

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

Dataset imagined as a point cloud in a *high-dimensional* vector space

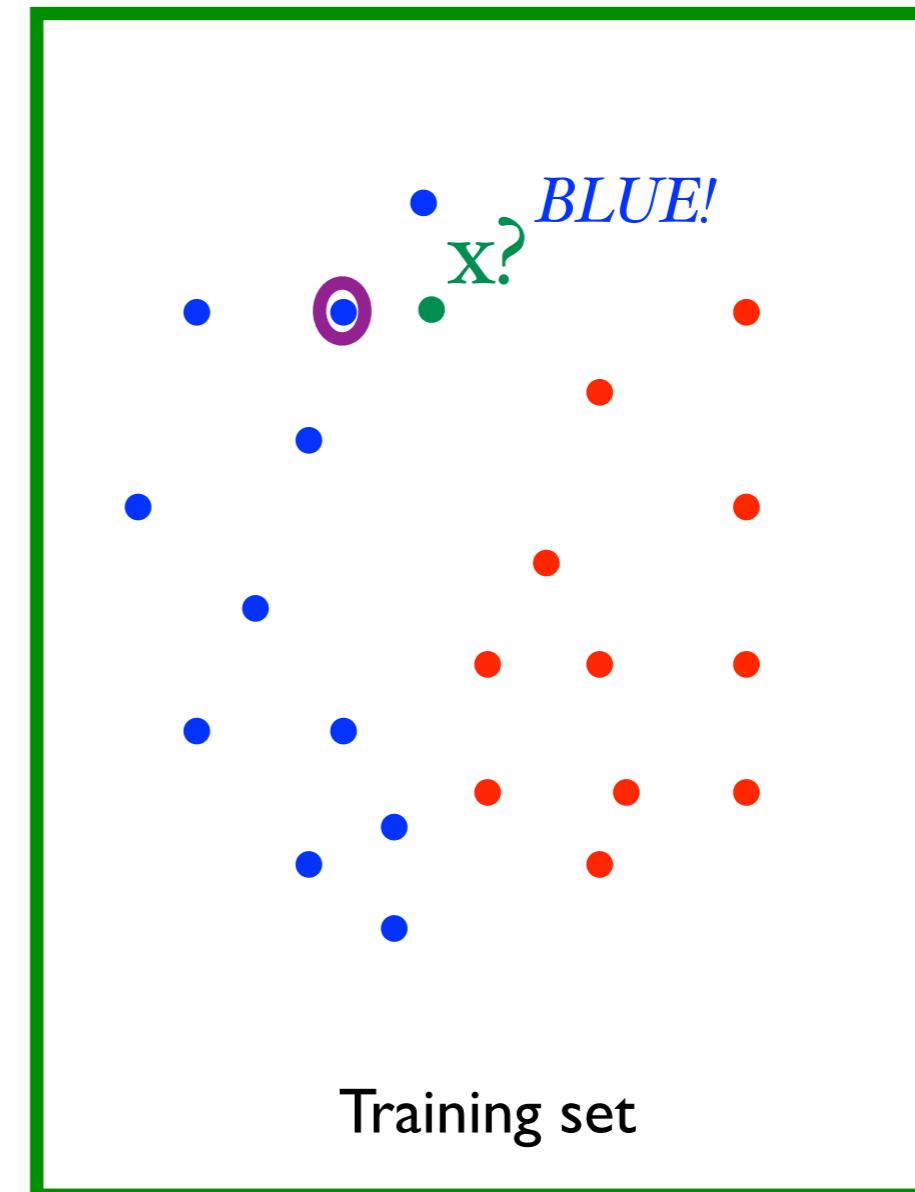


Ex: nearest-neighbor classifier

Algorithm:

For test point x :

- Find **nearest neighbor** of x among the training set according to some distance measure (eg: Euclidean distance).
- Predict that x has the same class as this nearest neighbor.



Machine learning tasks (problem types)

Supervised learning = predict a target y from input x
(and semi-supervised learning)

- y represents a category or “class”

 ⇒ classification

 binary : $y \in \{-1, +1\}$ or $y \in \{0, 1\}$

 multiclass : $y \in \{1, m\}$ or $y \in \{0, m - 1\}$

- y is a real-value number

 ⇒ regression

$y \in \mathbb{R}$ or $y \in \mathbb{R}^m$

} Predictive models

Unsupervised learning: no explicit prediction target y

- model the probability distribution of x

 ⇒ density estimation

- discover underlying structure in data

 ⇒ clustering

 ⇒ dimensionality reduction

 ⇒ (unsupervised) representation learning

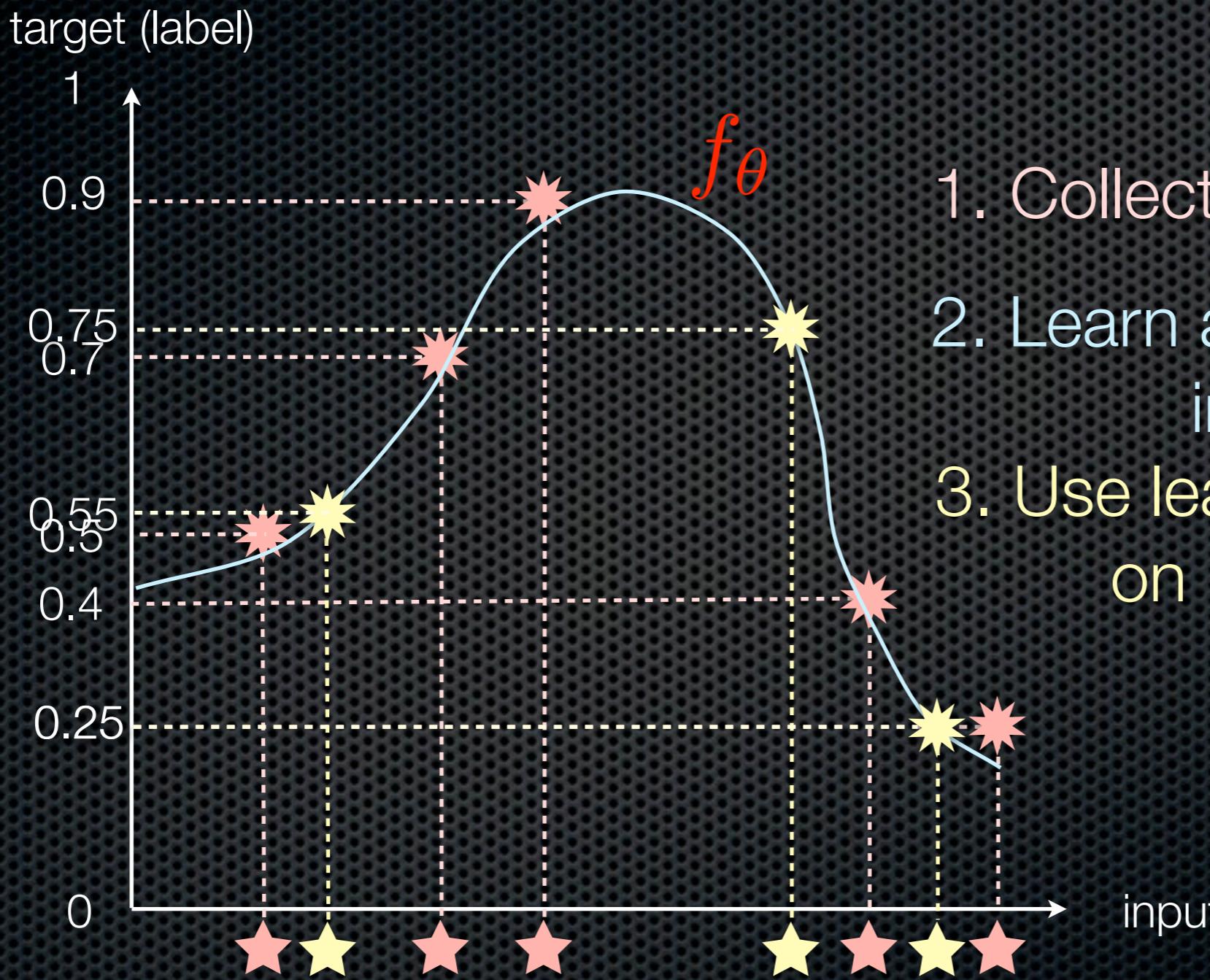
} Descriptive modeling

Reinforcement learning: taking good sequential decisions to maximize a reward
in an environment influenced by your decisions.

Learning phases

- **Training:** we **learn** a predictive function f_θ by optimizing it so that it predicts well on **the training set**.
 - **Use for prediction:** we can then use f_θ on new (test) inputs that were not part of the training set.
- ⇒ The GOAL of learning is *NOT* to learn perfectly (*memorize*) the training set.
- ⇒ What's important is the ability for the predictor to **generalize** well on new (future) cases.

Ex: 1D regression

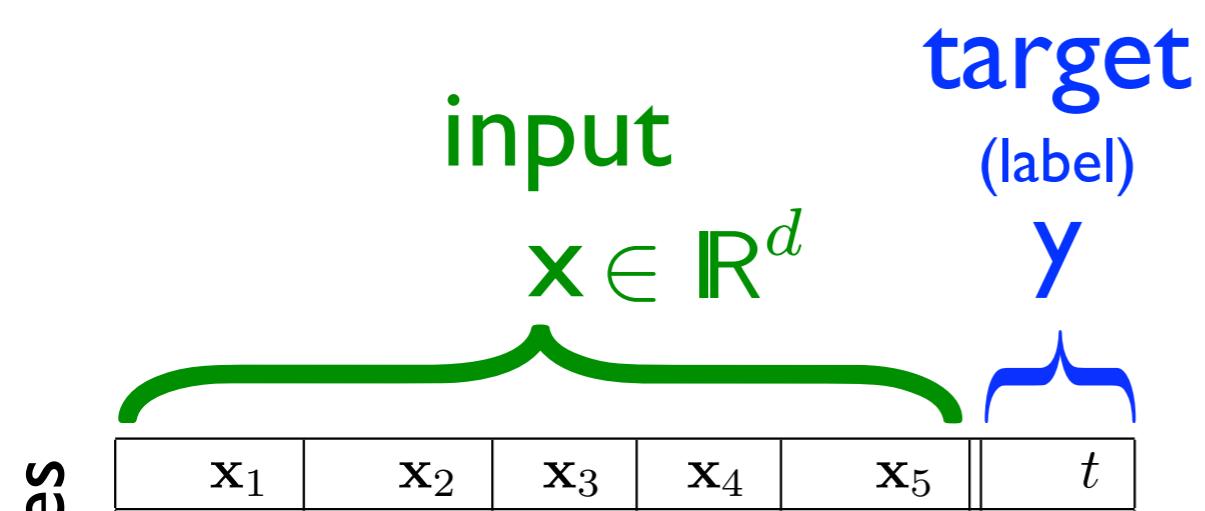


1. Collect training data
2. Learn a function (predictor)
input → target
3. Use learned function
on new inputs

Original slide by Olivier Delalleau

Supervised task:

predict y from x



Training set D_n

Learn a function f_θ that will minimize prediction errors as measured by cost (loss) L .

loss function

$$L(f_\theta(\mathbf{x}), y)$$

output

$$f_\theta(\mathbf{x})$$

$$f_\theta$$

: parameters

-0.12	0.42	-1	1	0.22
-------	------	----	---	------

34

target y

input \mathbf{x}

A machine learning algorithm
usually corresponds to a combination of
the following 3 elements:
(either explicitly specified or implicit)

- ✓ the choice of a specific **function family**: F
(often a parameterized family)
- ✓ a **way to evaluate the quality** of a function $f \in F$
(typically using a **cost** (or **loss**) function L
measuring how wrongly f predicts)
- ✓ a **way to search for the «best»** function $f \in F$
(typically an optimization of function parameters to
minimize the overall loss over the training set).

Evaluating the quality of a function $f \in F$

and

Searching for the «best» function $f \in F$

Evaluating a predictor $f(x)$

The performance of a predictor is often **evaluated using several different evaluation metrics**:

- Evaluations of **true quantities of interest** (\$ saved, #lifes saved, ...) when using predictor inside a more complicated system.
- «Standard» evaluation metrics in a specific field (e.g. BLEU (Bilingual Evaluation Understudy) scores in translation)
- Misclassification error rate for a classifier (or precision and recall, or F-score, ...).
- **The loss actually being optimized** by the ML algorithm (often different from all the above...)

Standard loss-functions

- **For a density estimation task:** $f : \mathbb{R}^d \rightarrow \mathbb{R}^+$ a proper probability mass or density function
negative log likelihood loss: $L(f(x)) = -\log f(x)$
- **For a regression task:** $f : \mathbb{R}^d \rightarrow \mathbb{R}$
squared error loss: $L(f(x), y) = (f(x) - y)^2$
- **For a classification task:** $f : \mathbb{R}^d \rightarrow \{0, \dots, m-1\}$
misclassification error loss: $L(f(x), y) = I_{\{f(x) \neq y\}}$

Surrogate loss-functions

- For a classification task: $f : \mathbb{R}^d \rightarrow \{0, \dots, m-1\}$
 misclassification error loss: $L(f(x), y) = I_{\{f(x) \neq y\}}$

Problem: it is hard to optimize the misclassification loss directly
 (gradient is 0 everywhere. NP-hard with a linear classifier) Must use a surrogate loss:

	Binary classifier	Multiclass classifier
Probabilistic classifier	<p>Outputs probability of class 1 $g(x) \approx P(y=1 x)$ Probability for class 0 is $1-g(x)$</p> <p><u>Binary cross-entropy loss:</u> $L(g(x), y) = -(y \log(g(x)) + (1-y) \log(1-g(x)))$</p> <p>Decision function: $f(x) = I_{g(x)>0.5}$</p>	<p>Outputs a vector of probabilities: $g(x) \approx (P(y=0 x), \dots, P(y=m-1 x))$</p> <p><u>Negated conditional log likelihood loss</u> $L(g(x), y) = -\log g(x)_y$</p> <p>Decision function: $f(x) = \text{argmax}(g(x))$</p>
Non-probabilistic classifier	<p>Outputs a «score» $g(x)$ for class 1. score for the other class is $-g(x)$</p> <p><u>Hinge loss:</u> $L(g(x), t) = \max(0, 1-tg(x))$ where $t=2y-1$</p> <p>Decision function: $f(x) = I_{g(x)>0}$</p>	<p>Outputs a vector $g(x)$ of real-valued scores for the m classes.</p> <p><u>Multiclass margin loss</u> $L(g(x), y) = \max(0, 1 + \max_{k \neq y} (g(x)_k - g(x)_y))$</p> <p>Decision function: $f(x) = \text{argmax}(g(x))$</p>

Expected risk v.s. Empirical risk

Examples (\mathbf{x}, \mathbf{y}) are supposed drawn i.i.d. from an **unknown true distribution** $p(\mathbf{x}, \mathbf{y})$ (from nature or industrial process)

- **Generalization error** = Expected risk (or just «Risk»)
«how poorly we will do on average on the infinity of future examples from that unknown distribution»

$$R(f) = \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [L(f(\mathbf{x}), \mathbf{y})]$$

- **Empirical risk** = average loss on a finite dataset
«how poorly we're doing on average on this finite dataset»

$$\hat{R}(f, D) = \frac{1}{|D|} \sum_{(\mathbf{x}, \mathbf{y}) \in D} L(f(\mathbf{x}), \mathbf{y})$$

where $|D|$ is the number of examples in D

Empirical risk minimization

Examples (x,y) are supposed drawn i.i.d. from an unknown true distribution $p(x,y)$ (nature or industrial process)

- We'd **love** to find a predictor that **minimizes the generalization error** (the expected risk)
- But **can't even compute it!** (expectation over unknown distribution)
- Instead: **Empirical risk minimization principle**
«Find predictor that minimizes average loss over a trainset»

$$\hat{f}(D_{\text{train}}) = \underset{f \in F}{\operatorname{argmin}} \hat{R}(f, D_{\text{train}})$$

This is the training phase in ML

Evaluating the generalization error

- ▶ We can't compute expected risk $R(f)$
- ▶ But $\hat{R}(f, D)$ is a good estimate of $R(f)$ provided:
 - D was not used to find/choose f otherwise estimate is biased \Rightarrow can't be the training set!
 - D is large enough (otherwise estimate is too noisy); drawn from p

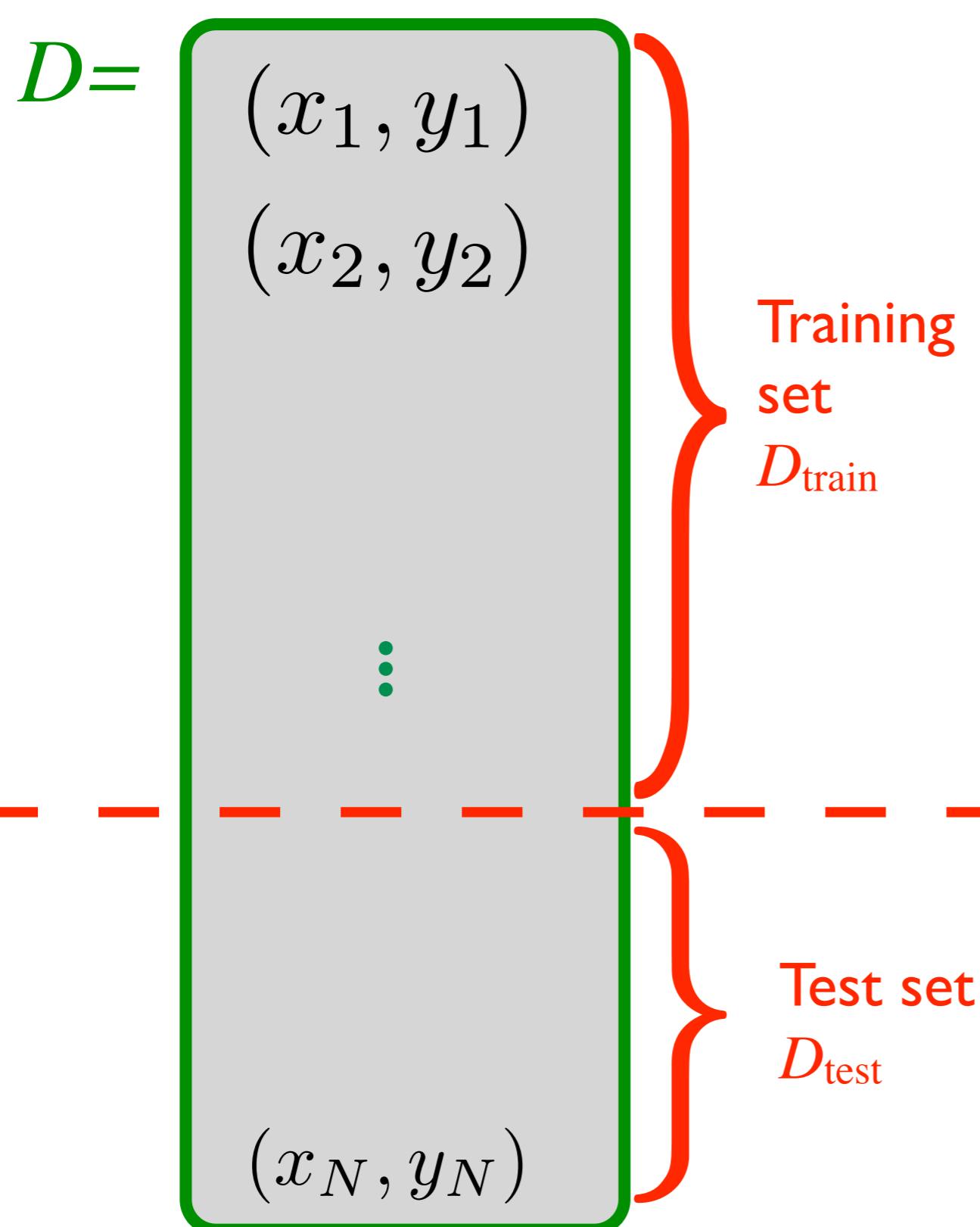
→ Must keep a separate test-set $D_{\text{test}} \neq D_{\text{train}}$ to properly estimate generalization error of $\hat{f}(D_{\text{train}})$:

$$R(\hat{f}(D_{\text{train}})) \approx \hat{R}(\hat{f}(D_{\text{train}}), D_{\text{test}})$$

generalization average error on
error test-set (never used for training)

This is the test phase in ML

Simple train/test procedure



- Provided large enough dataset D drawn from $p(x,y)$
- Make sure examples are in random order.
- Split dataset in **two**: D_{train} and D_{test}
- Use D_{train} to choose/ optimize/find best predictor $f = \hat{f}(D_{\text{train}})$
- Use D_{test} to evaluate generalization performance of predictor f .

Model selection

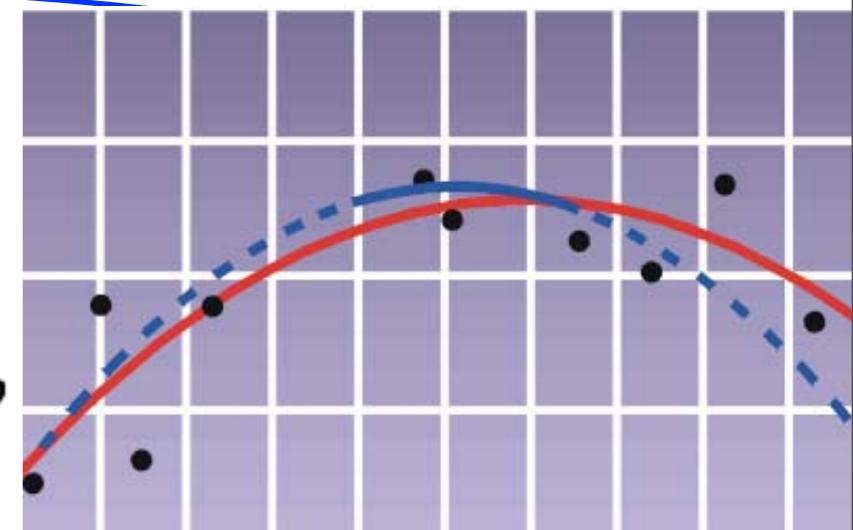
Choosing a specific
function family F

Ex. of parameterized function families

$F_{\text{polynomial } p}$

Polynomial predictor (of degree p):

$$f(x) = b + a_1x + a_2x^2 + a_3x^3 + \dots + a_p x^p$$



F_{linear}

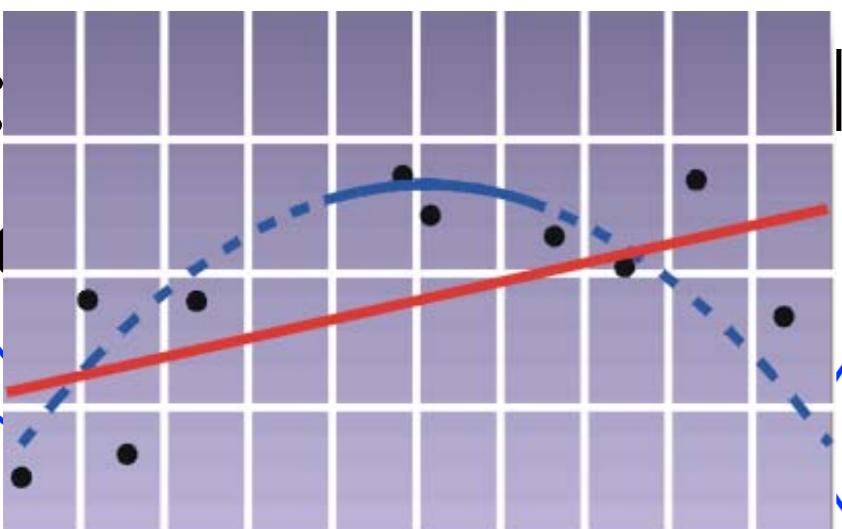
Linear (affine) predictor:
«linear regression»

$$\theta = \{w \in \mathbb{R}^d, b \in \mathbb{R}\}$$

$$f_\theta(x) = wx + b \quad (\text{in 1 dimension})$$

$$f_\theta(x) = w^T x + b \quad (\text{in } d \text{ dimensions})$$

Q:
pre
test



F_{const}

Constant predictor: $f_\theta(x) = b$
where $\theta = \{b\}$
(always predict the same value or class!)

Capacity of a learning algorithm

- Choosing a specific Machine Learning algorithm means choosing a specific function family F .
- How «big, rich, flexible, expressive, complex» that family is, defines what is informally called the «capacity» of the ML algorithm.
Ex: $\text{capacity}(F_{\text{polynomial } 3}) > \text{capacity}(F_{\text{linear}})$
- One can come up with several formal measures of «capacity» for a function family / learning algorithm (e.g. **VC-dimension** Vapnik–Chervonenkis)
- One rule-of-thumb estimate, is the **number of adaptable parameters**: i.e. how many scalar values are contained in θ .
Notable exception: chaining many linear mappings is still a linear mapping!

Effective capacity, and capacity-control hyper-parameters

The «effective» capacity of a ML algo is controlled by:

- Choice of ML algo, which determines big family F
- Hyper-parameters that further specify F
e.g.: degree p of a polynomial predictor; Kernel choice in SVMs;
#of layers and neurons in a neural network
- Hyper-parameters of «regularization» schemes
e.g. constraint on the norm of the weights w
(\Rightarrow ridge-regression; L_2 weight decay in neural nets);
Bayesian prior on parameters; noise injection (dropout); ...
- Hyper-parameters that control early-stopping of the iterative search/optimization procedure.
(\Rightarrow won't explore as far from the initial starting point)

Popular classifiers

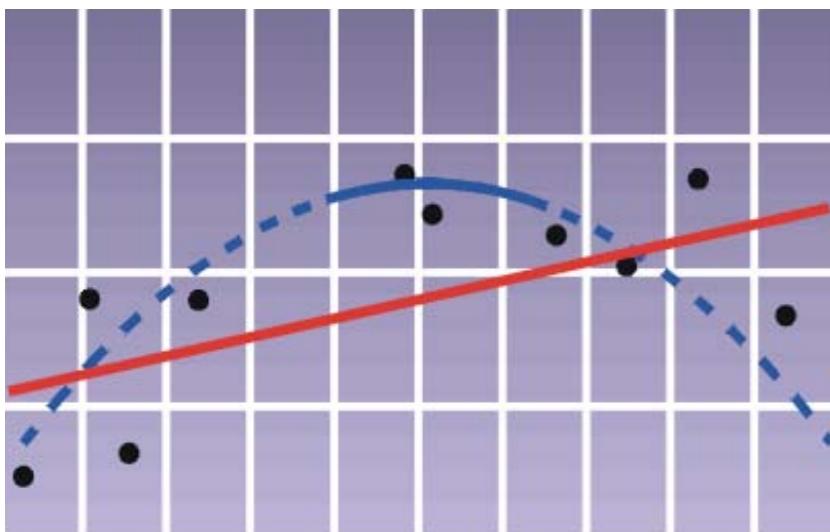
their parameters and hyper-parameters

Algo	Capacity-control hyperparameters	Learned parameters
logistic regression (L ₂ regularized)	strength of L ₂ regularizer	w,b
linear SVM	C	w,b
kernel SVM	C; kernel choice & params (σ for RBF; degree for polynomial)	support vector weights: α
neural network	layer sizes; early stop; ...	layer weight matrices
decision tree	depth	the tree (with index and threshold of variables)
k-nearest neighbors	k; choice of metric	memorizes trainset

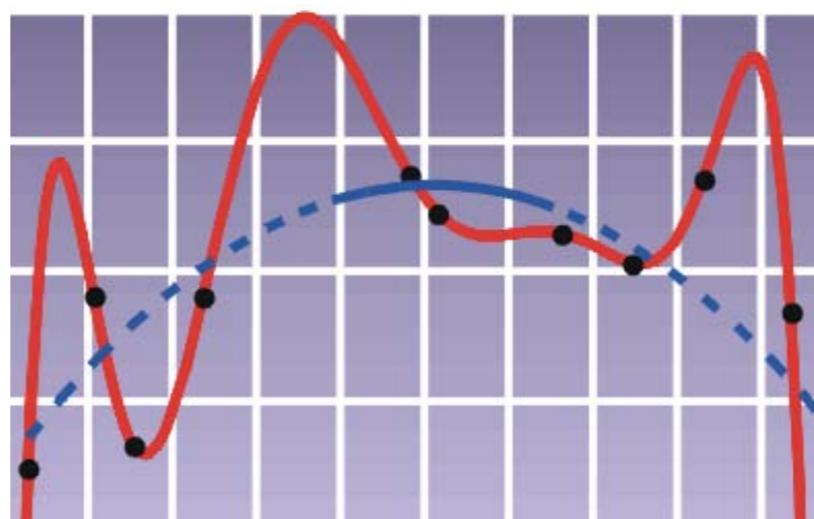
Tuning the capacity

- Capacity must be optimally tuned to ensure good generalization
- by choosing Algorithm and hyperparameters
- to avoid under-fitting and over-fitting.

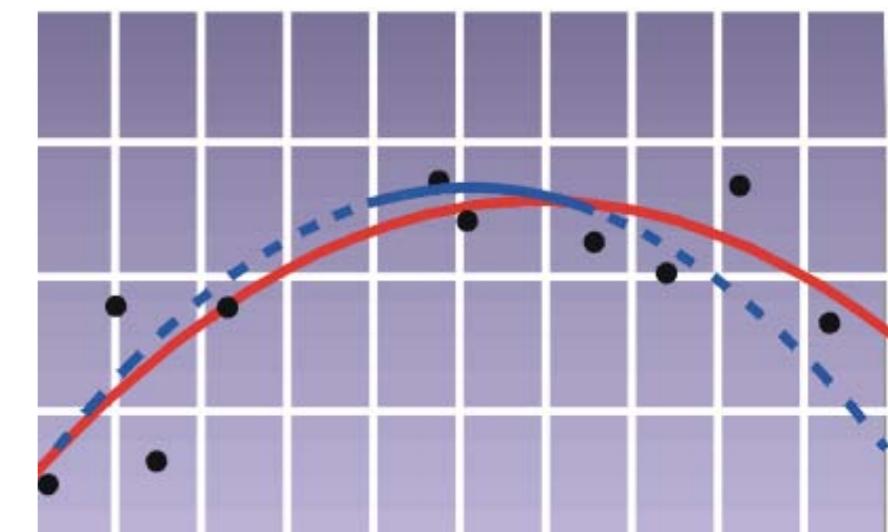
Ex: 1D regression with polynomial predictor



capacity too low
→under-fitting



capacity too high
→over-fitting



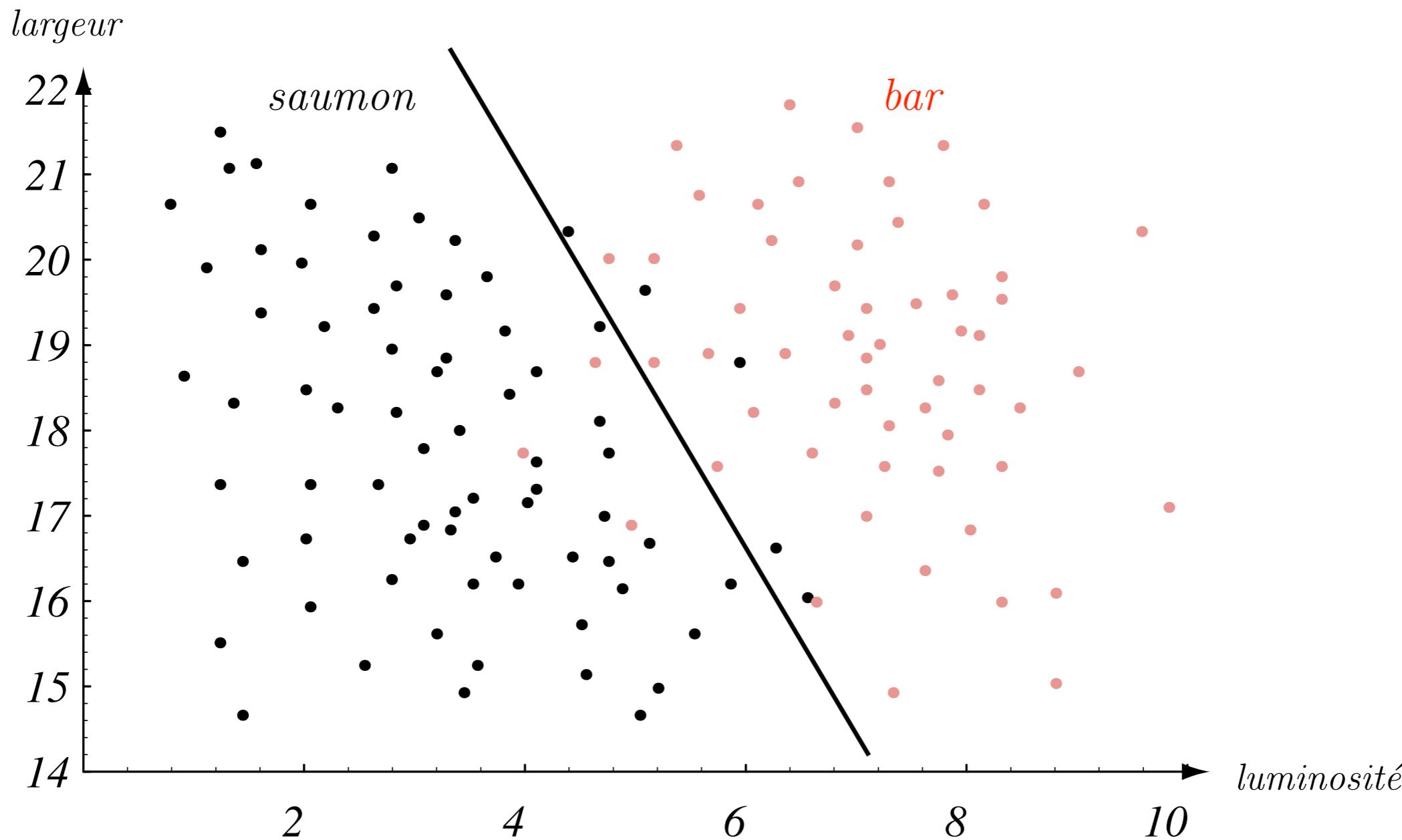
optimal capacity
→good generalisation

performance on training set is not a good estimate of generalization

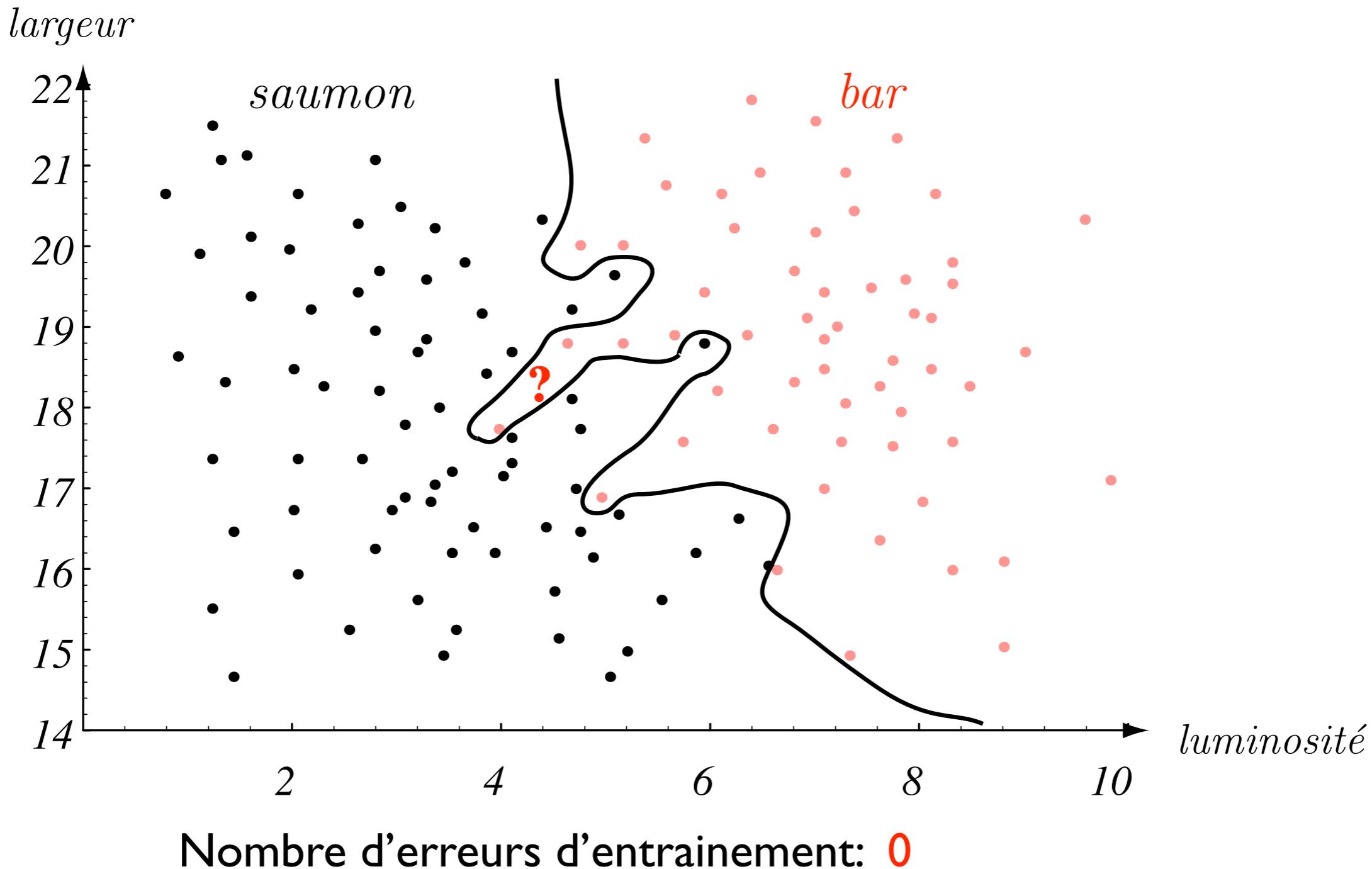
Ex: 2D classification

Linear classifier

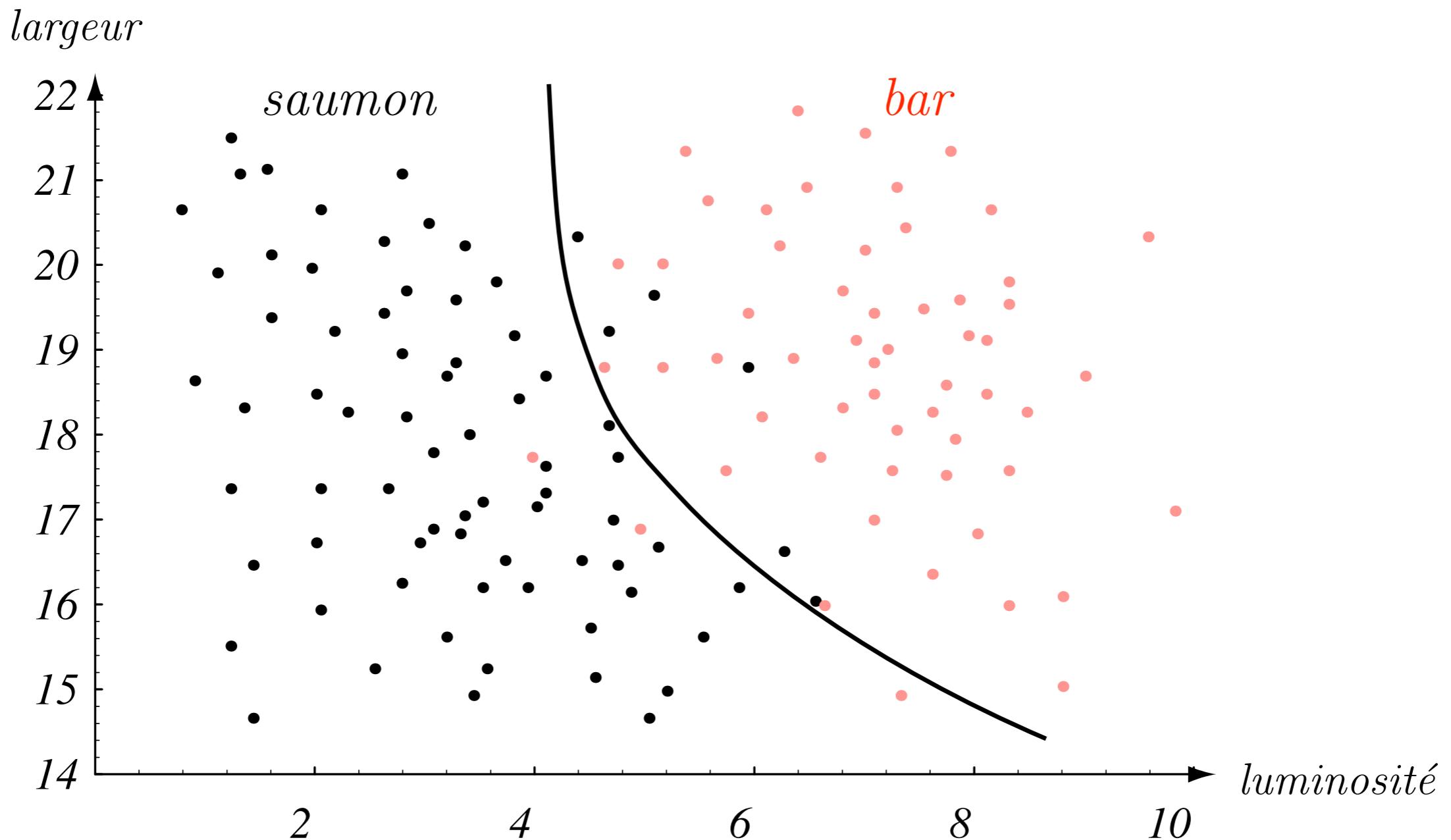
- Function family too poor
(too inflexible)
- = **Capacity too low** for this problem
(relative to number of examples)
- => Under-fitting



- Function family too rich
(too flexible)
- = **Capacity too high** for this problem
(relative to the number of examples)
- => Over-fitting

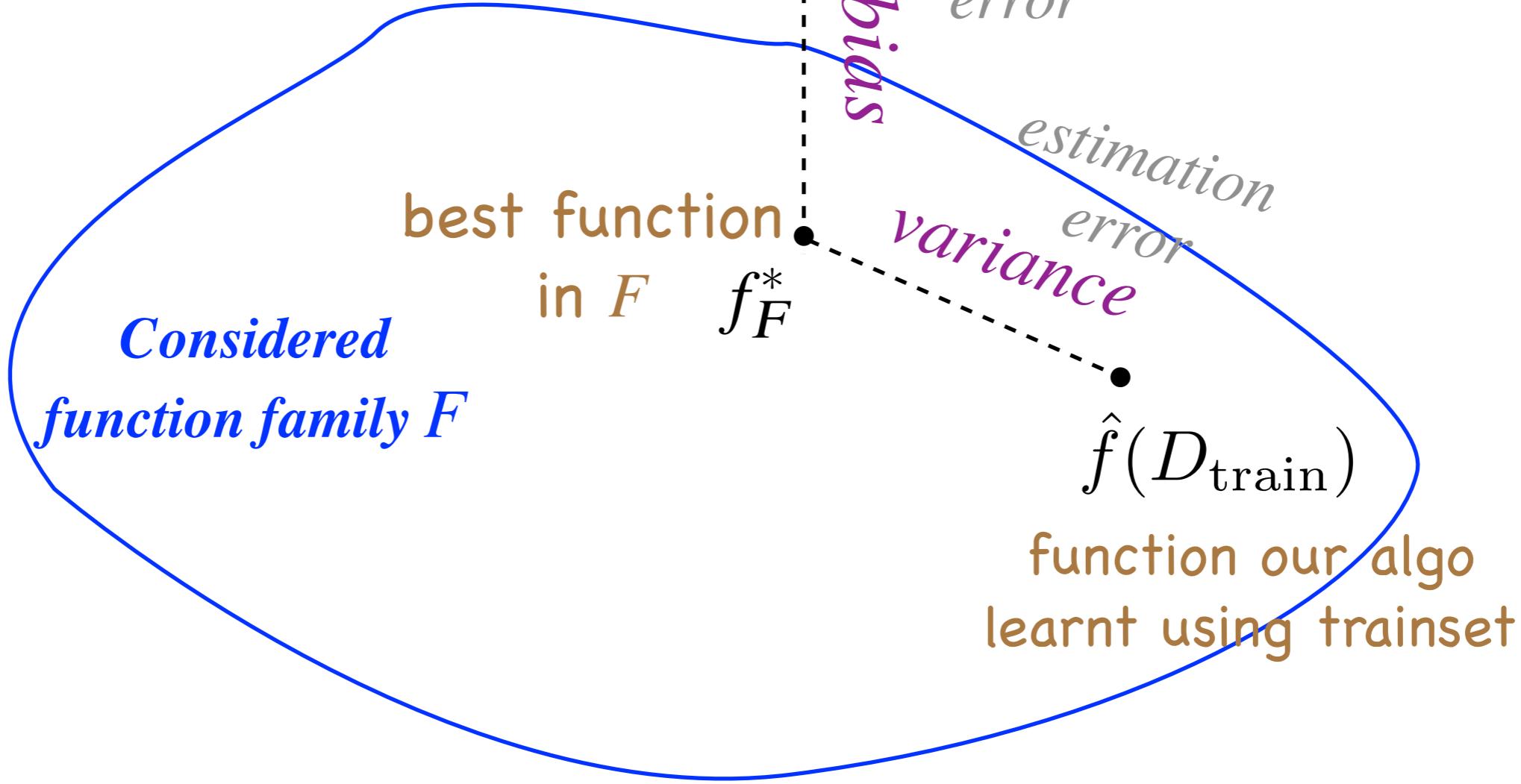


- **Optimal capacity** for this problem
(par rapport à la quantité de données)
- => Best generalization
(on future test points)



Decomposing the generalization error

*Set of all possible
functions
in the universe*

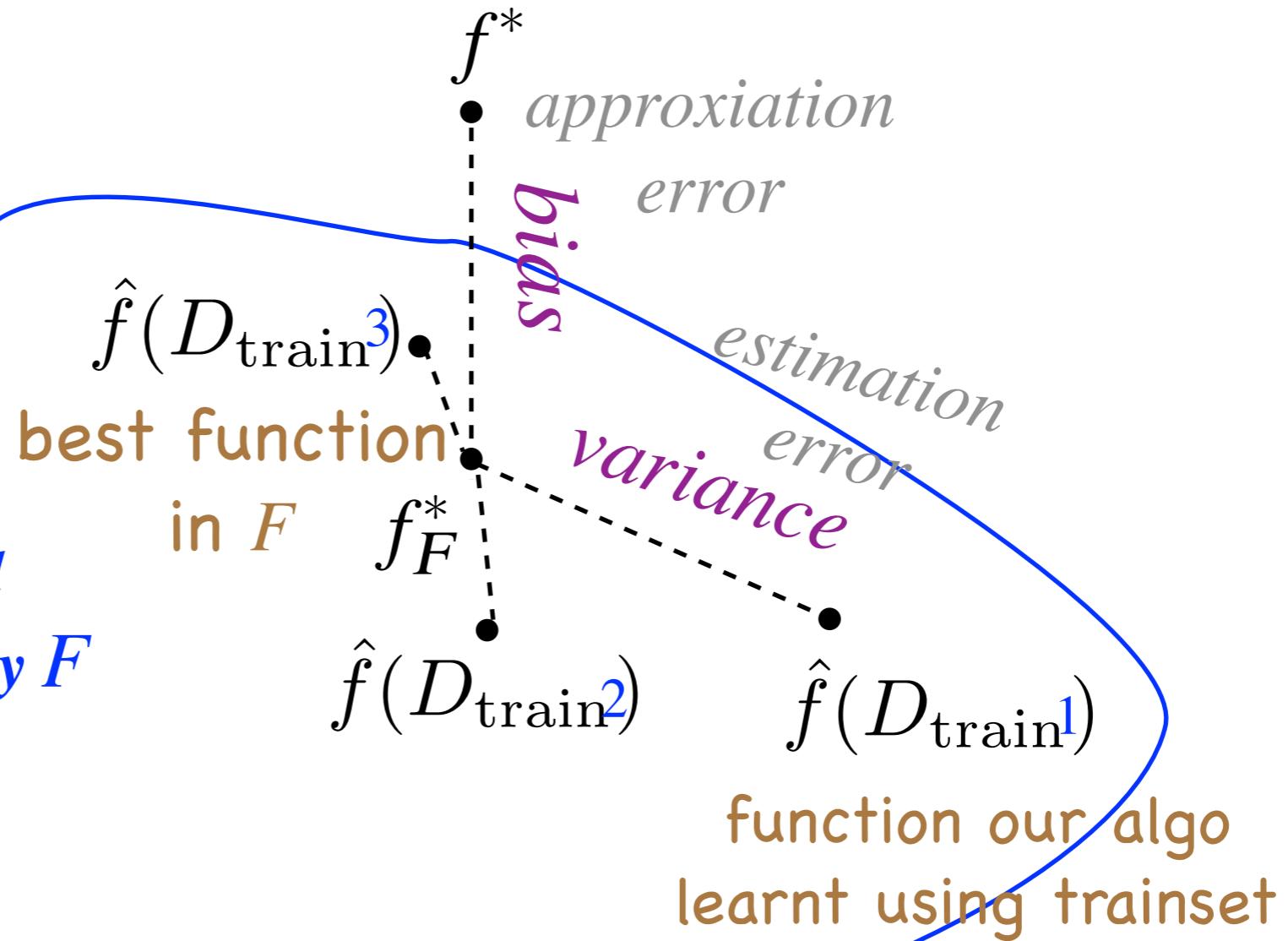


What is responsible for the variance?

*Set of all possible
functions
in the universe*

*Considered
function family F*

best possible
function



Optimal capacity & the bias-variance dilemma

- Choosing richer F : capacity \uparrow
 ⇒ bias \downarrow but variance \uparrow .
- Choosing smaller F : capacity \downarrow
 ⇒ variance \downarrow but bias \uparrow .
- Optimal compromise... will depend on number of examples n
- Bigger n ⇒ variance \downarrow
 So we can afford to increase capacity (to lower the bias)
 ⇒ can use more expressive models
- The best regularizer is more data!

Model selection how to

$D =$

(x_1, y_1)

(x_2, y_2)

— — —

⋮

— — —

(x_N, y_N)

Training set
 D_{train}

Validation set
 D_{valid}

Test set
 D_{test}

Make sure examples are in random order

Split data D in 3: D_{train} D_{valid} D_{test}

Model selection meta-algorithm:

For each considered model (ML algo) A :

For each considered hyper-parameter config λ :

- train model A with hyperparams λ on D_{train}

$$\hat{f}_{A_\lambda} = A_\lambda(D_{\text{train}})$$

- evaluate resulting predictor on D_{valid}
(with preferred evaluation metric)

$$e_{A_\lambda} = \hat{R}(\hat{f}_{A_\lambda}, D_{\text{valid}})$$

Locate A^*, λ^* that yielded best e_{A_λ}

Either return $f^* = f_{A_{\lambda^*}}$

Or retrain and return

$$f^* = A_{\lambda^*}^*(D_{\text{train}} \cup D_{\text{valid}})$$

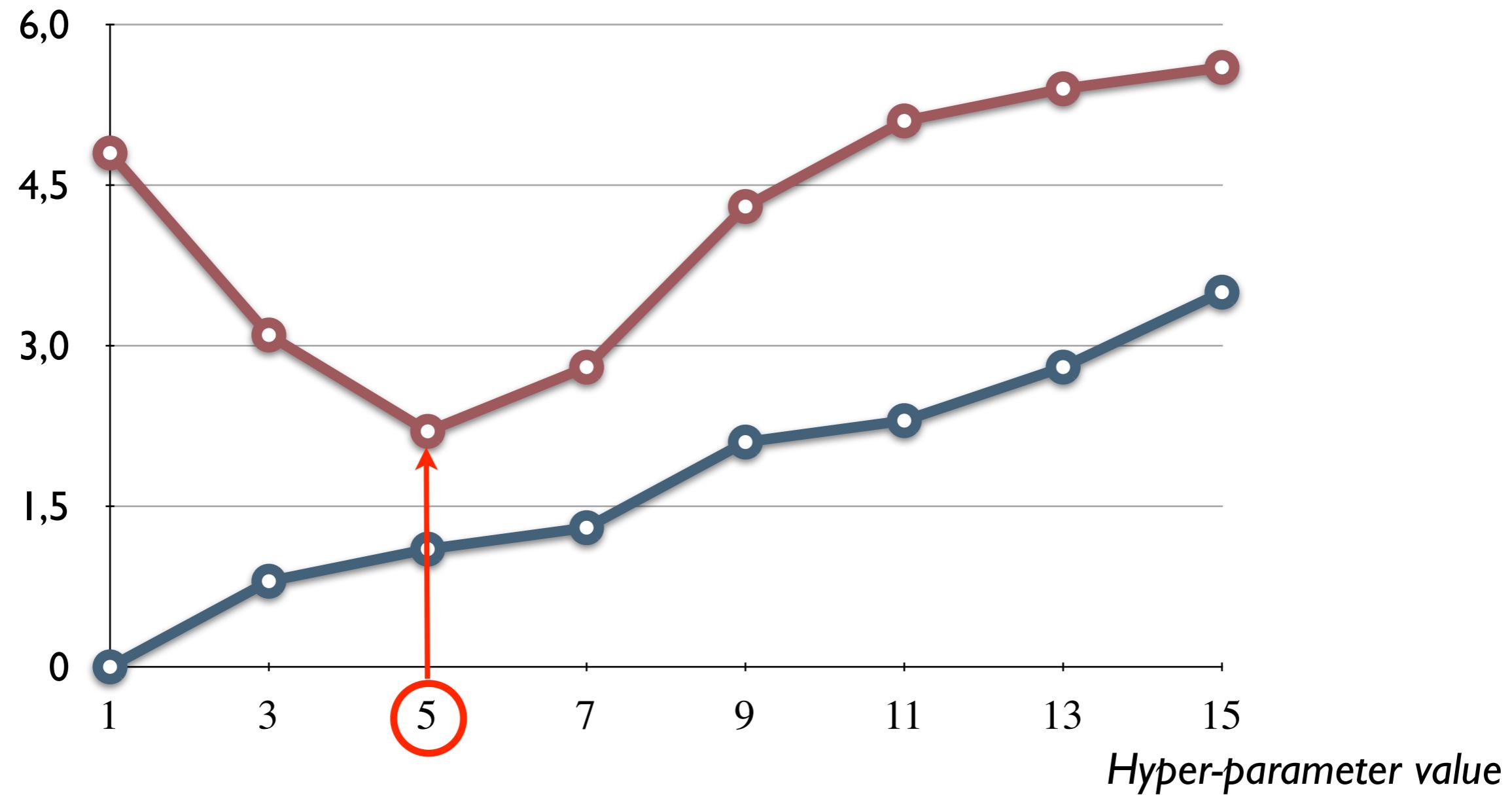
Finally: compute unbiased estimate of generalization performance of f^* using D_{test}

$$\hat{R}(f^*, D_{\text{test}})$$

D_{test} must never have been used during training or model selection to select, learn, or tune anything.

Ex of model hyper-parameter selection

- Training set error
- Validation set error



Hyper-parameter value which yields smallest error on validation set is 5
(it was 1 for the training set)

Question

What if we selected capacity-control hyper-parameters that yield best performance on the training set?

What would we tend to select?

Is it a good idea? Why?

Model selection procedure summary:

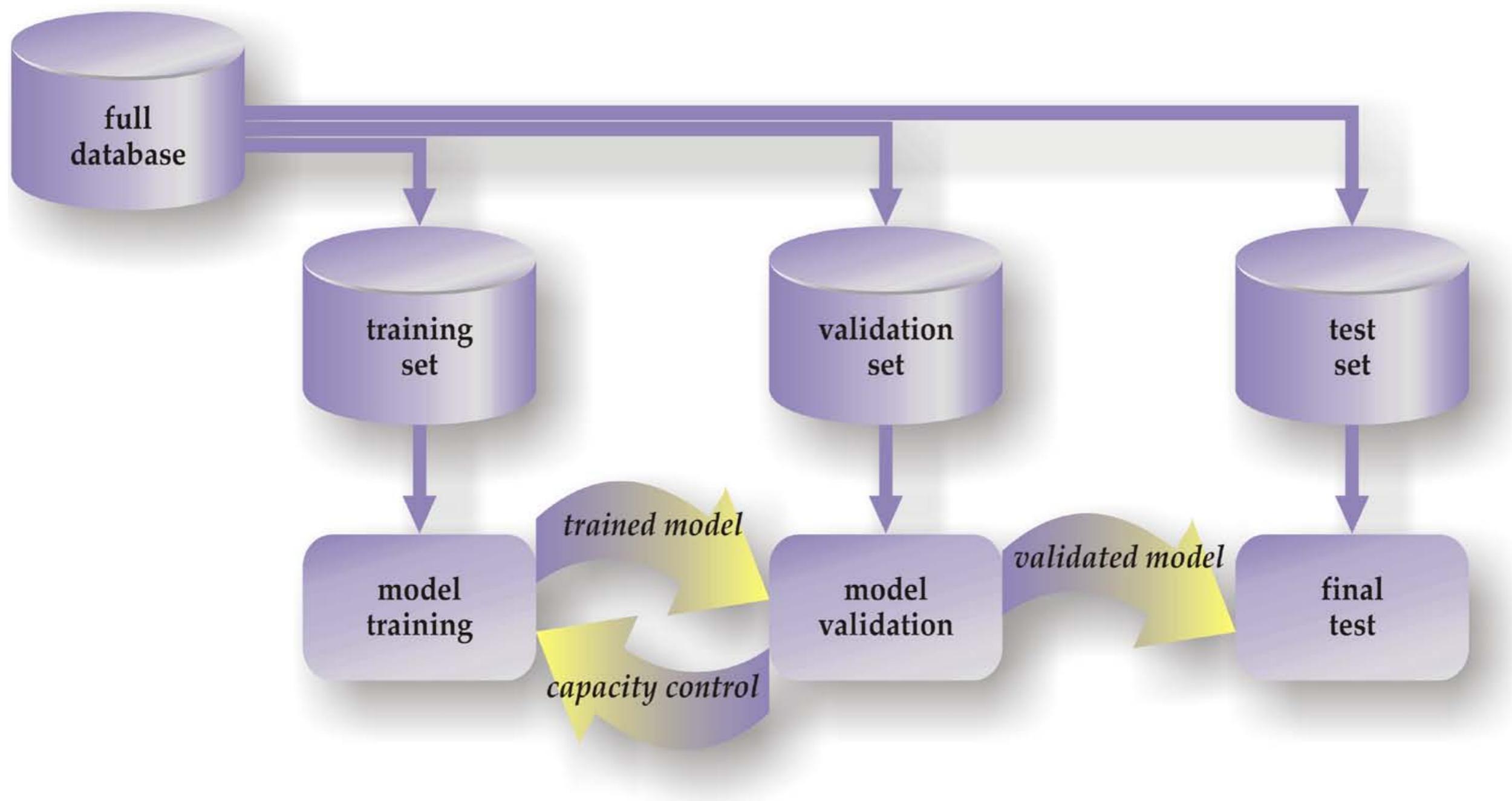
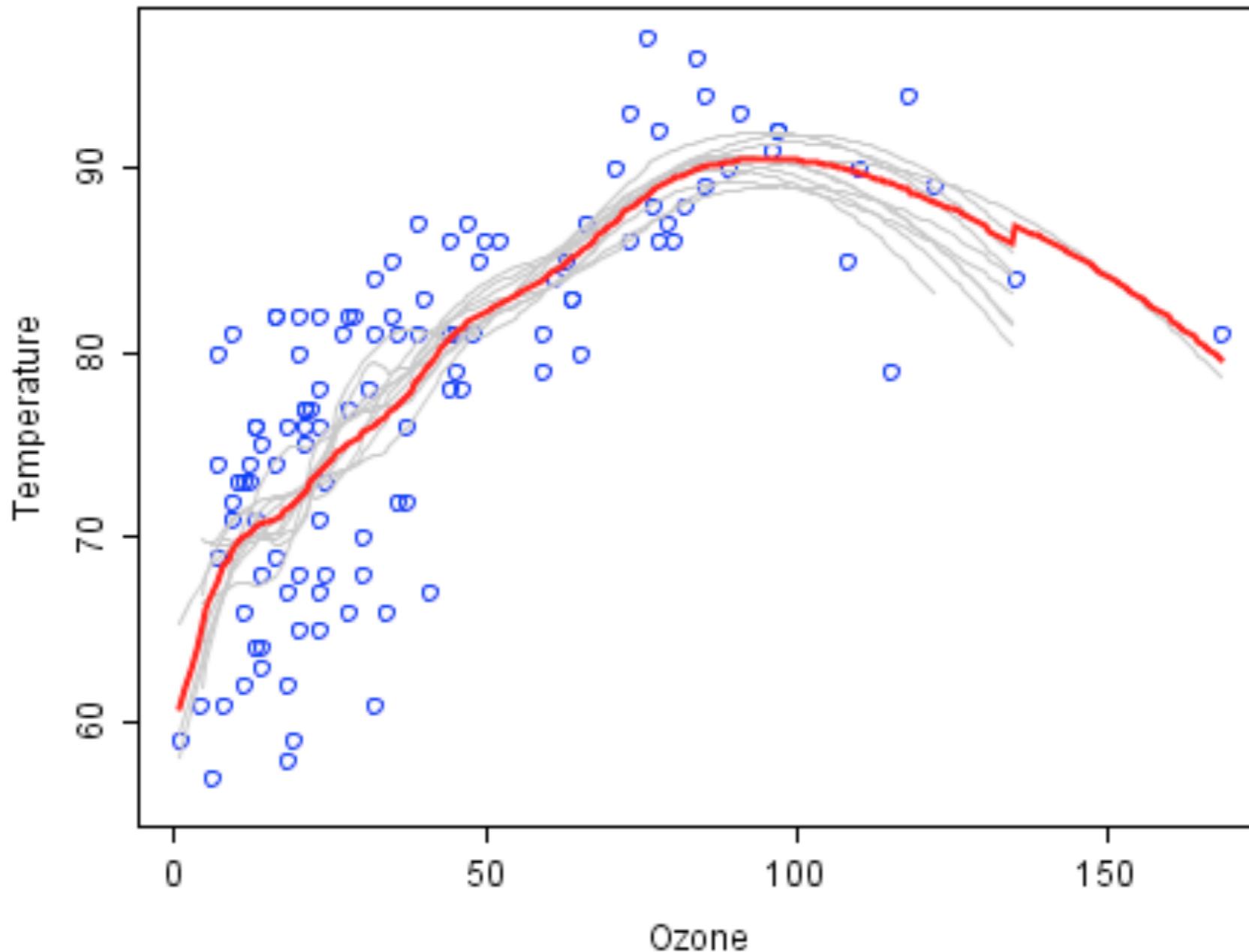


Figure by Nicolas Chapados

Ensemble methods

- Principle: train and combine multiple predictors to good effect
- Bagging: average many high-variance predictors
 - variance ↓
(e.g.: average deep trees → Random decision forests)
- Boosting: build weighted combination of low-capacity classifiers
 - bias ↓ and capacity ↑
(e.g. boosting shallow trees; or linear classifiers)

Bagging for reducing variance on a regression problem



How to obtain non-linear predictor with a linear predictor

Three ways to map x to a feature representation $\tilde{x} = \phi(x)$

- Use an **explicit fixed mapping** (ex: hand-crafted features)
- Use an **implicit fixed mapping**
 - ➡ Kernel Methods (SVMs, Kernel Logistic Regression ...)
- **Learn a parameterized mapping**
 - (i.e. let the ML algo learn the new representation)
 - ➡ **Multilayer feed-forward Neural Networks**
such as **Multilayer Perceptrons (MLP)**

Levels of representation



very high level representation:

CAT

JUMPING

... etc ...

slightly higher level representation

raw input vector representation:

$$\mathcal{X} = \boxed{23 \quad 19 \quad 20} \quad \cdots \quad \boxed{\quad \quad 18}$$

x_1 x_2 x_3 \cdots x_n





Questions ?