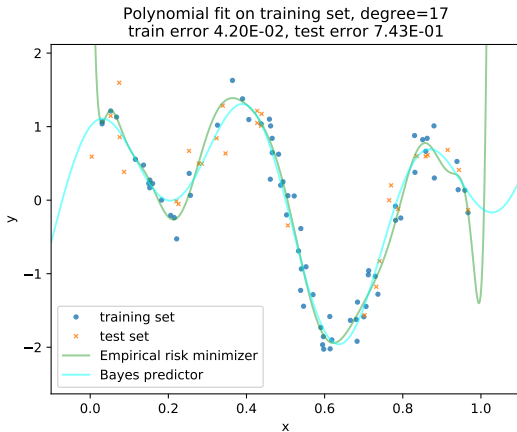


Machine learning I, supervised learning: risks



Notion of risks

- ▶ We are ready to introduce an important notion that is specific to machine learning and optimization : the risk
- ▶ there are several types of risks and several denominations for each.
- ▶ this denomination "risk" might seem counter-intuitive at first, as there is no notion of danger involved. However, this is a classical term in optimization and ML.

Setting

We consider

- ▶ an input space \mathcal{X} (e.g. $\mathcal{X} = \mathbb{R}^d$)
- ▶ an output space \mathcal{Y} .

In supervised learning, we predict outputs $y \in \mathcal{Y}$ from inputs $x \in \mathcal{X}$.

- ▶ classification : discrete \mathcal{Y} , e.g. $\mathcal{Y} = \{0, 1\}$, $\mathcal{Y} = \{-1, 1\}$, $\mathcal{Y} = \{0, 1, 2\}$.
- ▶ regression : continuous \mathcal{Y} , e.g. $\mathcal{Y} = \mathbb{R}$, $\mathcal{Y} = [a, b]$.

The couples (x, y) are called **samples** and are considered to be sampled from a joint random variable (X, Y) .

Supervised learning

- ▶ Assumption : there exists a joint probability law ρ , such that $(X, Y) \sim \rho$. However, ρ is **unknown**.
- ▶ Hence there exists a map $f : \mathcal{X} \mapsto \mathcal{Y}$, such that $Y = f(X)$.
- ▶ f is most of the time non deterministic.

Supervised Learning : from a finite dataset of samples, produce an estimate \tilde{f} of f .

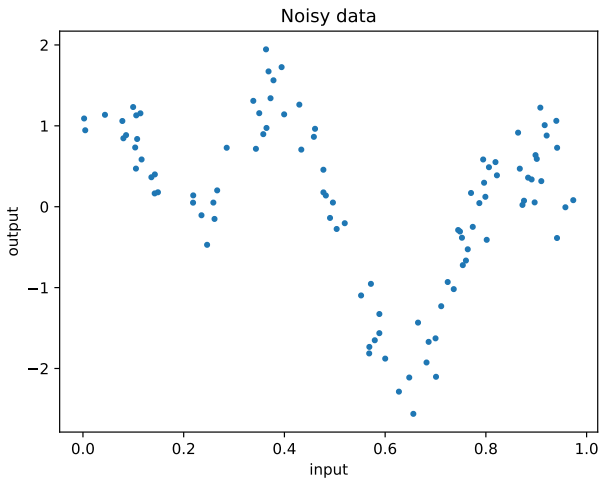


Figure – Finite dataset in 1 dimension

Loss functions

A **loss function** l is a map that measures the discrepancy between to elements of a set.

$$l : \begin{cases} \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+ \\ (y, y') \mapsto l(y, y') \end{cases}$$

We use it in order to evaluate the quality of our prediction $\tilde{f}(x)$, that should be close to the label y that corresponds to x .

Common loss functions

- ▶ "0-1" loss (for **classification**.)

$$l(y, z) = 1_{y \neq z} \quad (1)$$

- ▶ squared loss (for **regression**)

- ▶ $\mathcal{Y} = \mathbb{R}$.

$$l(y, z) = (y - z)^2 \quad (2)$$

- ▶ $\mathcal{Y} = \mathbb{R}^d$

$$l(y, z) = \|y - z\|_2^2 \quad (3)$$

- ▶ absolute loss (for **regression**). $\mathcal{Y} = \mathbb{R}$.

- ▶ $\mathcal{Y} = \mathbb{R}$.

$$l(y, z) = |y - z| \quad (4)$$

- ▶ $\mathcal{Y} = \mathbb{R}^d$

$$l(y, z) = \|y - z\|_1 \quad (5)$$

Prerequisite : expected value

Let Z be a real random variable. If it is correctly defined, the expected value is

- ▶ for a discrete random variable (that can take the values $\{z_i, i \in \mathbb{N}\}$).

$$E[Z] = \sum_{i=1}^{+\infty} z_i P(Z = z_i) \quad (6)$$

- ▶ for a continuous random variable

$$E[Z] = \int_{-\infty}^{+\infty} zp(z)dz \quad (7)$$

$p(z)$ is the density of probability of Z , assumed to exist.

Expected values

Expected value of an unbiased dice game :

$$E[Z] = \frac{1}{6}[1 + 2 + 3 + 4 + 5 + 6] = 3.5 \quad (8)$$

Expected value of a cheated dice game :

$$E[Z] = \frac{1}{100}(1 + 2 + 3 + 4) + \frac{48}{100}(5 + 6) = 5.38 \quad (9)$$

Risks

- ▶ We call "estimator" a map $\mathcal{X} \mapsto \mathcal{Y}$
- ▶ We note $D_n = \{(x_1, y_1), \dots, (x_n, y_n), i \in [1, n]\}$ the dataset.
From D_n , we want to estimate f .

To measure the quality of some estimator g , we consider the **risks** :

- ▶ Risk / generalization error ("risque réel" in french)

$$R(g) = E_{(X,Y) \sim \rho}[l(Y, g(X))] \quad (10)$$

- ▶ Empirical risk ("risque empirique" in french)

$$R_n(g) = \frac{1}{n} \sum_{i=1}^n l(y_i, g(x_i)) \quad (11)$$

Both risks depend on the loss function l !

Risks

Risk / generalization error :

$$R(g) = E_{(X,Y) \sim \rho}[l(Y, g(X))] \quad (12)$$

Problem : we cannot compute $R(g)$!

Risks

Risk / generalization error :

$$R(g) = E_{(X,Y) \sim \rho}[l(Y, g(X))] \quad (13)$$

Problem : we cannot compute $R(g)$!

We **only** have access to the empirical risk.

$$R_n(g) = \frac{1}{n} \sum_{i=1}^n l(y_i, g(x_i)) \quad (14)$$

given the finite dataset $D_n = \{(x_1, y_1), \dots, (x_n, y_n), i \in [1, n]\}$.

Optimization problem : empirical risk minimization

- ▶ The smaller the generalization error $R(g)$ is, the better g is.
- ▶ The situation is more tricky for $R_n(g)$: it is not obvious that as estimator that has a very small empirical risk $R_n(g)$ has a small generalization error $R(g)$! This is the problem of **overfitting**.

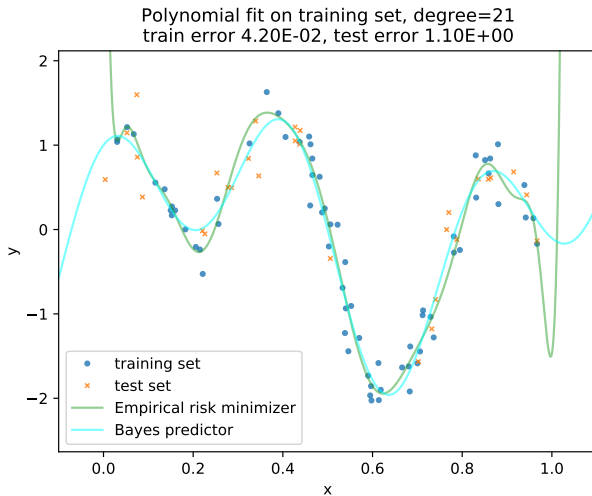


Figure – Overfitting : the green estimator has a small empirical risk, but it a large generalization.

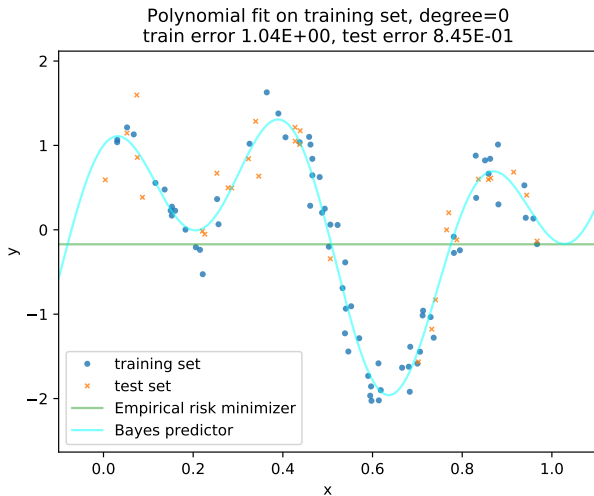


Figure – Very simple estimator

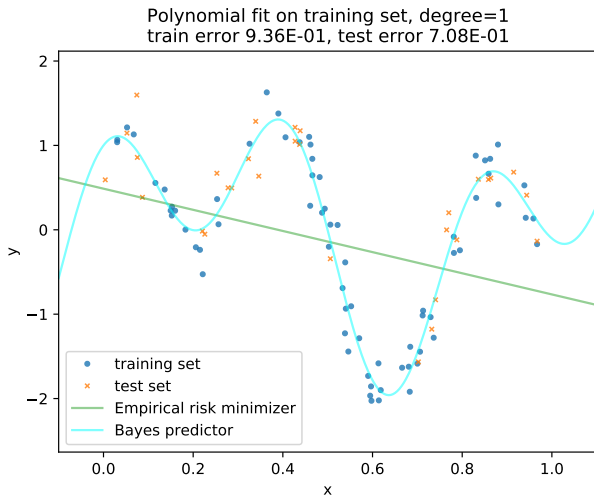
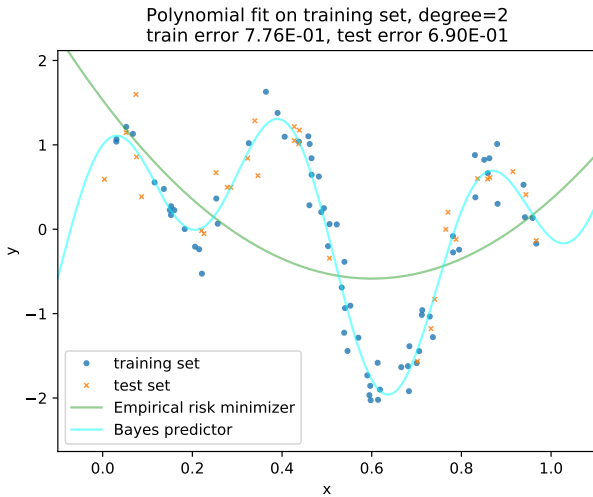
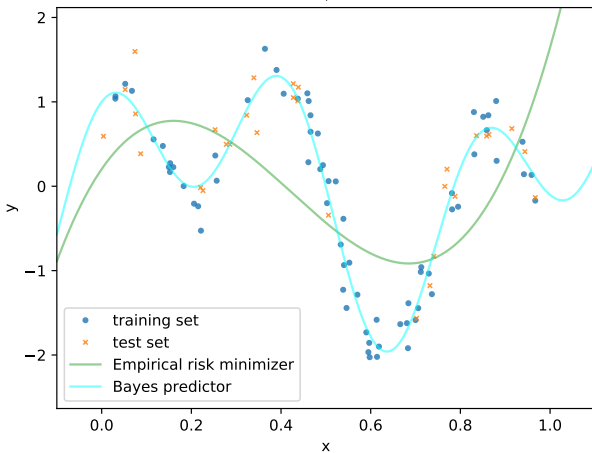


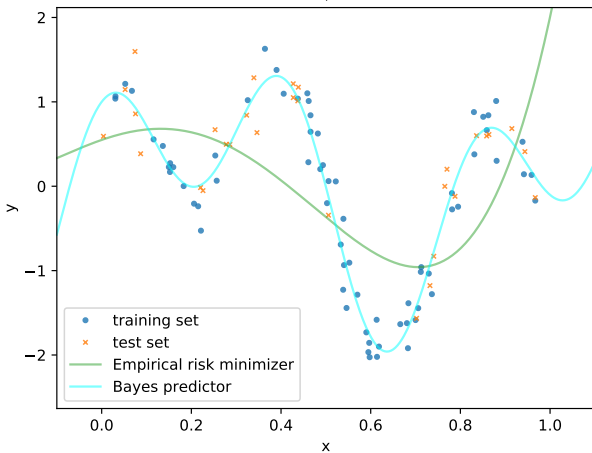
Figure – Very simple estimator



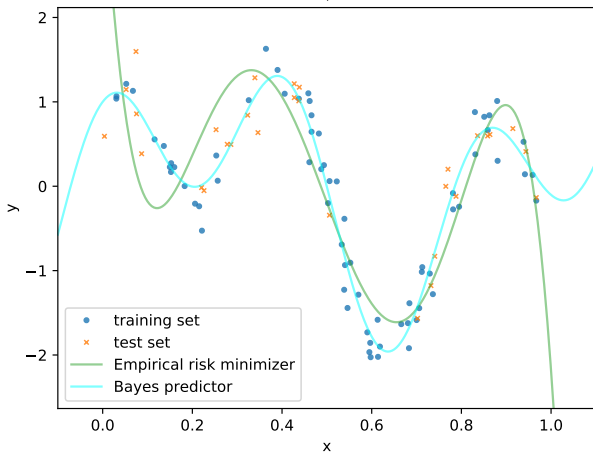
Polynomial fit on training set, degree=3
train error 6.46E-01, test error 5.51E-01



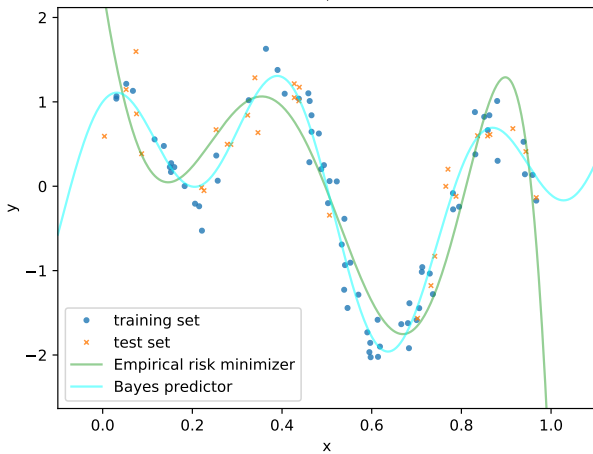
Polynomial fit on training set, degree=4
train error 6.40E-01, test error 5.64E-01



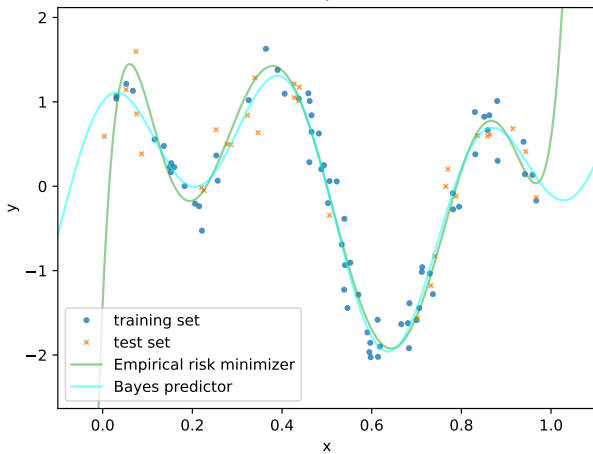
Polynomial fit on training set, degree=5
train error 1.99E-01, test error 5.52E-01



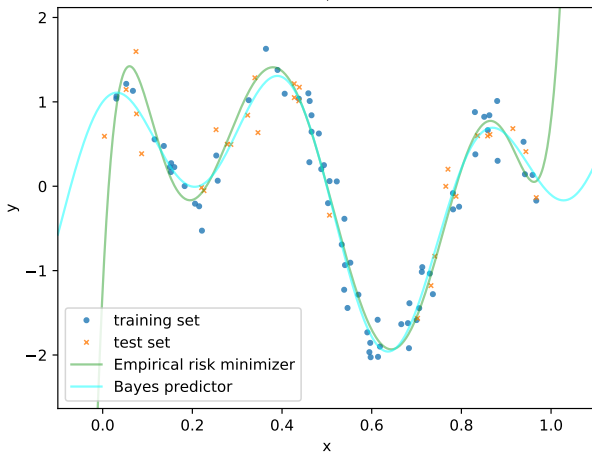
Polynomial fit on training set, degree=6
train error 1.62E-01, test error 2.83E-01



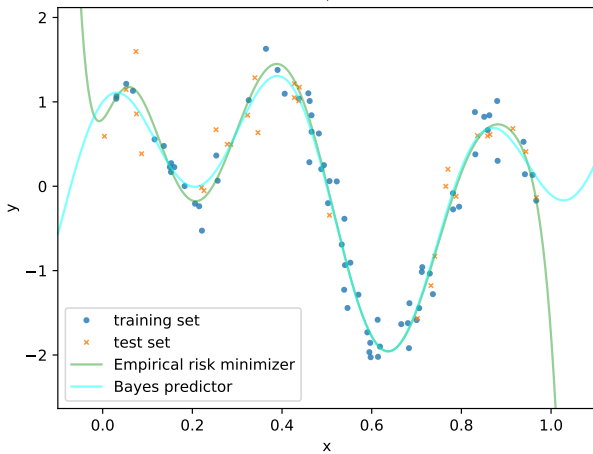
Polynomial fit on training set, degree=7
train error $5.71\text{E-}02$, test error $1.65\text{E-}01$

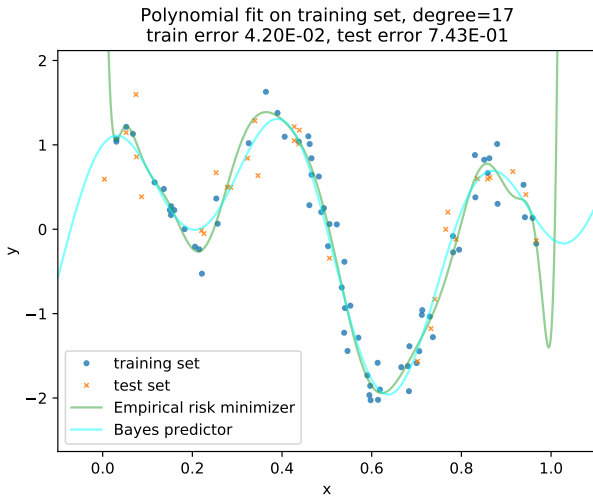


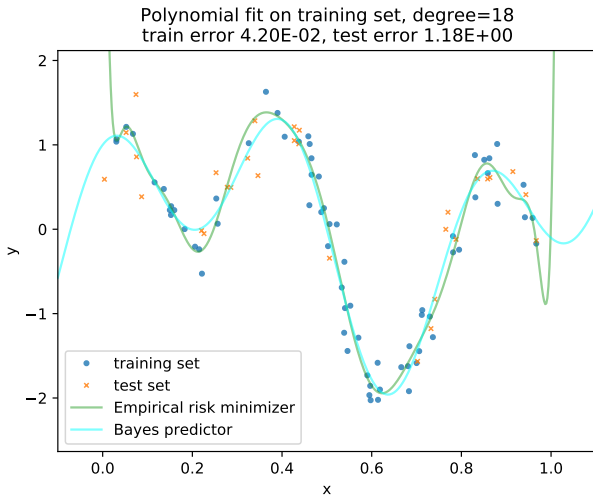
Polynomial fit on training set, degree=8
train error 5.70E-02, test error 1.50E-01

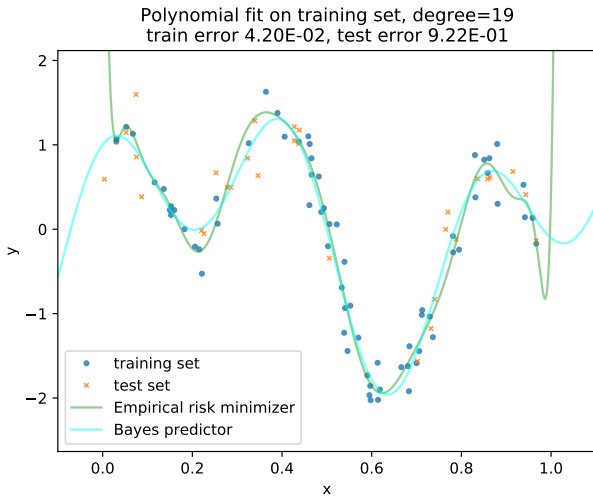


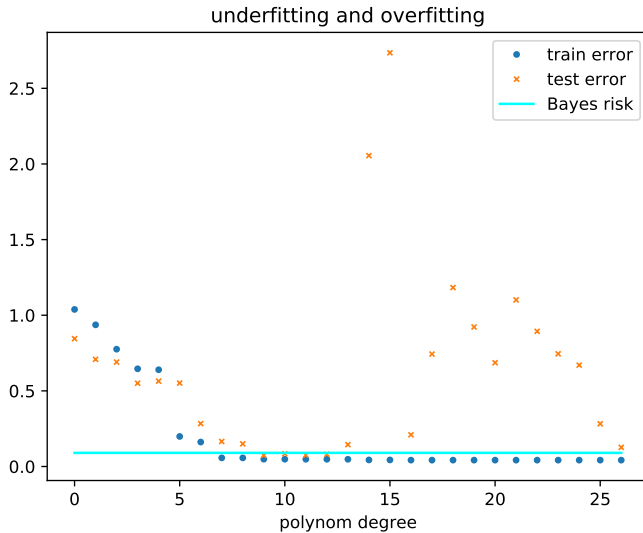
Polynomial fit on training set, degree=9
train error 4.84E-02, test error 7.33E-02











Randomness

If the data were deterministic ($Y = f(X)$ is deterministic), there would be no overfitting!

Randomness might come from several sources, such as :

- ▶ measurement errors
- ▶ hidden variables (not represented in X)

Optimization problem : empirical risk minimization

Empirical risk minimization (ERM) : finding the estimator f_n that minimizes the empirical risk R_n .

This raises important questions :

- ▶ 1) does f_n have a good generalization error $R(f_n)$?
- ▶ 2) how can we have guarantees on the generalization error $R(f_n)$?
- ▶ 3) how can we find the empirical risk minimizer f_n ?
- ▶ 4) is it even interesting to strictly minimize R_n ?

Generalization error

Question 1) Does f_n have a good generalization error $R(f_n)$?

This will depend on :

- ▶ the number of samples n
- ▶ the shape of f (the map such that $Y = f(X)$), in particular on its **regularity**
- ▶ the distribution ρ
- ▶ the dimensions of the input space and of the output space.
- ▶ the space of functions where f_n is taken from.

Statistical bounds

Question 2) How can we have guarantees on the generalization error $R(f_n)$?

By making **assumptions** on the problem (learning is impossible without making assumptions), for instance assumptions on ρ .

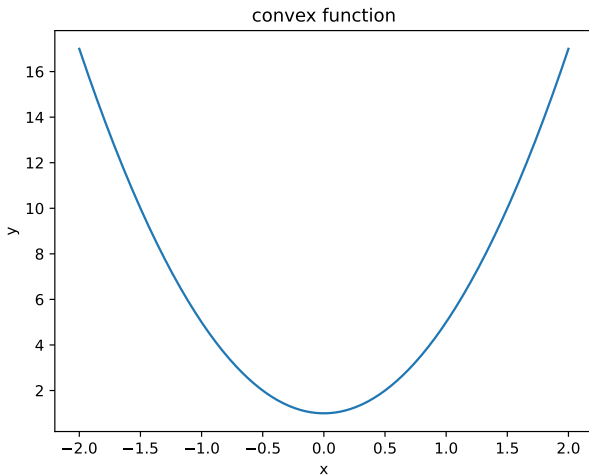
Optimization

Question 3) how can we find the empirical risk minimizer f_n ?

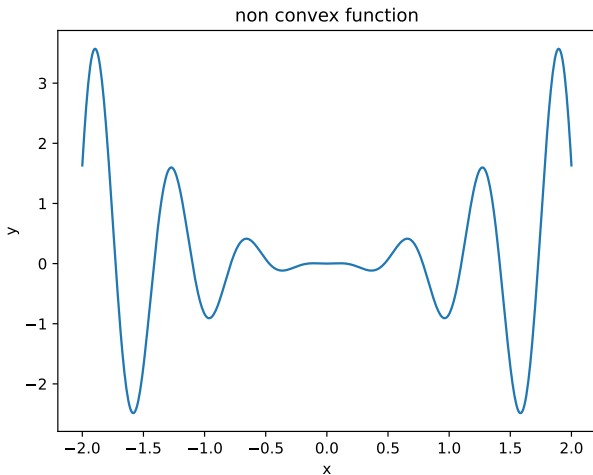
By using an optimization algorithm or by solving the minimization in closed-form.

Convex functions

Convex functions are easier to minimize.



Non convex functions



What is convex here?

In this context, the convexity that is involved is the dependence of R_n in g . More precisely, for instance if g depends on $\theta \in \mathbb{R}^d$, e.g. $g(x) = \langle \theta, x \rangle$, the convexity is that of

$$\theta \mapsto R_n(\theta) \tag{15}$$

Example (ordinary least squares) :

$$R_n(\theta) = \frac{1}{n} \sum_{i=1}^n (\langle \theta, x_i \rangle - y_i)^2 \tag{16}$$

with $x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$.

Optimization error

Question 4) is it even interesting to strictly minimize R_n ?

Most of the time it is **not**, as we are interested in R , not in R_n , so we should not try to go to machine precision in the minimization of a quantity that is itself an approximation !

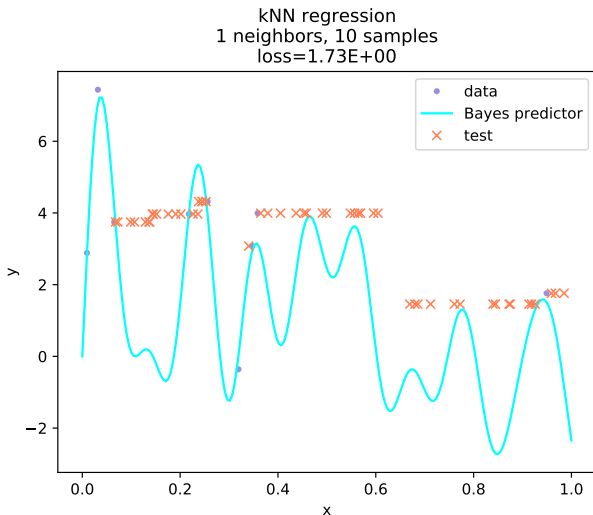
This is linked to the **estimation error** (advanced concept) that is often of order $\mathcal{O}(1/\sqrt{n})$.

Nearest neighbors algorithms

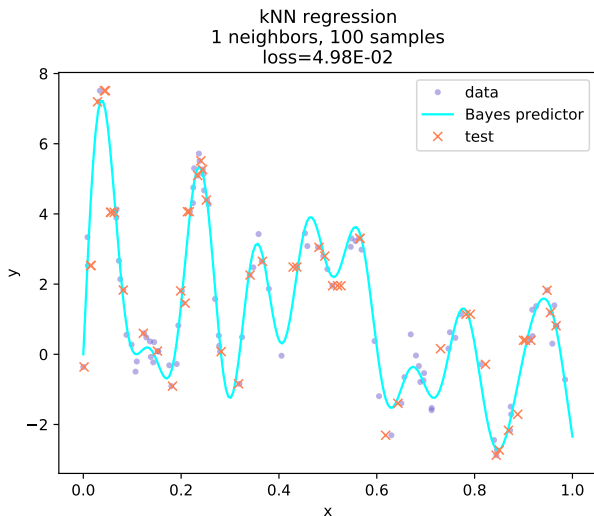
Not all supervised learning methods consist in Empirical risk minimization (ERM).

For instance the nearest neighbors algorithm is not an ERM.

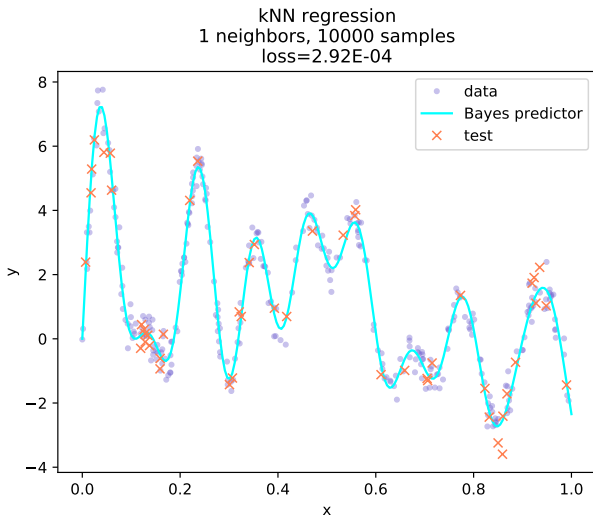
kNN algorithm



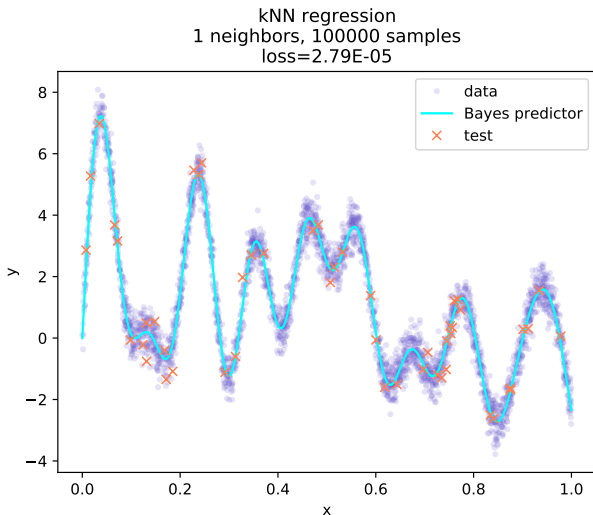
kNN algorithm



kNN algorithm



kNN algorithm



Bayes risk

Notion of bayes risk on the penalty shootout example.