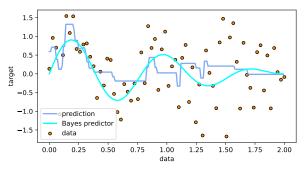
Fondamentaux théoriques du machine learning

Random forest regression number of estimators: 5 max depth: 3 test error: 8.14E-01 Bayes risk: 7.00E-01



Ensemble learning

Bagging Random forest Boosting

Statistical learning

Bounding the estimation error Interpolation regime and double descent

Local averaging methods

Supervised learning Density estimation

Ensemble learning

- Bagging and boosting are methods that combine estimators (in parallel or sequentially)
- they are often applied to decision trees
- ▶ they reach state-of-the-art performance in several supervised learning tasks [Fernández-Delgado et al., 2014]
- they are an active area of research (both theoretically and for practical applications).

https://scikit-learn.org/stable/modules/ensemble.html

Aggregating to reduce the variance

- usual setup : input space \mathcal{X} , output space \mathcal{Y} .
- we note z_b a dataset sampled from the unknown distribution ρ . If we sample b different datasets, $b \in [1, B]$,

$$z_b = \{(x_{1b}, y_{1b}), \dots (x_{nb}, y_{nb})\}$$
 (1)

• we note \hat{f}_{z_b} the estimator obtained after learning with z_b .

Aggregate to reduce the variance

we note z_b a dataset sampled from the unknown distribution ρ . If we sample b different datasets, $b \in [1, B]$,

$$z_b = \{(x_{1b}, y_{1b}), \dots (x_{nb}, y_{nb})\}$$
 (2)

• we note \hat{f}_{z_b} the estimator obtained after learning with z_b .

Aggregating consists in using as an estimator :

for regression

$$\hat{f}_B = \frac{1}{B} \sum_{b=1}^{B} \hat{f}_{z_b}$$
 (3)

for classification

$$\hat{f}_B(x) = \arg\max_{j} |\{b, \hat{f}_{z_b}(x) = j\}|$$
 (4)

Bootstrapping and bagging

If B is large, it is not possible to sample B independent datasets with n samples, from a finite dataset.

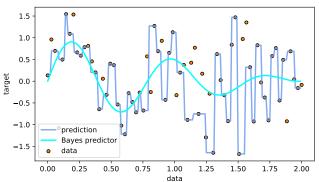
Bootstrapping consists in sampling B times a sample dataset with n elements with replacement.

Bagging is the combination of bootstrapping and aggregating.

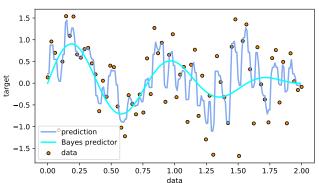
Out-of-bag error

Vocabulary: for an observation (x_i, y_i) , the **out-of-bag error** is the mean error on this observation among the estimators that were trained on a bootstrap dataset **not** containing it.

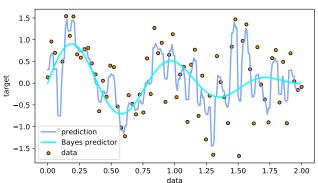
Bagging regression number of estimators: 1 test error: 1.22E+00 Bayes risk: 7.00E-01



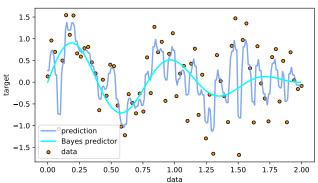
Bagging regression number of estimators: 10 test error: 9.09E-01 Bayes risk: 7.00E-01



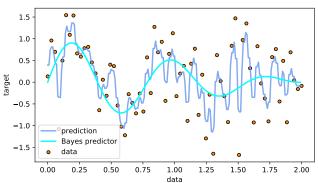
Bagging regression number of estimators: 20 test error: 9.39E-01 Bayes risk: 7.00E-01



Bagging regression number of estimators: 50 test error: 8.95E-01 Bayes risk: 7.00E-01



Bagging regression number of estimators: 100 test error: 8.81E-01 Bayes risk: 7.00E-01



Individual estimators

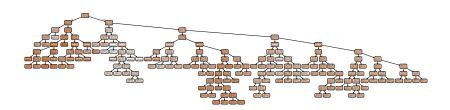


Figure – Estimator used in the averaging

Individual estimators

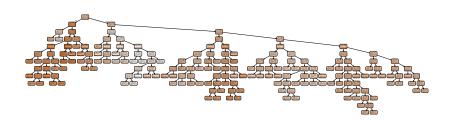


Figure – Other estimator used in the averaging

Possible issues

- number of trees to compute (B) in order to have a stable out-of-bag error.
- need for storing all the base estimators
- ▶ the final estimator is a *black-box* model.

Random forest

- Random forest [Breiman, 2001] is a bagging method with binary CART estimators, with additional randomness added to the choice of segmentation variables.
- the goal is to increase the independence between the base trees.
- although the theoretical properties of random forest are not yet fully understood, it is an important benchmark in supervised learning (but for instance not when the problem can be solved in a linear way).

Variance of an average : independent variables

Recall that if $(X_k)_{k\in\mathbb{N}}$ is a sequence of i.i.d. real variables that have a moment of order 2, and hence a variance σ^2 , and an expected value of m. Then if $S_B = \frac{1}{B} \sum_{i=1}^B X_i$

$$Var(S_B) = \sum_{i=1}^{B} Var(\frac{1}{B}X_i)$$

$$= \frac{1}{B^2} \sum_{i=1}^{B} Var(X_i)$$

$$= \frac{\sigma^2}{B}$$
(5)

Variance of the average : correlated variables

However, if the X_i are identically distributed but correlated with correlation c, then we admit that

$$Vac(S_B) = c\sigma^2 + \frac{1-c}{B}\sigma^2 \tag{6}$$

Random forest diminishes c by adding some randomness in the choice of the segmentation variables.

Random forest

Let d be the number of features. Let $m \le d$ be an integer.

Result: Random forest estimator

for $b \in [1, B]$ do

Sample a bootstrap dataset z_b ;

Estimate $\hat{f_{z_b}}$ with **randomization** of the variables. The search of the optimal segmentation is done on m randomly sampled variables.;

end

$$\hat{f}_B = \frac{1}{B} \sum_{b=1}^{B} \hat{f}_{z_b} \tag{7}$$

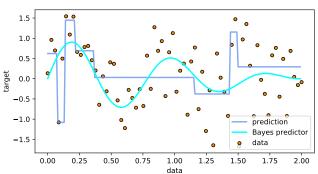
return $\hat{f_B}$

Algorithm 1: Random forest

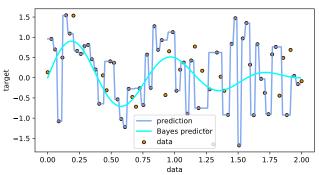
Remarks on random forest

- ► With random forest, it is possible to use a more brutal pruning (for instance using only trees of depth 2)
- ▶ To tune *m*, heuristics are used. Common ones include :
 - $ightharpoonup m = \sqrt{d}$ for classification (d is the number of features)
 - ightharpoonup m = d/3 for regression
 - cross validation.

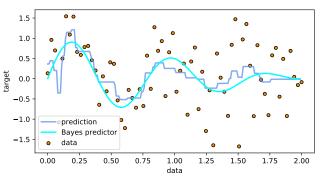
Random forest regression number of estimators: 1 max depth: 3 test error: 9.64E-01 Bayes risk: 7.00E-01



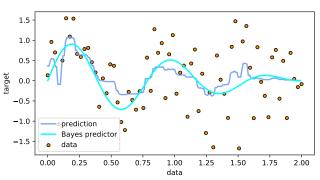
Random forest regression number of estimators: 1 max depth: 30 test error: 1.26E+00 Bayes risk: 7.00E-01



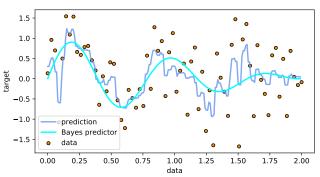
Random forest regression number of estimators: 10 max depth: 3 test error: 7.54E-01 Bayes risk: 7.00E-01



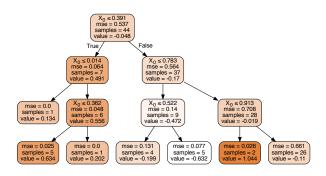
Random forest regression number of estimators: 50 max depth: 3 test error: 7.53E-01 Bayes risk: 7.00E-01



Random forest regression number of estimators: 20 max depth: 5 test error: 8.09E-01 Bayes risk: 7.00E-01



Estimators



Importance of the variables

Although the obtained estimator is hard to interpret, we can still study the statistical importance of variables for the prediction.

Mean decrease accuracy

Consider a variable j. For a given tree b

- compute the out-of-bag error.
- ▶ shuffle the values of *j* in the out-of-bag sample.
- compute the new out-of-bag error
- \triangleright store the decrease in prediction quality, D_j^b .

Average the D_j^b on b in order to measure the importance of the variable j.

Boosting

- ▶ While bagging is a random method, **boosting** is an adaptive method.
- Boosting builds on weak classifiers, and like bagging, aggregates them, for instance with a weighted sum.
- Mowever, the weak classifiers are built **sequentially**. The classifier b+1 adapts the classifier b by focusing on improving the prediction of incorrectly classified training samples.
- Several variants exists (in the weighting, loss function, aggregating method, etc).

Adaboost

- Original boosting algorithm [Freund and Schapire, 1996]
- $\mathcal{Y} = \{-1, 1\}$ (but can also be adapted to a regression problem)
- ▶ Given the sample dataset $z = \{(x_1, y_1), ...(x_n, y_n)\}$. We learn a sequence of estimators $(\delta_m)_{m \in [1..B]}$ (often CART).

FTML Ensemble learning Boosting

Initialize the weights $w = \{w_i = \frac{1}{n}, i = 1..n\};$ for m = 1..B do

Estimate δ_m on the weighted dataset. Compute the error rate

$$\hat{\epsilon_m} = \frac{\sum_{i=1}^n w_i 1(\delta_m(x_i) \neq y_i)}{\sum_{i=1}^n w_i}$$
 (8)

Compute the logit of $\hat{\epsilon_m}$

$$c_m = \log\left(\frac{1 - \hat{\epsilon_m}}{\hat{\epsilon_m}}\right) \tag{9}$$

Update the weights

$$w_i \leftarrow w_i \exp\left(c_m 1(\delta_m(x_i) \neq y_i)\right)$$
 (10)

end

$$\hat{\delta_B} = \operatorname{sign}\left(\sum_{k=1}^{B} \delta_{z_k}\right) \tag{11}$$

return $\hat{f_B}$

Algorithm 2: Adaboost for binary classification (adaptive boos-



Remarks

We need to enforce that $c_m \geq 0$. It is verified if $\hat{\epsilon_m} \geq \frac{1}{2}$.

$$c_m = \log\left(\frac{1 - \hat{\epsilon_m}}{\hat{\epsilon_m}}\right) \tag{12}$$

- ▶ It is experimentally shown that if the weak classifiers are *stump trees* (2 leaves), then adaBoost performs better than one given tree with a number of leaves equal to the number of iterations of adaBoost (hence a comparable computation time).
- ▶ A number of leaves *q* between 4 and 8 for the weak classifiers is often recommended.
- adaBoost can be adapted to multiclass classification and regression [Schapire, 2003].

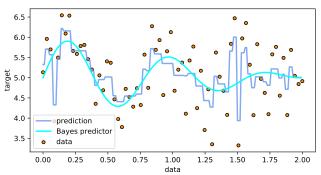
Gradient tree boosting

- differentiate with respect to the predictions of the tree (compute a gradient)
- approximate the gradient by a regression tree
- update the tree following this gradient.

Gradient tree boosting

Algorithm 3: Gradient tree boosting

Gradient boosting regression number of estimators: 40 max depth: 3 test error: 8.39E-01 Bayes risk: 7.00E-01



Gradient tree boosting

Many parameter can be tuned :

- maximum depth of the estimators
- ▶ shrinkage η ($\hat{f_m} = \hat{f_{m-1}} + \eta \gamma_m \delta_m$, with $0 \le \eta \le 1$). If η is low, e.g. ≤ 0.1 , it can lead to a slower convergence but might prevent overfitting.
- number of trees learned B

https://www.analyticsvidhya.com/blog/2016/02/ complete-guide-parameter-tuning-gradient-boosting-gbm-pytho It is necessary to experiment, read documentations, cross validate, use heuristics, etc.

Extrem gradient boosting

- ➤ XGBoost : variant of gradient boosting, very efficient on several benchmarks [Chen and Guestrin, 2016]
- can exploit parallelization
- https: //xgboost.readthedocs.io/en/latest/parameter.html
- https://github.com/dmlc/xgboost

Catboost

➤ See also : Catboost https://catboost.ai/

Ensemble learning

Bagging Random forest Boosting

Statistical learning

Bounding the estimation error Interpolation regime and double descent

Local averaging methods

Supervised learning Density estimation

Statistical learning

- ▶ We come back to the statistical analysis of supervised learning.
- ▶ More precisely, to that of empirical risk minimization.

Reminder on risks

Let I be a loss. Generalization error:

$$R(f) = E_{(X,Y)\sim\rho}[I(Y,f(X))] \tag{15}$$

The empirical risk (ER) of an estimator f writes

$$R_n(f) = \frac{1}{n} \sum_{i=1}^n I(y_i, f(x_i))$$
 (16)

Remember that the risks depends on the loss 1.

Risk decomposition

- ► f* : Bayes predictor
- F : Hypothesis space
- $ightharpoonup f_n$: estimated predictor (hence in F).

$$E\left[R(f_n)\right] - R^* = \left(E\left[R(f_n)\right] - \inf_{f \in F} R(f)\right) + \left(\inf_{f \in F} R(f) - R^*\right)$$
(17)

Risk decomposition

- ▶ f* : Bayes predictor
- F : Hypothesis space
- $ightharpoonup f_n$: estimated predictor.

When doing empirical risk minimization, f_n is obtained by minimization of the empirical risk.

However:

- we have seen that in many cases, finding the exact minimizer of the empirical risk might be computationnally hard.
- also, a natural question is whether it is sufficient to have an approximate minimizer of the empirical risk, as the empirical risk is an approximation of the generalization error.

Risk decomposition

Estimation error (variance term, fluctuation error, stochastic error) : depends on D_n , F, f_n .

$$E\Big[R(f_n)\Big]-\inf_{f\in F}R(f)\geq 0$$

Approximation error (bias term): depends on f^* and F, not on f_n , D_n .

$$\inf_{f\in F}R(f)-R^*\geq 0$$

It is also possible to consider the **Optimization error** : depends on D_n , F, f_n .

$$E\left[R(\hat{f}_n) - R(f_n)\right] \tag{18}$$

where \hat{f}_n is an approximate solution to the optimization problem.

We will now focus on the estimation error

$$E\Big[R(f_n)\Big]-\inf_{f\in F}R(f)\geq 0$$

 f_n is the empirical risk minimizer We consider the best estimator in hypothesis space

$$f_a = \underset{h \in F}{\operatorname{arg min}} R(h)$$

We have seen that

$$R(f_n) - R(f_a) \le 2 \sup_{h \in F} |R(h) - R_n(h)|$$
 (19)

Deterministic bound on the estimation error

$$f_a = \underset{h \in F}{\arg \min} R(h) \tag{20}$$

$$f_n = \underset{h \in F}{\arg\min} \, R_n(h) \tag{21}$$

$$R(f_n) - R(f_a) = (R(f_n) - R_n(f_n)) + (R_n(f_n) - R_n(f_a)) + (R_n(f_a) - R(f_a))$$
(22)

Deterministic bound on the estimation error

$$f_{a} = \underset{h \in F}{\operatorname{arg min}} R(h)$$

$$f_{n} = \underset{h \in F}{\operatorname{arg min}} R_{n}(h)$$
(23)

$$R(f_{n}) - R(f_{a}) = (R(f_{n}) - R_{n}(f_{n})) + (R_{n}(f_{n}) - R_{n}(f_{a})) + (R_{n}(f_{a}) - R(f_{a}))$$
(24)

But by definition f_n minimizes R_n , so $(R_n(f_n) - R_n(f_a)) \le 0$.

Deterministic bound on the estimation error

$$f_a = \underset{h \in F}{\arg \min} R(h)$$

$$f_n = \underset{h \in F}{\arg \min} R_n(h)$$
(25)

$$R(f_{n}) - R(f_{a}) = (R(f_{n}) - R_{n}(f_{n})) + (R_{n}(f_{n}) - R_{n}(f_{a})) + (R_{n}(f_{a}) - R(f_{a}))$$
(26)

But by definition f_n minimizes R_n , so $(R_n(f_n) - R_n(f_a)) \le 0$. Finally:

$$R(f_n) - R(f_a) \le 2 \sup_{h \in F} |R(h) - R_n(h)|$$
 (27)

If we are able to bound $\sup_{h\in F}|R(h)-R_n(h)|$, then we have a bound on the estimation error.

Theorem

Weak law of large numbers

Let $(X_i)_{i\in\mathbb{N}}$ be a sequence of i.i.d. variables that have a moment of order 2. We note m their expected value. Then

$$\forall \epsilon > 0, \lim_{n \to +\infty} P\left(\left|\frac{1}{n} \sum_{i=1}^{n} X_i - m\right| \ge \epsilon\right) = 0$$
 (28)

We say that we have convergence in probability.

Theorem

Weak law of large numbers

Let $(X_i)_{i\in\mathbb{N}}$ be a sequence of i.i.d. variables that have a moment of order 2. We note m their expected value. Then

$$\forall \epsilon > 0, \lim_{n \to +\infty} P\left(\left|\frac{1}{n} \sum_{i=1}^{n} X_i - m\right| \ge \epsilon\right) = 0$$
 (29)

We say that we have convergence in probability.

However, this is only an asymptotical result : it is a limit for $n \to +\infty$.

If we are able to bound $\sup_{h\in F} |R(h) - R_n(h)|$, then we have a bound on the estimation error.

To do this, we will use some other mathematical results :

- Boole's inequality
- Hoeffding's inequality (non-asymptotical probabilistic bound)

Boole's inequality

Proposition

Let A_1, A_2, \ldots , be accountable set of events of a probability space $\{\Omega, \mathcal{F}, P\}$. Then.

$$P\Big(\cup_{i\geq 1}A_i\Big)\leq \sum_{i\geq 1}P(A_i) \tag{30}$$

This set might be infinite.

Boole's inequality

Proposition

Let A_1, A_2, \ldots , be accountable set of events of a probability space $\{\Omega, \mathcal{F}, P\}$. Then.

$$P\Big(\cup_{i\geq 1} A_i\Big) \leq \sum_{i\geq 1} P(A_i) \tag{31}$$

This set might be infinite.

Exercice 1: Proove the proposition.

Hoeffding's inequality

Theorem

Hoeffding's inequality

Let $(X_i)_{1 \le i \le n}$ be n i.i.d real random variables such that $\forall i \in [1, n]$, $X_i \in [a, b]$ and $E(X_i) = \mu \in \mathbb{R}$. Let $\bar{\mu} = \frac{1}{n} \sum_{i=1}^n X_i$. Then $\forall \epsilon > 0$.

$$P(|\bar{\mu} - \mu| \ge \epsilon) \le 2 \exp\left(-\frac{2n\epsilon^2}{(b-a)^2}\right)$$

We admit this theorem.

Setting

- Supervised learning.
- Finite space of estimator F.
- ▶ The loss I is uniformly bounded : $I(\hat{y}, y) \in [a, b]$ with a and b real numbers.

Step 1

We have seen that

$$R(f_n) - R(f_a) \le 2 \sup_{h \in F} |R(h) - R_n(h)|$$
 (32)

As a consequence, for all $t \geq 0$:

$$P\Big(R(f_n) - R(f_a) \ge t\Big) \le P\Big(2\sup_{h \in F} |R(h) - R_n(h)| \ge t\Big)$$
 (33)

Step 2

The fact that

$$2\sup_{h\in F}|R(h)-R_n(h)|\geq t\tag{34}$$

is equivalent to:

$$\cup_{h\in F} \Big(2|R(h)-R_n(h)|\geq t\Big) \tag{35}$$

Conclusion

Exercice 2: Using Boole's inequality and Hoeffding's inequality, show that

$$P(R(f_n) - R(f_a) \ge t) \le 2|F| \exp\left(-\frac{nt^2}{2(b-a)^2}\right)$$
 (36)

Conclusion

We write

$$\delta = 2|F| \exp\left(-\frac{nt^2}{2(b-a)^2}\right) \tag{37}$$

Exercice 3:

We assume that b-a=1. Show that with probability $\geq 1-\delta$,

$$R(f_n) \le R(f_a) + 2\sqrt{\frac{\log(|F|) + \log(\frac{2}{\delta})}{2n}}$$
 (38)

Generalization

It is possible to generalize to infinite sets :

- by sampling F
- by using Rademacher complexity and Vapnik Vapnik-Chervonenkis theory.

This classical bound on the statistical error does not guarantee that the generalization error is small when log(|F|) is large.

Interpolation regime

- ▶ If the number of parameters is sufficient, it is possible to have $R_n(f_n)=0.$
- In that case, it seems that the statistical error might to be way smaller than the previous bound.

For instance, for Wide Resnet, $\frac{p}{n} = 179$ with

- p: number of parameters in the netwok
- n : number of samples

Double descent

[Belkin et al., 2019]

Ensemble learning

Bagging Random forest Boosting

Statistical learning

Bounding the estimation error Interpolation regime and double descent

Local averaging methods Supervised learning Density estimation

Local averaging methods

Local averaging methods : approximation $\mbox{without}$ optimization of an empirical risk.

Setting

Let I be a loss. Generalization error :

$$R(f) = E_{(X,Y)\sim\rho}[I(Y,f(X))] \tag{39}$$

Bayes estimator :

$$f^* = \underset{f \text{ measurable}}{\text{arg min }} R(f)$$
 (40)

Bayes risk:

$$R^* = E_X \left[\inf_{y \in \mathcal{Y}} E_{Y \sim dP(Y|X)} [I(Y,y)|X] \right]$$
 (41)

As always, the law of (X, Y) and of Y|X are unknown.

Local averaging

Local averaging:

- **b** based on the dataset D_n , compute an approximation of the law Y|X or (Y,X), without optimization of an empirical risk.
- use this approximation in the estimator (e.g. to compute an approximate conditional expected value.)

Local averaging : regression

 $\tilde{f}(\boldsymbol{x})$: local averaging estimator, in the case of regression, squared loss, we can use

$$\tilde{f}(x) = \int_{Y \in \mathbb{R}} y \, d\hat{P}(Y|X = x) \tag{42}$$

Local averaging : classification

 $\tilde{f}(\boldsymbol{x})$: local averaging estimator, in the case of binary classification, squared loss, we can use

$$\tilde{f}(x) = \arg\max_{z \in \mathcal{Y}} \hat{P}(Y = z | X = x)$$
(43)

Linear estimators

The question is then : how to choose the approximation $d\hat{P}(Y|X=x)$? **Linear estimators**

$$\hat{dP}(Y|X=x) = \sum_{i=1}^{n} \hat{w}_i(x)\delta_{y_i}$$
 (44)

 δ_{y_i} is the Dirac mass in y_i .

- $\forall i, \hat{w}_i(x) \geq 0$
- $\sum_{i=1}^{n} \hat{w}(x) = 1$

Linear estimators

Linear estimators

$$\hat{dP}(Y|X=x) = \sum_{i=1}^{n} \hat{w}_i(x)\delta_{y_i}$$
 (45)

 δ_{y_i} is the Dirac mass in y_i .

- $\forall i, \hat{w}_i(x) \geq 0$
- $\sum_{i=1}^{n} \hat{w}(x) = 1$

Application to regression:

$$\tilde{f}(x) = \sum_{i=1}^{n} \hat{w}_i(x) y_i \tag{46}$$

Linear estimators

Linear estimators

$$\hat{dP}(Y|X=x) = \sum_{i=1}^{n} \hat{w}_i(x)\delta_{y_i}$$
(47)

 δ_{y_i} is the Dirac mass in y_i .

- $\forall i, \hat{w}_i(x) \geq 0$
- $\sum_{i=1}^{n} \hat{w}(x) = 1$

Application to classification:

$$\tilde{f}(x) = \underset{j \in \{0,1\}}{\arg \max} \sum_{i=1}^{n} \hat{w}_{i}(x) 1_{y_{i}=j}$$
(48)

Choice of the weights

Linear estimators

$$\hat{dP}(Y|X=x) = \sum_{i=1}^{n} \hat{w}_i(x)\delta_{y_i}(y)$$
(49)

 δ_{y_i} is the Dirac mass in y_i .

- $\forall i, \hat{w}_i(x) \geq 0$
- $\sum_{i=1}^{n} \hat{w}(x) = 1$

For any sample i, the weight function $\hat{w}_i(x)$ should be

- \triangleright closer to 1 for training point x_i that are close to x.
- closer to 0 for training point x_i that are far from x.

Choice of the weights

Linear estimators

$$\hat{dP}(Y|X=x) = \sum_{i=1}^{n} \hat{w}_{i}(x)\delta_{y_{i}}(y)$$
 (50)

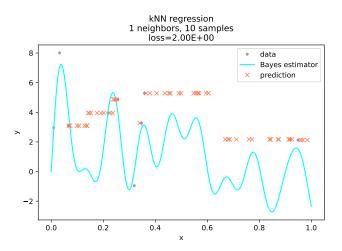
 δ_{y_i} is the Dirac mass in y_i .

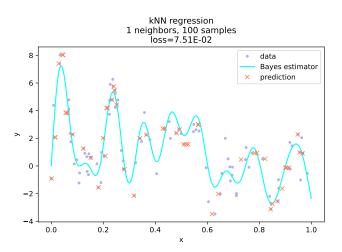
- $\forall i, \hat{w}_i(x) \geq 0$
- $\sum_{i=1}^n \hat{w}(x) = 1$

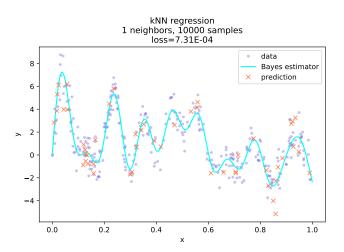
Three possibilities:

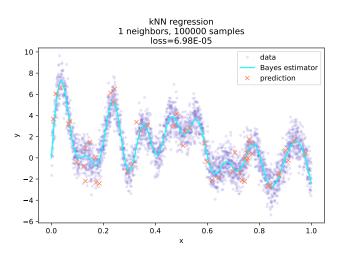
- partition estimators
- nearest neighbors
- ► Nadaraya-Watson (kernel regression)

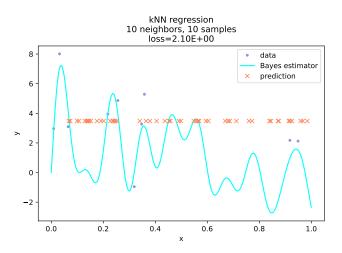
Given $k \geq 1$, and a metric d on \mathcal{X} , average the predictions of the k nearest neighbors (for regression) or take the majority vote (for classification).

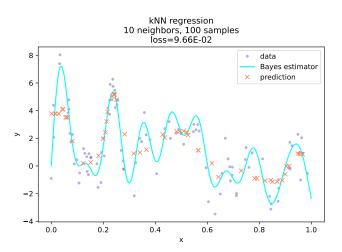


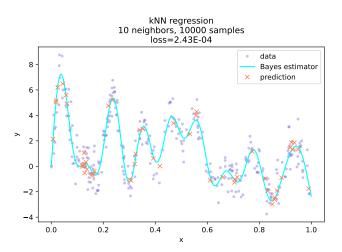


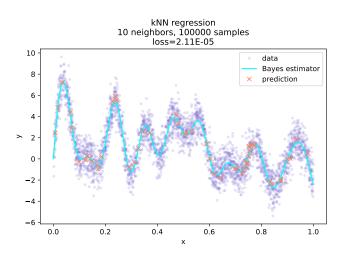












Given $k \ge 1$, and a metric d on \mathcal{X} , average the predictions of the k nearest neighbors (for regression) or take the majority vote (for classification).

Exercice 4 : What is $\hat{w}_i(x)$?

k is a hyperparameter, hence it must be tuned, for instance with cross validation.

- too small k : underfitting
- ▶ too large k : overfitting

Nearest neighbors search

```
The search for nearest neighbors is a problem itsself! https://scikit-learn.org/stable/modules/neighbors.html https://en.wikipedia.org/wiki/K-d_tree https://en.wikipedia.org/wiki/Ball_tree
```

Partition estimators

$$\mathcal{X} = \cup_{j \in J} A_j$$
.

A_1	A_2	A_3	A_4	A_5
A_6	A_7	A_8	A_9	A_{10}
A_{11}	A_{12}	A_{13}	A_{14}	A_{15}
A_{16}	A_{17}	A_{18}	A_{19}	A_{20}
A_{21}	A_{22}	A_{23}	A_{24}	A_{25}

For each x, average the predictions of the samples that are in the same A_i as x. We can note it A(x).

Partition estimators

$$\mathcal{X}=\cup_{j\in J}A_j.$$

A_1	A_2	A_3	A_4	A_5
A_6	A_7	A_8	A_9	A_{10}
A_{11}	A_{12}	A_{13}	A_{14}	A_{15}
A_{16}	A_{17}	A_{18}	A_{19}	A_{20}
A_{21}	A_{22}	A_{23}	A_{24}	A_{25}

For each x, average the predictions of the samples that are in the same A_j as x. We can note it A(x).

Exercice 5: What is $\hat{w}_i(x)$?

Partition estimator

$$\hat{w}_i(x) = \frac{1_{x_i \in A(x)}}{\sum_{k=1}^{n} 1_{x_k \in A(x)}}$$
 (51)

Partition estimator

Exercice 6: We have seen in previous classes one example of partition estimator. What is it?

Kernel regression (Nadaraya-Watson)

We consider a non-negative kernel function $k: \mathcal{X} \times \mathcal{X} \to \mathbb{R}_+$ and

$$\hat{w}_i(x) = \frac{k(x, x_i)}{\sum_{i=1}^{n} k(x, x_i)}$$
 (52)

Non-negative kenrels

Often

$$k(x, x') = \frac{1}{h^d} q(\frac{x - x'}{h})$$
 (53)

with d the dimension, h a bandwidth parameter.

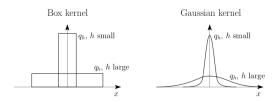


Image from [Bach, 2021].

Non-negative kenrels

$$k(x,x') = \frac{1}{h^d}q(\frac{x-x'}{h}) \tag{54}$$

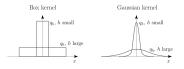


Image from [Bach, 2021].

- $\blacktriangleright \text{ Box kernel}: q(x) = 1_{||x|| \le 1}$
- Gaussian kernel : $q(x) = e^{-\frac{||x||^2}{2}}$

Remark

These kernels are not exactly the same as the ones we mentioned earlier (positive-definite kernels).

These kernels are more simply non-negative (less specific).

Estimator:

$$f(x) = \frac{\sum_{i=1}^{n} k(x, x_i) y_i}{\sum_{i=1}^{n} k(x, x_i)}$$
 (55)

Curse of dimensionality

It is posible to show, that under some simple regularity assumptions on the target, the convergence rate of the error of these estimators, as a function of n, is $\mathcal{O}(n^{-\frac{2}{d+2}})$, where d is the underlying dimension.

In order to have an error smaller than ϵ , we need to have

$$n \ge \left(\frac{1}{\epsilon}\right)^{\frac{d+2}{2}} \tag{56}$$

- ▶ It is not easy to exploit a higher regularity of the target function (no adaptivity to the regularity)
- ▶ It is not possible to learn with these methods in high dimension.

Kernel density estimation

```
It is possible to use similar ideas to perform Kernel density estimation (KDE). (Again, here it is juste a non-negative kernel) https://francisbach.com/cursed-kernels/https://seaborn.pydata.org/generated/seaborn.jointplot.html https://fr.wikipedia.org/wiki/Estimation_par_noyau https://en.wikipedia.org/wiki/Kernel_density_estimation
```

References I



 $\label{learning Theory from First Principles Draft.} Learning Theory from First Principles Draft.$

Book Draft, page 229.

Belkin, M., Hsu, D., Ma, S., and Mandal, S. (2019).

Reconciling modern machine-learning practice and the classical bias-variance trade-off.

Proceedings of the National Academy of Sciences of the United States of America, 116(32):15849–15854.

Breiman, L. (2001).

Random Forests.

Machine Learning, 45(1):5-32.

References II



Chen, T. and Guestrin, C. (2016).

XGBoost: A scalable tree boosting system.

Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 13-17-Augu: 785-794.



Fernández-Delgado, M., Cernadas, E., Barro, S., and Amorim, D. (2014).

Do we need hundreds of classifiers to solve real world classification problems?

Journal of Machine Learning Research, 15:3133–3181.



Freund, Y. and Schapire, R. E. (1996).

Experiments with a New Boosting Algorithm.

Proceedings of the 13th International Conference on Machine Learning, pages 148-156.

References III



Schapire, R. E. (2003).

The Boosting Approach to Machine Learning : An Overview BT - Nonlinear Estimation and Classification.

Nonlinear Estimation and Classification, 171(Chapter 9):149–171.