

ML, Decision Trees and Random Forests Data science Certificate

Aymeric DIEULEVEUT

April 2025

1 Introduction to Supervised Learning



what you are not allowed to forget

for any algorithm $X \rightarrow \boxed{\quad} \rightarrow Y$

how do I
train / what training is

2 Decision Trees

Dec Tree

3 Random Forests

Neural Networks

Math

Lin Regress°

Log Regress°

Outline

1 Introduction to Supervised Learning

2 Decision Trees

3 Random Forests

A Reminder on Supervised Learning

X features

goal : predict 7

- Observations: $(X_i)_{i=1,\dots,n}$. Each X_i has d components.

- Outputs: $(Y_i)_{i=1,\dots,n}$

3 Examples:

- Iris Dataset:

prediction iris type from n features

$\underbrace{d}_{\text{number}} = \text{number}$
of features

$n =$ number of obs.

=

- MNIST:

predictions a digit from its image / pixels

s features

- House price prediction



$X = (s_f \text{ Poo} \text{tage}; s_f \text{ Poo} \text{tage balcony},$
 $\text{floor}; \text{neighborhood})$
 (days of sun)

Example 1: House price Prediction - Regression

Goal: observe the features for a house, predict its price.

House price Prediction dataset

- Output : house value

$$\rightarrow Y \in \mathbb{R} \rightarrow \text{regression task}$$

- 1500 examples

$$\rightarrow n = 1500$$

- 81 features: Construction year, Neighborhood, Surface, Floor, Balcony: both quantitative and qualitative, some ordinal variables.

$$\rightarrow X \in \mathbb{R}^{81}.$$

Example 2: Iris Dataset - Classification



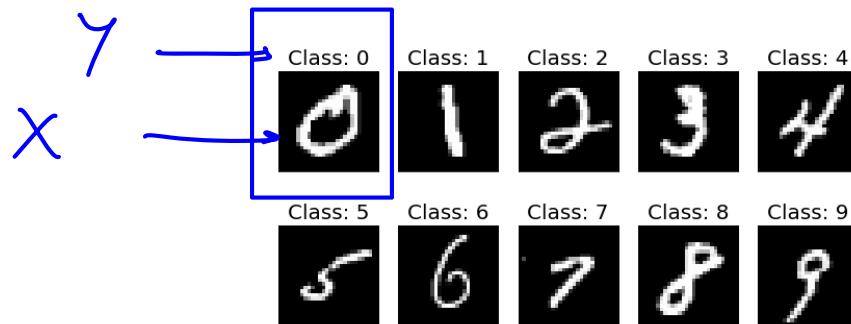
Figure: Iris: setosa, versicolor, virginica

Goal: identify flower species based on petal/sepal characteristics.

Iris Dataset:

- 3 classes : `['setosa', 'versicolor', 'virginica']`
→ $Y \in \{0, 1, 2\}$ encoding each class.
- 50 points per class
→ $n = 150$
- 4 features: `['sepal length', 'sepal width', 'petal length', 'petal width']`
 $X \in \mathbb{R}^4$, → $d = 4$

Example: MNIST Dataset



Goal: recognize an object on an image, e.g., a handwritten digit on its image.

Digit recognition from an image:

- 10 classes :

$$Y \in \{0, 1, 2, \dots, 9\}$$

- 60k images, (50k train, 10k test), balanced

60k

$$n = 60000$$

- Observation: pixels. 28*28 pixels per image (gray-scale).

$$X \in \mathbb{R}^{784}, d = 784 \text{ features.}$$

Summary of those three datasets

	Iris	MNIST	House price prediction
Number of examples n	150	60k	1560
Number of the features d	4	784	81
Type of the features	quantitative	pixels → quantitative	quantitative / qualitative
Space \mathcal{Y}	3 classes	$\{0, \dots, 9\}$	\mathbb{R}
Task	classification	classification	regression

size dataset

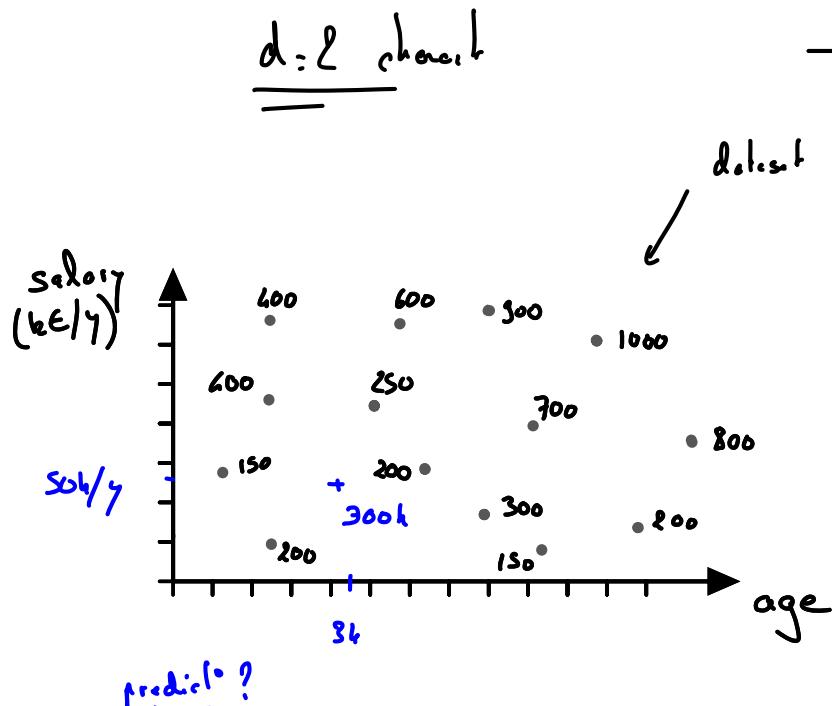
≈ 600

$\approx 60M$

$\approx 100,000$

Toy datasets

- In regression:



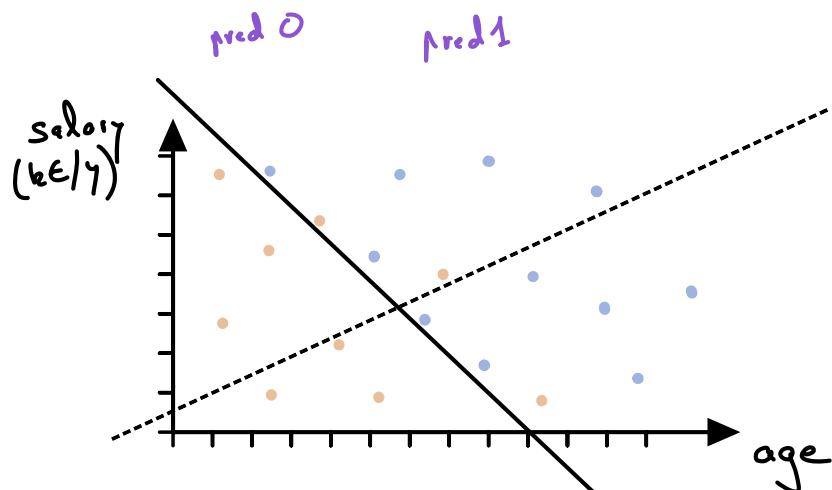
\rightarrow $d=1$ characteristic

$y = \text{yield}$

Learn the link X and Y .

$X_1 = \text{Temp}$

- In classification:

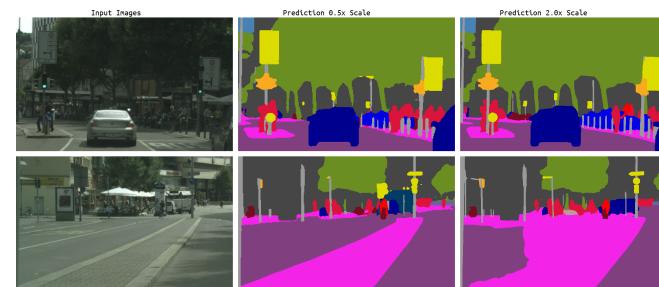


→ used a let : $d = 2$, $y \in \{0, 1\}$

More examples



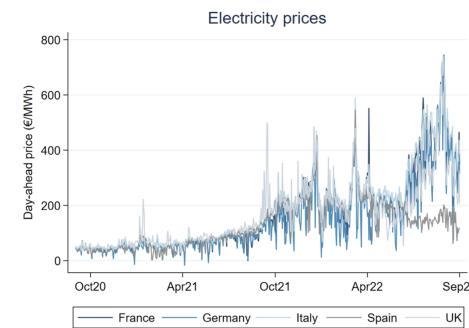
Spam detection



Semantic image segmentation



Sentiment classification



Electricity price forecasting

Supervised Learning

Supervised Learning -

- Attributes / features / characteristics / input:
 $X \in \mathcal{X} = \mathbb{R}^d$
- Output / Responses / label:
 - $Y \in \{-1, 1\}$ - binary classification,
 - $Y \in \{0, 1, 2, \dots, K-1\}$ - multi-class classification,
 - $Y \in \mathbb{R}$ - regression.
- Phenomenon distribution $(X, Y) \sim P$.

Main ML assumptions

Dataset

(X_i, Y_i)

iid

indp. identically
distributed.



Be careful if distribution of the phenomena changes

Take always carefull

Supervised Learning

Supervised Learning

- Attributes / features / characteristics / input:

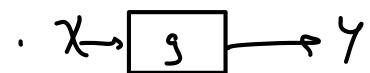
$$X \in \mathcal{X} = \mathbb{R}^d$$

- Output / Responses / label :

- $Y \in \{-1, 1\}$ - binary classification,
- $Y \in \{0, 1, 2, \dots, K-1\}$ - multi-class classification,
- $Y \in \mathbb{R}$ - regression.

- Phenomenon distribution $(X, Y) \sim \mathbb{P}$.

- A classifier or predictor is a measurable function $g : \mathcal{X} \rightarrow \mathcal{Y}$.



How to quantify the quality of a predictor?

→ introduce a risk function.

Loss function, Generalization Risk and data-dependent predictors

Loss Function

- Loss function: $\ell(y, g(x))$ quantifies the quality of the prediction $g(x)$ of y , e.g.:
 - ▶ 0-1 Loss (classification): $\ell(y, g(x)) = \begin{cases} 1 & y \neq g(x) \\ 0 & y = g(x) \end{cases}, y \in \mathcal{Y} = \{-1, 1\}$
 - ▶ Quadratic Loss (regression): $\ell(y, g(x)) = (y - g(x))^2, y \in \mathbb{R}^d$

$$\ell(y, g(x)) = (y - g(x))^2$$

(true - predict)²

difference

nb of error you make

Loss function, Generalization Risk and data-dependent predictors

Loss Function

- Loss function: $\ell(y, g(x))$ quantifies the quality of the prediction $g(x)$ of y , e.g.:
 - ▶ 0-1 Loss (classification): $\ell(y, g(x)) = \mathbb{1}_{y \neq g(x)}$, $y \in \mathcal{Y} = \{-1, 1\}$
 - ▶ Quadratic Loss (regression): $\ell(y, g(x)) = (y - g(x))^2$, $y \in \mathbb{R}^d$

Risk of a Decision Rule = average loss g classifier / regressor.

- Risk: $R(g) = \mathbb{E}[\ell(Y, g(x))] = \int \ell(y, g(x)) d\mathbb{P}(x, y)$, e.g.: $(x, y) \sim \mathbb{P}$
 - ▶ 0-1 Risk (classification): $R(g) = \mathbb{P}(Y \neq g(x))$
 - ▶ Quadratic Risk (regression): $R(g) = \mathbb{E}[|Y - g(x)|^2]$

↪ R is also referred to as *Generalization* risk, or *True* risk.

Find g that has a small Generalization risk.

Goal in ML: predict well on new data points - unseen in the dataset → this is what the generalization risk measures.

Loss function, Generalization Risk and data-dependent predictors

Loss Function

- Loss function: $\ell(y, g(x))$ quantifies the quality of the prediction $g(x)$ of y , e.g.:
 - ▶ 0-1 Loss (classification): $\ell(y, g(x)) = \mathbb{1}_{y \neq g(x)}$, $y \in \mathcal{Y} = \{-1, 1\}$
 - ▶ Quadratic Loss (regression): $\ell(y, g(x)) = (y - g(x))^2$, $y \in \mathbb{R}^d$

Risk of a Decision Rule = average loss

- Risk: $R(g) = \mathbb{E}[\ell(Y, g(x))] = \int \ell(y, g(x)) d\mathbb{P}(x, y)$, e.g.:
 - ▶ 0-1 Risk (classification): $R(g) = \mathbb{P}(Y \neq g(x))$
 - ▶ Quadratic Risk (regression): $R(g) = \mathbb{E}[|Y - g(x)|^2]$
- ↪ R is also referred to as Generalization risk, or True risk.

⚠ The distribution \mathbb{P} is unknown.

Data-dependent predictors

- Training set: $\mathcal{D}_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ (i.i.d. $\sim \mathbb{P}$) —
- Goal: build a **data dependent predictor / classifier** \hat{g}_n based on the training data
- That has a **minimal generalization risk** $R(\hat{g}_n)$. ↗ **inductive part**

$R(\hat{g}_n)$ = average loss on a “new point” $(X, Y) \sim \mathbb{P}$, independent from D_n

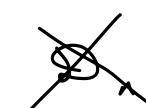
↔ Our goal is to predict well on inputs/outputs pairs that **we have not seen in the dataset**, i.e. *generalize* well.

Some high level messages:

- * what ML is !
- * a few examples
- * inductive thinking is — what learning is

Key assumption in ML modelling → data is iid, follows a fixed distrib

Approach: use a dataset D_n to learn a predictor \hat{g} .

overfitting } → we want \hat{g} to be good  ← not enough on a new point.

Supervised Learning: Summary

Supervised Learning: Framework Summary

- Training set: $D_n = \{(x_1, Y_1), \dots, (x_n, Y_n)\}$ (*i.i.d.* $\sim \mathbb{P}$)
- Decision rule (classifier, predictor): $g : \mathcal{X} \rightarrow \mathcal{Y}$
- Loss: $\ell(Y, g(x))$
- Risk: $R(g) = \mathbb{E}[\ell(Y, g(x))]$

Numerous methods:

- ① In regression:
 - ▶ linear, polynomial, kernel regression
 - ▶ regularized methods (Lasso, Ridge)
- ② In classification:
 - ▶ Logistic regression
 - ▶ SVM
- ③ For both
 - ▶ Nearest neighbors
 - ▶ Random forests
 - ▶ Neural networks
 - ▶ etc.

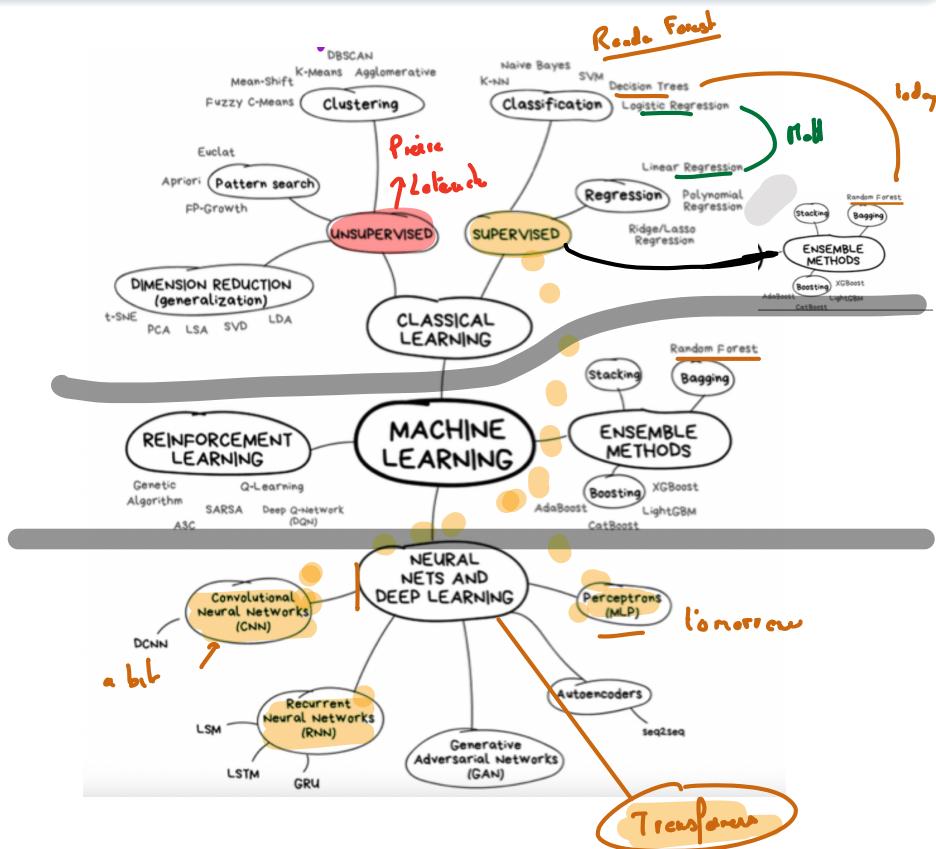


Figure:

Roadmap towards Understanding machine learning!

1. Understand conceptually how to train a model.



→ how to learn "weights"

→ what f_{Model} is

→ overfitting - generalization

single
validation

or cross validation

valid for any

2. Opening some



Linear Reg
DT

RF
NN

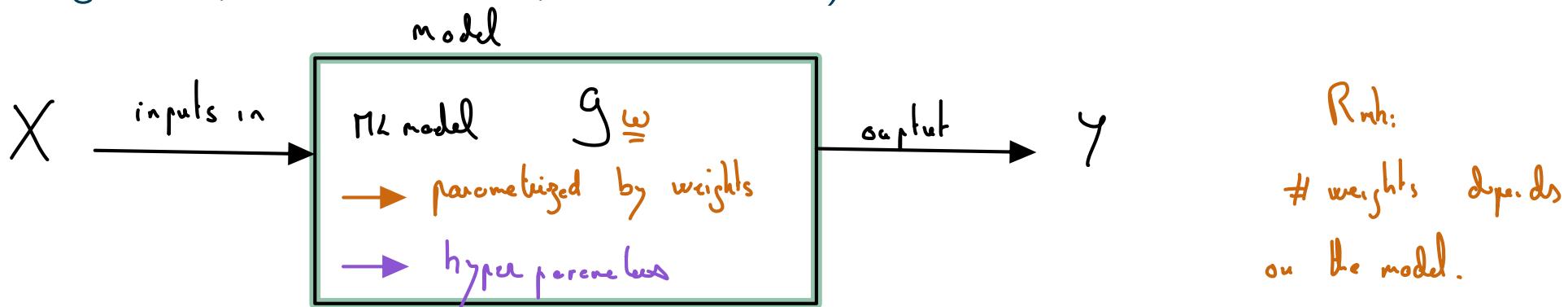
useful logic

S/A/I

Training Pipeline

① Fitting

💡 “Learn” the dependency between input and output. Each model (LR, Logistic regression, Neural networks, Decision Trees) → obtain the model !

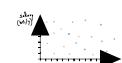


Linear regression on House price prediction

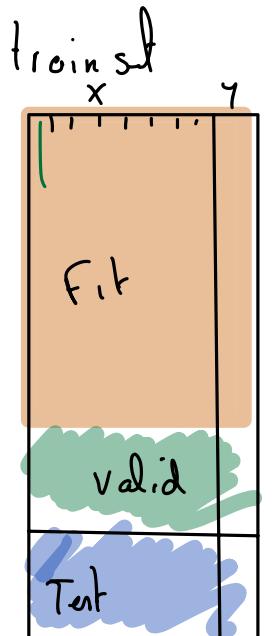
$$g(X_{\text{surf}}, X_{\text{bdc surf}}, X_{\text{Pcor}}, \dots) = b_0 + \omega_1 X_{\text{surf}} + \omega_2 X_{\text{bdc}} + \dots$$

weights to be learned.
→ #: 81

Fitting is done a part of the training dataset
Find a function that behaves well on it



$$(X_{p,t}, Y_{p,t})$$



Def

Fitting =

finding

ω

such that

g_ω is good on $(X_{p,t}, Y_{p,t})$

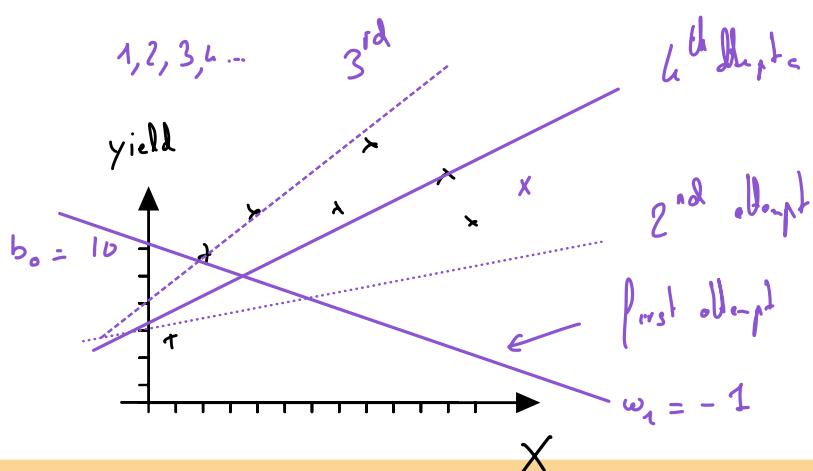
→ typically minimize the empirical risk

How to fit? main algorithm

is "incremental improvement"

$$\omega \mapsto R(g_\omega) = \frac{1}{n_{\text{fit}}} \sum_{i=1}^{n_{\text{fit}}} l(g_\omega(X_{p,t,i}), Y_{p,t,i})$$

quality of prediction (w.r.t. l)
on point i of fit set



Fit: → 95% time : "incremental improvements"

e.g. NN;
linear regression;
logistic regression; SVM; etc

→ 5% time : specific algorithm.

→ DT!, RF!

(
* or learned from
formula $(X^T X)^{-1} X^T Y$)

Fitting a method

Demo url (Lab): https://colab.research.google.com/github/adieulev/Certificate-DS-2025/blob/main/Labs/MWE_ML_init_student_version.ipynb

```
# Fetching the MNIST dataset
mnist = fetch_openml('mnist_784', version=1) .
X, y = mnist["data"], mnist["target"]
# Setting aside 10% of the data
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.9, random_state=42)
# Scaling the data
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train.astype(np.float64))
X_test_scaled = scaler.transform(X_test.astype(np.float64))
```

Preprocessing

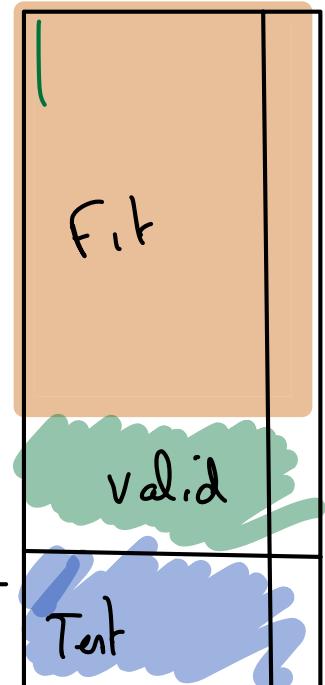
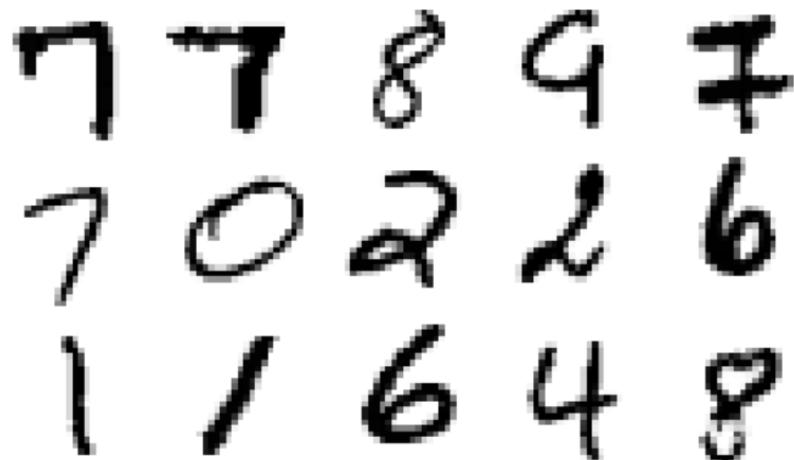


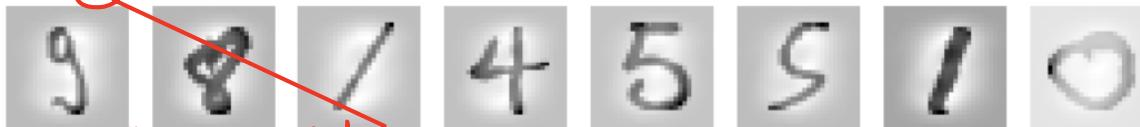
Figure: Importing a dataset



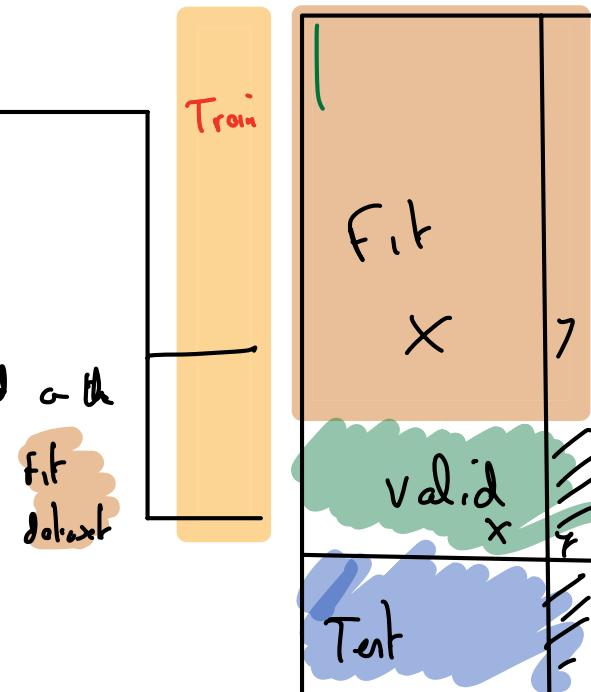
Fitting 3 methods

```
X_fit, X_val, y_fit, y_val = train_test_split(X_train, y_train,  
                                              test_size=0.2, random_state=42)  
  
from sklearn.tree import DecisionTreeClassifier  
  
dt_clf = DecisionTreeClassifier()  
dt_clf.fit(X_fit, y_fit)  
  
DecisionTreeClassifier()  
  
predictions = dt_clf.predict(X_train_scaled[:16])  
fig, axes = plt.subplots(1, 8, figsize=(10, 2)) # Adjust the figure size as needed  
for i, ax in enumerate(axes):  
    ax.imshow(X_train_scaled[i].reshape(28, 28), cmap=plt.cm.gray_r, interpolation='nearest')  
    ax.set_title(f"Predicted: {predictions[i]}\nTrue: {y_train.iloc[i]}")  
    ax.axis('off')
```

Predicted: 7 Predicted: 7 Predicted: 7 Predicted: 4 Predicted: 7 Predicted: 7 Predicted: 7 Predicted: 7
True: 9 True: 8 True: 1 True: 4 True: 5 True: 5 True: 1 True: 0



I on this image I predict 7 but it is 9.



bch
x,y used to fit
x used to predict
y to check if Id good

```
Os from sklearn.ensemble import RandomForestClassifier  
  
Os rf_clf = RandomForestClassifier(n_estimators=10, random_state=42) → define my object  
→ rf_clf.fit(X_fit, y_fit) → fit!  
  
RandomForestClassifier(n_estimators=10, random_state=42)  
  
from sklearn.linear_model import LogisticRegression  
  
log_reg = LogisticRegression(multi_class="multinomial", solver="lbfgs", max_iter=100, random_state=42)  
log_reg.fit(X_fit, y_fit) → fit!  
  
LogisticRegression(multi_class='multinomial', random_state=42)
```

define my object

all happens in "fit" part

HYPER PARAMETERS

↑↑

You need to
choose them

Looking at the accuracy

- Obtain Accuracies on dataset to fit (train accuracy)

```
[25] print(  
    f"On the dataset used to fit, the decision tree has "  
    f"{100 * dt_clf.score(X_fit, y_fit):.1f}% accuracy"  
)  
print(    ↑ DT    ↑ gets accuracy    ↑ specific dataset  
    f"On the dataset used to fit, the random forest has "  
    f"{100 * rf_clf.score(X_fit, y_fit):.1f}% accuracy"  
)  
print(  
    f"On the dataset used to fit, the logistic regression has "  
    f"{100 * log_reg.score(X_fit, y_fit):.2f}% accuracy"  
)
```

→ On the dataset used to fit, the decision tree has 100.0% accuracy
On the dataset used to fit, the random forest has 99.8% accuracy
On the dataset used to fit, the logistic regression has 100.00% accuracy

Train
accuracy
(fit accuracy)

not reliable
metric of
your generalization
risk!!

- Obtain Validation Accuracies

```
[26] print(  
    f"On the dataset left aside for the fitting part, the decision tree has "  
    f"{100 * dt_clf.score(X_val, y_val):.1f}% accuracy"  
)  
print(    ← Look on the validation dataset!  
    f"On the dataset left aside for the fitting part, the random forest has "  
    f"{100 * rf_clf.score(X_val, y_val):.1f}% accuracy"  
)  
print(  
    f"On the dataset left aside for the fitting part, the logistic regression has "  
    f"{100 * log_reg.score(X_val, y_val):.2f}% accuracy"  
)
```

→ On the dataset left aside for the fitting part, the decision tree has 77.6% accuracy
On the dataset left aside for the fitting part, the random forest has 88.8% accuracy
On the dataset left aside for the fitting part, the logistic regression has 87.57% accuracy

→ 💯%. To estimate the generalization risk; I use

DT is worse
than RF
or Logistic R.

Wooclap time!

a few points (velodol^o act) that the model
has not seen during fit.



- 1
- 2

Allez sur wooclap.com

Entrez le code d'événement dans le bandeau supérieur

Code d'événement
TAMXAU



- 1
- 2

Envoyez **@TAMXAU** au
06 44 60 96 62

 Désactiver les réponses par SMS

Vous pouvez participer

What is true about SML

① Supervised Machine Learning does not have inputs

② Supervised Machine Learning can be combined with unsupervised in some contexts

③ Supervised Machine Learning does not have outputs

④ The goal is to predict well on new unseen points

⑤ It is necessary to perform well on the train (fit) set for that

⑥ The goal is to predict well on the train (fit) points

⑦ The balance interpretability/speed/ performance depends on the method

⑧ I assume my training data (fit + validation + test) to be representative (iid) of a phenomenon with distribution P

you typically
often necessary,
not sufficient

net only
and applied

! ! !

What is the role of fitting, e.g., in

rf_clf = RandomForestClassifier(n_estimators=10,random_state=4...)

rf-clf . fit (x_{fit} , y_{fit})

① Finding the best parameters of the Forest/Trees

rf (ω) ✓

② Finding the best hyperparameters (n_estimators)

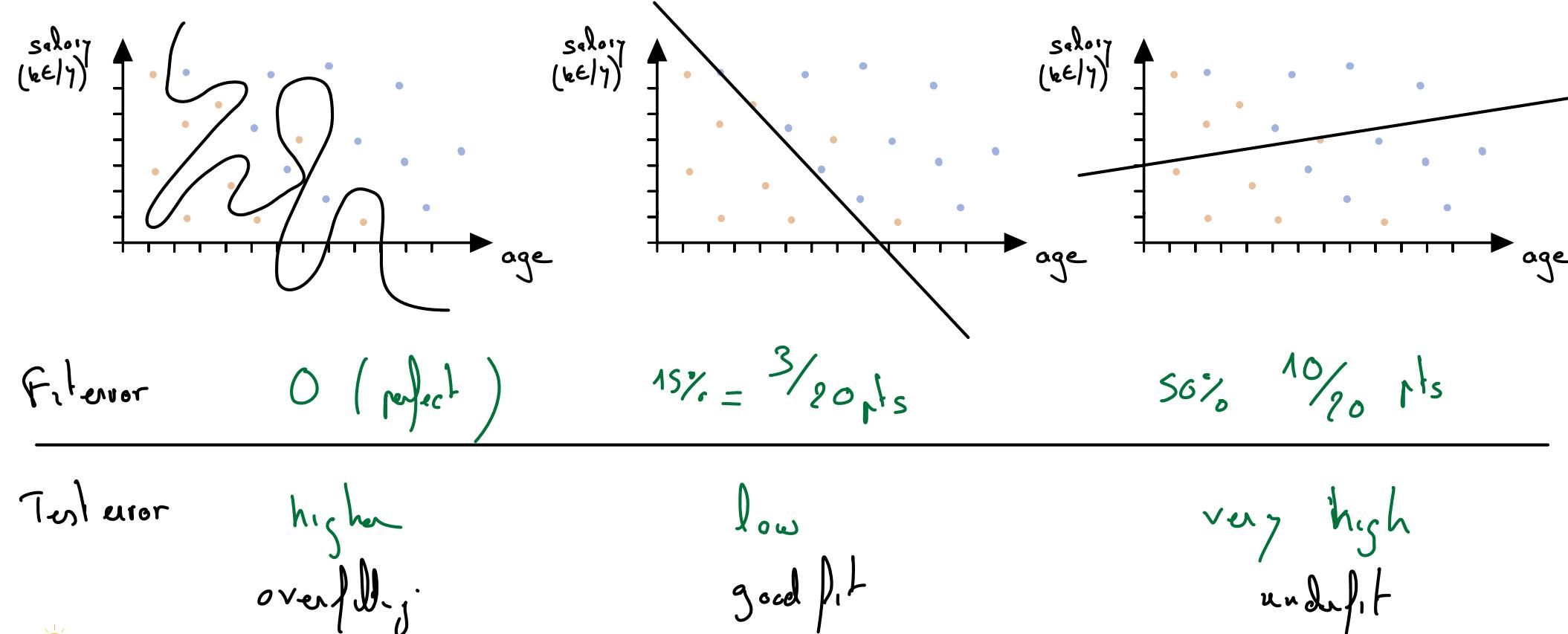
? ↗
use cross val part

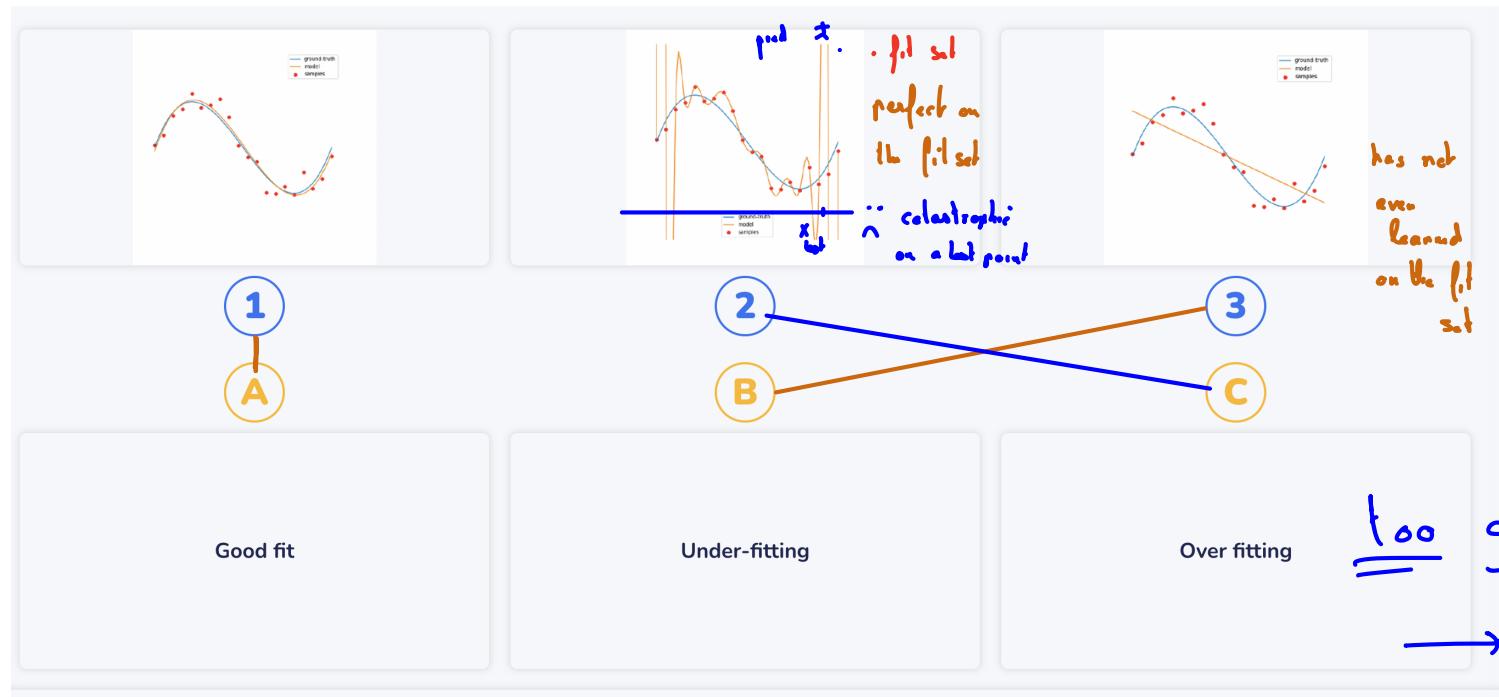
③ Predicting on a test point

rf-clf . predict (x_{test}) to predict

Caveat: Overfitting - Toy datasets

- ⚠ The performance on the training set may not reflect the ability to predict on a future point !!
- 💡 Our Goal is to perform well on unseen outputs!!





Caveat: Overfitting - Toy datasets

⚠ The performance on the training set may not reflect the ability to predict on a future point !!



Ou



Need to focus on validation/test performance

Underfitting and Overfitting depending on the choice of \mathcal{C}

$\mathcal{C} =$

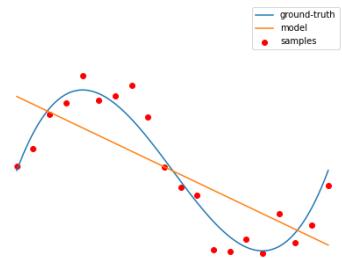
$\hat{R}_n(\theta_n^*)$

$R(\theta_n^*)$

Problem

Dim. of Θ

Underfitting and Overfitting depending on the choice of \mathcal{C}



$\mathcal{C} =$ Affine functions

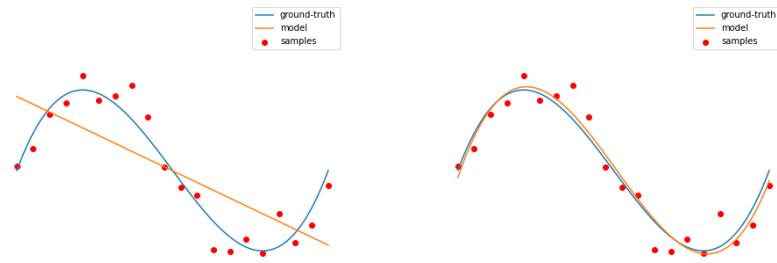
$\hat{R}_n(\theta_n^*)$ high

$R(\theta_n^*)$ high

Problem Underfitting

Dim. of Θ 2

Underfitting and Overfitting depending on the choice of \mathcal{C}



$\mathcal{C} =$	Affine functions	$\mathbb{R}_3[X]$
-----------------	------------------	-------------------

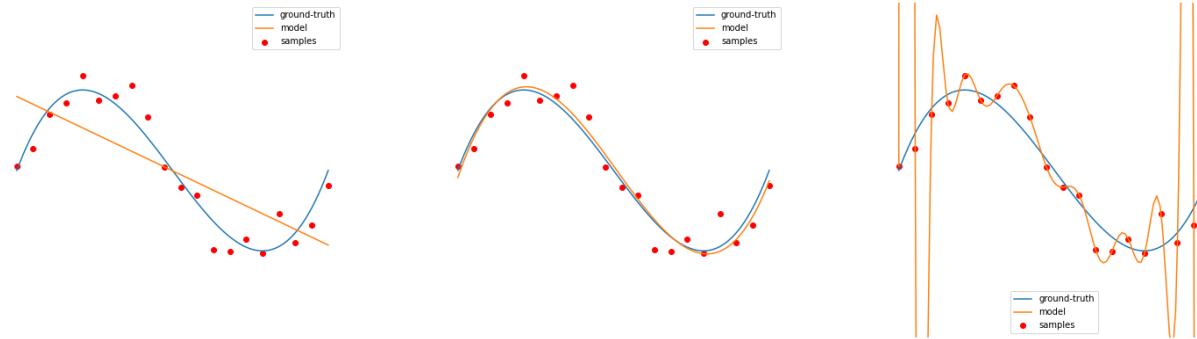
$\hat{R}_n(\theta_n^*)$	high	low
-------------------------	------	-----

$R(\theta_n^*)$	high	low
-----------------	------	-----

Problem	Underfitting	All good
---------	--------------	----------

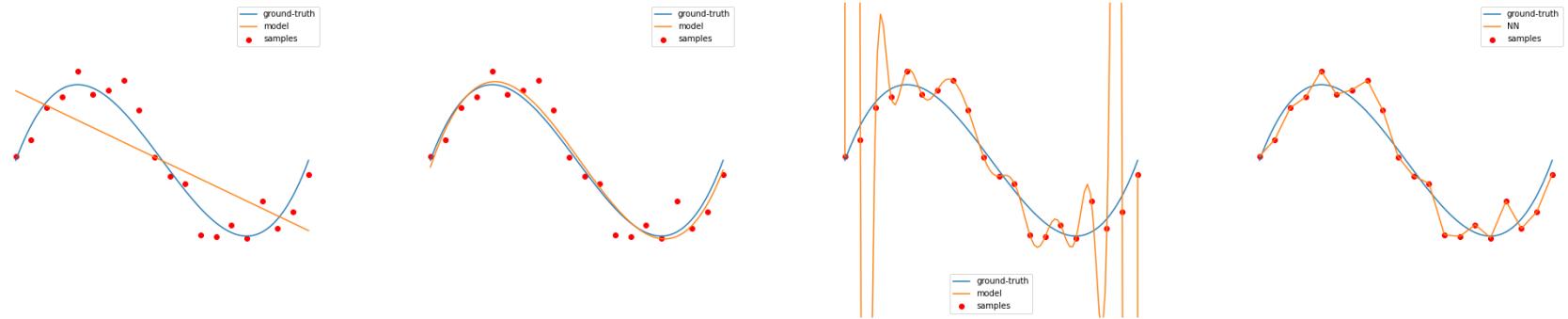
Dim. of Θ	2	4
------------------	---	---

Underfitting and Overfitting depending on the choice of \mathcal{C}



$\mathcal{C} =$	Affine functions	$\mathbb{R}_3[X]$	$\mathbb{R}_{19}[X]$
$\hat{R}_n(\theta_n^*)$	high	low	zero
$R(\theta_n^*)$	high	low	high
Problem	Underfitting	All good	Overfitting
Dim. of Θ	2	4	20

Underfitting and Overfitting depending on the choice of \mathcal{C}



$\mathcal{C} =$	Affine functions	$\mathbb{R}_3[X]$	$\mathbb{R}_{19}[X]$	FCNN
$\hat{R}_n(\theta_n^*)$	high	low	zero	zero
$R(\theta_n^*)$	high	low	high	mid low
Problem	Underfitting	All good	Overfitting	
Dim. of Θ	2	4	20	17375

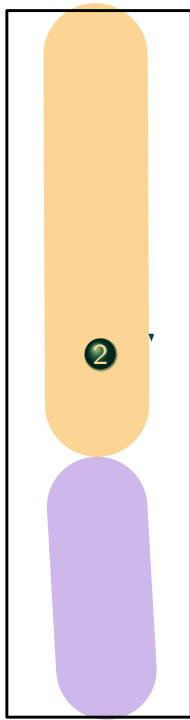
Demo Colab

Learning Pipeline

① Training phase

Training set
 $(X_{\text{train}}, Y_{\text{train}})$

T
R
A
I
N



1 model. how good is it

① Fitting Phase

F
I
T

Goal: learn the parameters
of the model on the
 $(X_{\text{fit}}, Y_{\text{fit}})$

② Validation phase

V
A
L
I
D

Possible models
& Hyperparameters

Find an
accurate predictio

Goal

of the generalizat

error by estimatio on validat

ish data

choose
the best
model

Test phase

T
E
S
T

New point.
 X_{new}



Predicted
output
 y_{new} .

③ Deployment

deploy the best model.

Fitting algo.

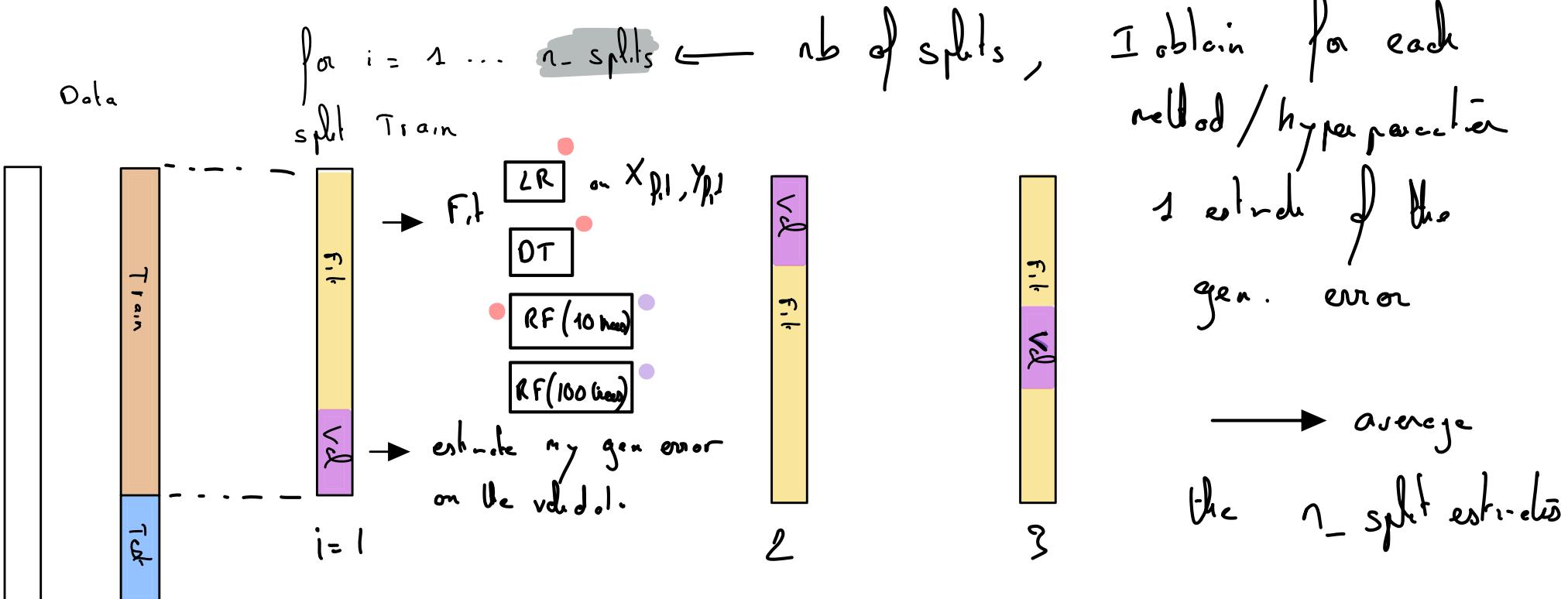
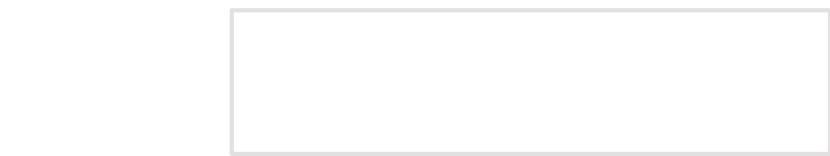
Learns funct $\hat{f}_n = \boxed{\quad}$
i.e. weight \downarrow
bias \uparrow

+ Cross validation

→ estimation
of the generalisat

Cross Validation

Goal: Select the best model / best hyperparam.



One algorithm to rule SML: Cross validation

💡 Enables to select the best method within multiple methods

```
Entrée [7]: 1 models = {  
2     "Decision Tree": DecisionTreeClassifier(random_state=42),  
3     "Random Forest": RandomForestClassifier(n_estimators=10,random_state=42),  
4     "Logistic Regression": LogisticRegression(multi_class="multinomial", solver="lbfgs", max_iter=100, random_st  
5     "MLP Classifier": MLPClassifier(random_state=42),  
6     "SVM": SVC()  
7 }
```

```
Entrée [12]: 1 cv_results = {}  
2  
3 for name, model in models.items():  
4     kfold = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)  
5     cv_score = cross_val_score(model, X_train, y_train, cv=kfold, scoring="accuracy")  
6     cv_results[name] = cv_score  
7     print(f"{name}: Mean CV accuracy={np.mean(cv_score)}, Std Deviation={np.std(cv_score)}")
```

```
Decision Tree: Mean CV accuracy=0.77125, Std Deviation=0.013163512495701512  
Random Forest: Mean CV accuracy=0.8826785714285714, Std Deviation=0.008957097288560329  
Logistic Regression: Mean CV accuracy=0.875, Std Deviation=0.007804204779790029  
MLP Classifier: Mean CV accuracy=0.8842857142857143, Std Deviation=0.0031237242293813998  
SVM: Mean CV accuracy=0.9505357142857143, Std Deviation=0.005463235193135154
```

```
Entrée [18]: 1 # Convert results to DataFrame for better visualization  
2 results_df = pd.DataFrame(cv_results)  
3 results_df
```

Out[18]:

	Decision Tree	Random Forest	Logistic Regression	MLP Classifier	SVM
0	0.775893	0.885714	0.880357	0.885714	0.953571
1	0.752679	0.875000	0.866071	0.882143	0.942857
2	0.776786	0.898214	0.883929	0.885714	0.958929
3	0.790179	0.881250	0.865179	0.888393	0.950000
4	0.760714	0.873214	0.879464	0.879464	0.947321
av:

← split 1 ← split 5

```
Entrée [17]: 1 # Descriptive statistics can also be helpful  
2 results_df.describe()
```

Out[17]:

	Decision Tree	Random Forest	Logistic Regression	MLP Classifier	SVM
count	5.000000	5.000000	5.000000	5.000000	5.000000
mean	0.771250	0.882679	0.875000	0.884286	0.950536
std	0.014717	0.010014	0.008725	0.003492	0.006108
min	0.752679	0.873214	0.865179	0.879464	0.942857
25%	0.760714	0.875000	0.866071	0.882143	0.947321
50%	0.775893	0.881250	0.879464	0.885714	0.950000
75%	0.776786	0.885714	0.880357	0.885714	0.953571
max	0.790179	0.898214	0.883929	0.888393	0.958929

Figure: Cross Validation

One algorithm to rule SML: Cross validation

💡 Enables to select the best hyperparameters for a given methods

Leveraging grid search CV

Entrée []:

```
1 from sklearn.model_selection import GridSearchCV
2
3 # Example for SVM
4 param_grid = {
5     'C': [1, 10, 100],
6     'gamma': [0.001, 0.01, 0.1],
7 }
8 grid_search = GridSearchCV(SVC(), param_grid, cv=3, verbose=2)
9 grid_search.fit(X_train_scaled[:10000], y_train[:10000]) # Reduced data size for quicker execution
10
11 print("Best parameters:", grid_search.best_params_)
12 print("Best cross-validation accuracy:", grid_search.best_score_)
13 best_model = grid_search.best_estimator_
14 final_predictions = best_model.predict(X_test_scaled)
15 print("Test set accuracy:", accuracy_score(y_test, final_predictions))
```

Logistic Regression

Entrée [43]:

```
1 pd.DataFrame(grid_search.cv_results_)
```

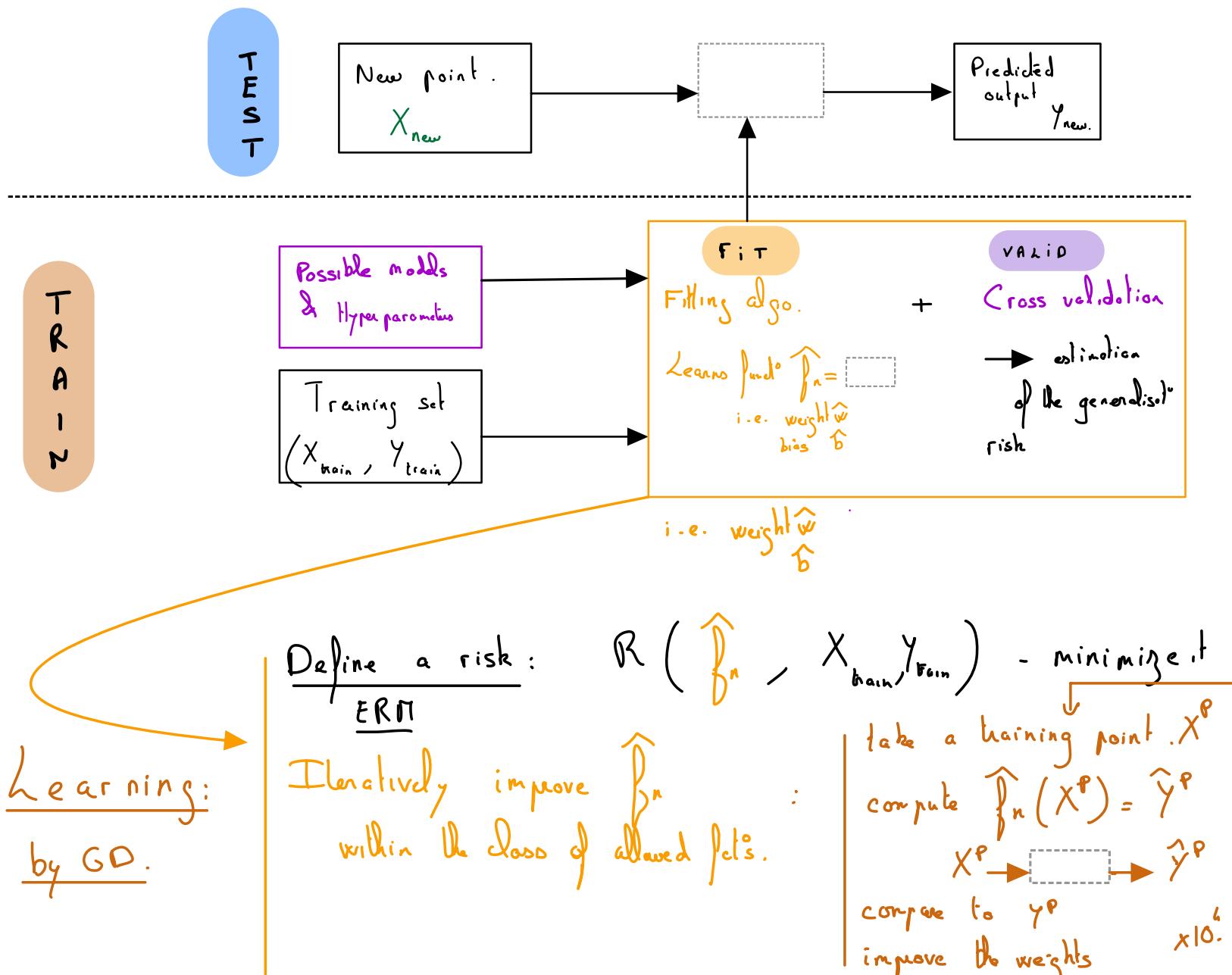
Out[43]:

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_C	param_gamma	params	split0_test_score	split1_test_score	split2_test_score	mean_score
0	2.354261	0.049712	1.981289	0.027614	1	0.001	{'C': 1, 'gamma': 0.001}	0.918050	0.930370	0.927653	0.927653
1	7.614626	0.022478	2.995119	0.006781	1	0.01	{'C': 1, 'gamma': 0.01}	0.693626	0.718800	0.724009	0.724009
2	8.525636	0.034385	3.381043	0.004120	1	0.1	{'C': 1, 'gamma': 0.1}	0.175147	0.187467	0.181136	0.181136
3	1.957532	0.022052	1.767812	0.001655	10	0.001	{'C': 10, 'gamma': 0.001}	0.928227	0.933048	0.937835	0.937835
4	7.731749	0.122276	3.002695	0.004842	10	0.01	{'C': 10, 'gamma': 0.01}	0.708623	0.735404	0.746517	0.746517
5	8.540356	0.023448	3.378721	0.004322	10	0.1	{'C': 10, 'gamma': 0.1}	0.176219	0.187467	0.184887	0.184887
6	1.959067	0.025085	1.774580	0.014967	100	0.001	{'C': 100, 'gamma': 0.001}	0.928763	0.933048	0.937835	0.937835
7	7.861557	0.143017	3.017998	0.025798	100	0.01	{'C': 100, 'gamma': 0.01}	0.708623	0.735404	0.746517	0.746517
8	8.649604	0.021181	3.455342	0.009529	100	0.1	{'C': 100, 'gamma': 0.1}	0.176219	0.187467	0.184887	0.184887

best hyperparameters!

Figure: Grid search CV

Supervised Machine Learning pipeline



SML -Summary

① Key concepts

- ▶ Tasks
- ▶ Methods
- ▶ Overfitting and cross validation

SML -Summary

① Key concepts

- ▶ Tasks
- ▶ Methods
- ▶ Overfitting and cr

② Black box solution:

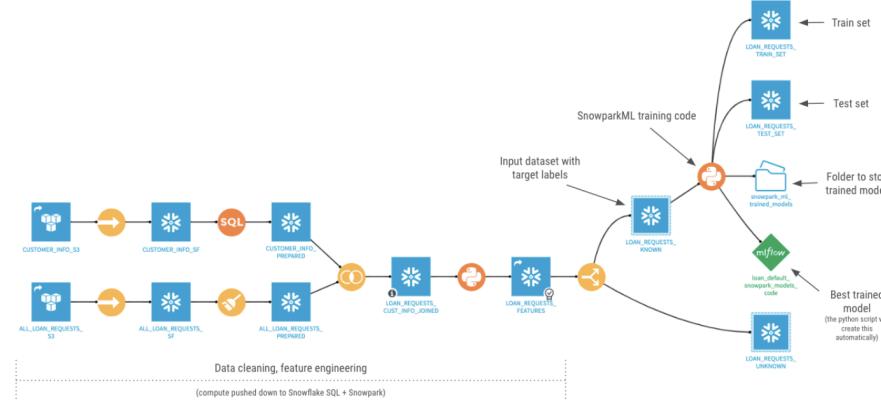


Figure: Dataiku

<https://blog.dataiku.com/training-ml-models-with-dataiku-and-snowpark-ml-a-code-approach>

SML -Summary

① Key concepts

- ▶ Tasks
- ▶ Methods
- ▶ Overfitting and cr

② Black box solution:

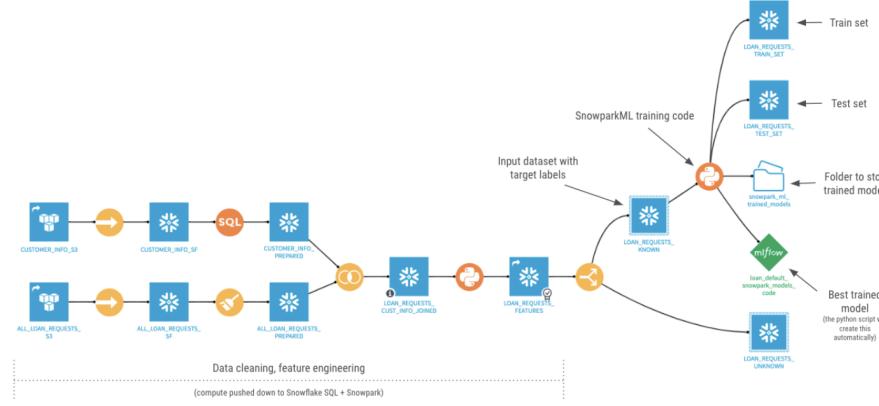


Figure: Dataiku

<https://blog.dataiku.com/training-ml-models-with-dataiku-and-snowpark-ml-a-code-approach>

③ Our goal for today and tomorrow: **Opening the black box**

- Understanding methods enables to better anticipate the best choice, limits, potential, etc.
- Which method is suitable depending on Performance/Speed/Interpretability trade-offs
- Better choices for hyper-parameters to be chosen

Wooclap time!



Code d'événement
TAMXAU

- 1 Allez sur wooclap.com
- 2 Entrez le code d'événement dans le bandeau supérieur

- 1 Envoyez **@TAMXAU** au
06 44 60 96 62
- 2 Vous pouvez participer

Désactiver les réponses par SMS

What is the problem of overfitting

- ① I will have a bad error on the train set

X

- ② I will have a bad error on the test set

✓

- ③ My predictor may change a lot if I change the training set

✓ instability

What is the role of cross validation

- ① Find the best model in terms of validation error

✓

- ② Find the best parameters in terms of training error

→ fit

at the end of
CV, you get a

- ③ Find the best hyperparameters in terms of validation error

✓

fitted model for
the best model + hyperpara
that is the best in
terms of validation error

- ④ Fit the model

- ⑤ Find the best parameters in terms of validation error

→ ? X

```
Entrée [7]: 1 models = {
2     "Decision Tree": DecisionTreeClassifier(random_state=42),
3     "Random Forest": RandomForestClassifier(n_estimators=10, random_state=42),
4     "Logistic Regression": LogisticRegression(multi_class="multinomial", solver="lbfgs", max_iter=100, random_st
5     "MLP Classifier": MLPClassifier(random_state=42),
6     "SVM": SVC()
7 }
```

```
Entrée [12]: 1 cv_results = {}
2
3 for name, model in models.items():
4     kfold = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
5     cv_score = cross_val_score(model, X_train, y_train, cv=kfold, scoring="accuracy")
6     cv_results[name] = cv_score
7 print(f"\n{name}: Mean CV accuracy={np.mean(cv_score)}, Std Deviation={np.std(cv_score)}")
```

Decision Tree: Mean CV accuracy=0.77125, Std Deviation=0.013163512495701512
 Random Forest: Mean CV accuracy=0.8826785714285714, Std Deviation=0.008957097288560329
 Logistic Regression: Mean CV accuracy=0.875, Std Deviation=0.007804204779790029
 MLP Classifier: Mean CV accuracy=0.8842857142857143, Std Deviation=0.0031237242293813998
 SVM: Mean CV accuracy=0.9505357142857143 Std Deviation=0.005463235193135154

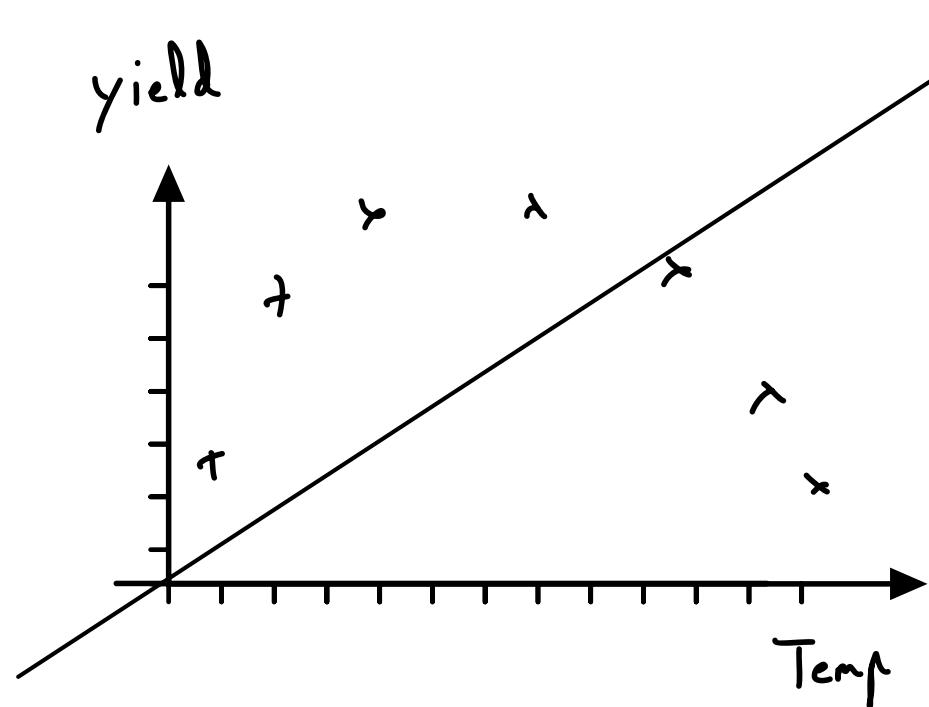
```
Entrée [18]: 1 # Convert results to DataFrame for better visualization
2 results_df = pd.DataFrame(cv_results)
3 results_df
```

	Decision Tree	Random Forest	Logistic Regression	MLP Classifier	SVM
0	0.775893	0.885714	0.880357	0.885714	0.953571 ✓
1	0.752679	0.875000	0.866071	0.882143	0.942857 •
2	0.776786	0.898214	0.883929	0.885714	0.958929 •
3	0.790179	0.881250	0.865179	0.888393	0.950000 •
4	0.780714	0.873214	0.879464	0.879464	0.947321 •

What does this code do

- 1 It enables to compare the five methods given above ✓
- 2 It computes Train score for each method ✓
- 3 It computes Validation score for each method ✓
- 4 It computes 5 scores for each method, splitting the test data ✓
- 5 It computes 5 scores for each method, splitting the train data ✓
- 6 The best method is SVM ✓
- 7 About 95% of the digits are well classified on the calibration set : validat. at least for SVM yes
- 8 Random forest outperform Decision Trees ✓
- 9 25 fits were performed overall ✓
- 10 5 fits were performed overall ✓

Limits of Linear and logistic regression - towards DT



Linear regression : is bad here.

- * The relationship between X (Temp) and Y is non linear
- * if I have House price dataset with 81 features , I
get a model w. 81 coefficients

DT

learn non linear
relationships !
strong interpretability

Outline

1 Introduction to Supervised Learning

2 Decision Trees

3 Random Forests

Decision Tree

Goals of Decision Trees

- ① Supervised Learning: solve the classification or regression task
- ② Provide some understanding: what on this example allows me to classify it? What are the important features ?
- ③ Go beyond linear functions → allow complex dependencies.

Decision Tree

Goals of Decision Trees

- ① Supervised Learning: solve the classification or regression task
- ② Provide some understanding: what on this example allows me to classify it? What are the important features ?
- ③ Go beyond linear functions → allow complex dependencies.

Two questions:

- How to train a decision tree? → **training phase (learning algorithm)**
- How to predict with a given decision tree? → **inference phase**

~~• fit
• predict~~

NN. → $\begin{cases} \text{• fit} = \text{incremental updates (Stoch gradient descent)} \\ \text{• predict} \end{cases}$ → what we have to understand tomorrow

Inference Phase

being able to predict from the tree.

$$X = (\text{age}, \text{precnd}, \text{medical doctor})$$

IR
 $\{0,1\}$
 $\{0,1\}$

$$Y = \{ \begin{array}{l} \text{got covid shot or not?} \end{array} \}$$

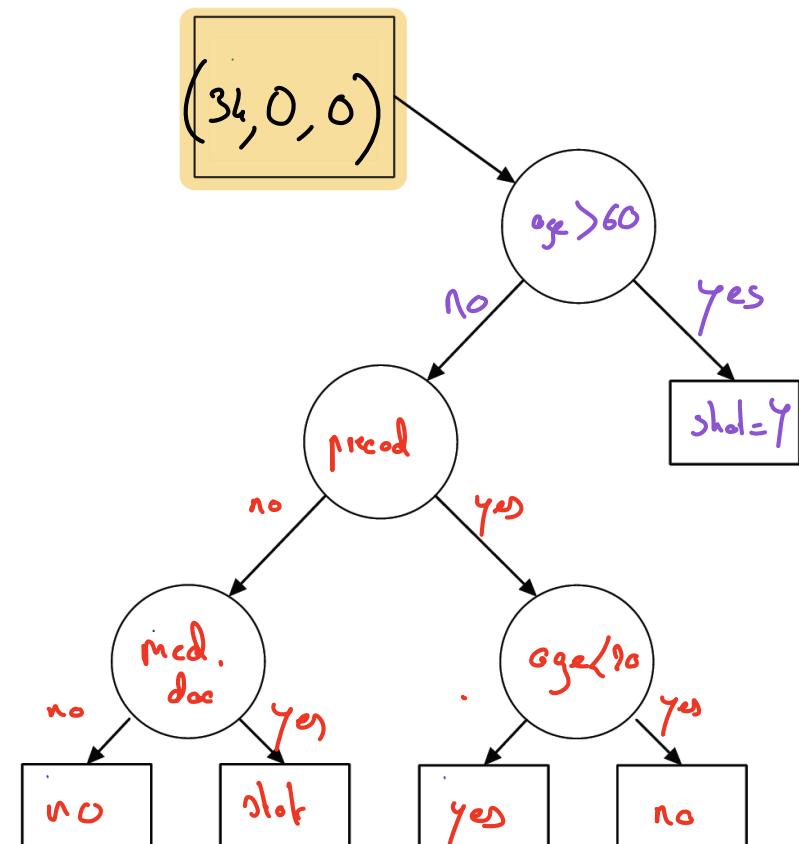
- Root = observation = input

$$(16, \text{yes}, \text{no})$$

- Nodes = tests

- Branches = possible outcomes

- Leaves = final decision = output



What do we have to learn →

list of questions
 each question = { Positive, Threshold }

Training: how to build a tree.

Key question !

- What is a good tree?
- What do I need to choose?
- How to choose?

What do you think?

1st point of view

one that performs well on $(X_{f,i}, Y_{f,i})$

choose

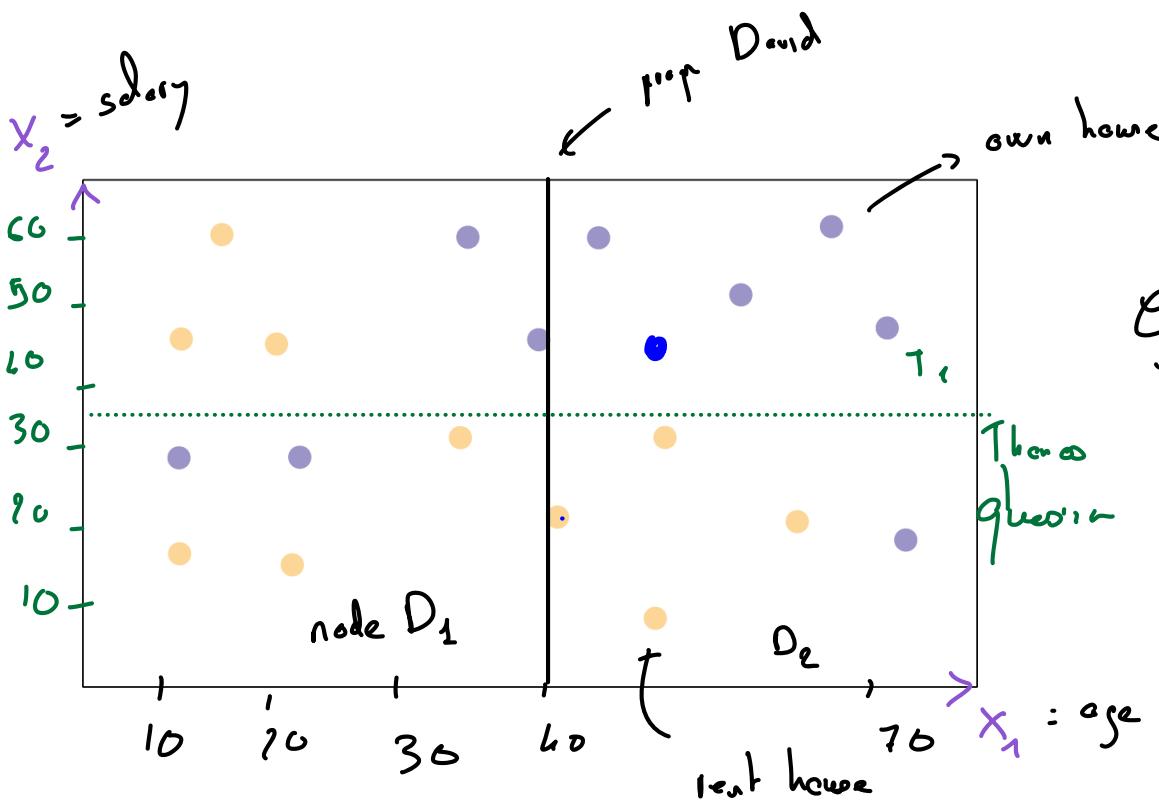
min $\left(\text{empirical risk } (T) \right)$

$T \in \text{Trees}$

→ lasted we incrementally build the tree!

one question at each increased step

Building trees : CART



homogeneity within the splits

Gini's impurity: within each node ; $P_i(1 - P_i)$

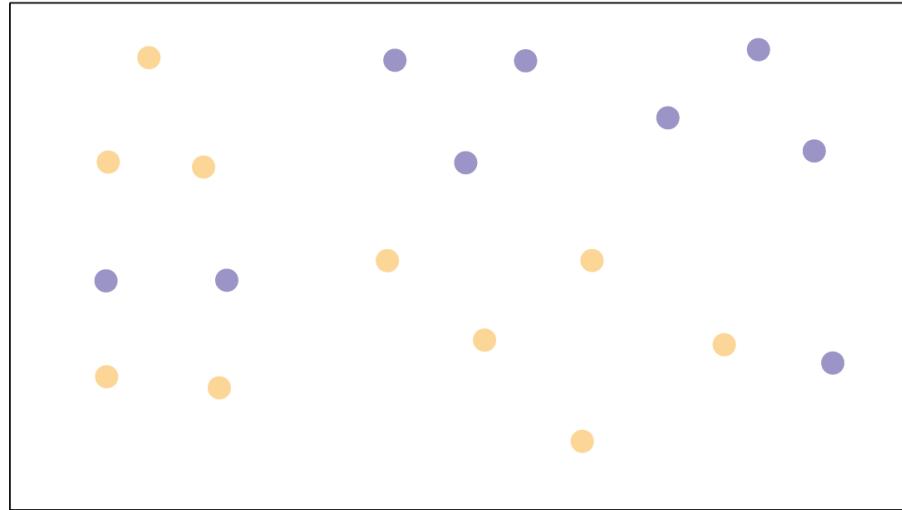
$\rightarrow P_i = \text{probability of color}$

Select minimal impurity
or

	David's ?		Thao ?	
D_1	6 •	6 •	T_1	3 • 7 •
D_2	4 •	6 •	T_2	8 • 8 •
$P_{uni} \rightarrow D_1$	0.4×0.6	$= 24\%$	0.3×0.7	$\approx 81\%$
$\rightarrow D_2$	0.4×0.6	$= 24\%$	0.8×0.2	$= 16\%$
		<hr/>		<hr/>
		24%		16%
		<hr/>		<hr/>
		48%		37%

Key : criterion of quality ↑

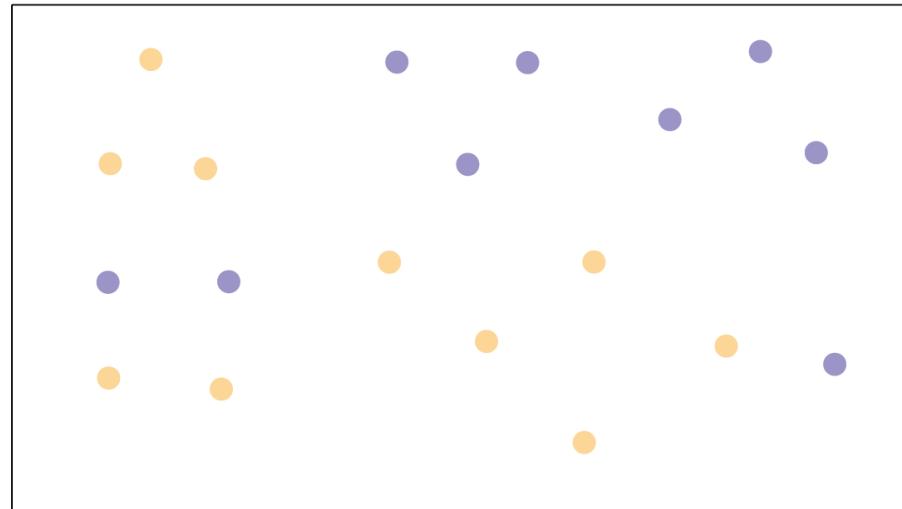
Building trees : CART



How to chose the best cut ?

- e.g., Given a possible cut, measure the reduction of « impurity »:
 - ▶ in regression: mean-squared-error $\sum_{i \in C} (Y_i - \bar{Y}_C)^2$
 - ▶ in classification: Gini's impurity.
- Find dimension and threshold with optimal impurity reduction.
- Iterate until stopping criterion is met.

Gini's impurity



Gini impurity measures **how often a randomly chosen element from the set would be incorrectly labeled if it was randomly labeled according to the distribution of labels in the subset.**

$$I_G(p) = \sum_{i=1}^J p_i(1 - p_i) = 1 - \sum_{i=1}^J p_i^2$$

For two classes, related to the variance if classified as Bernoulli r.v..

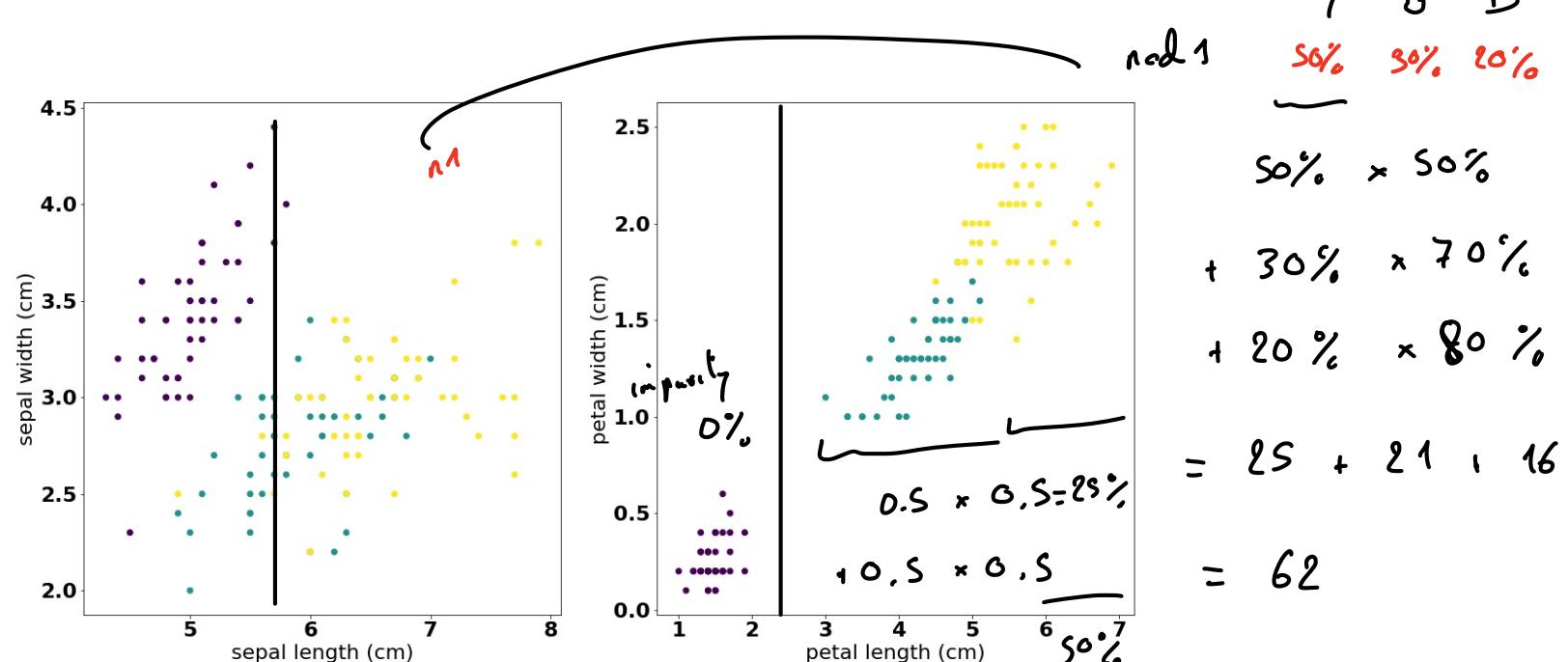
Example: Iris Dataset



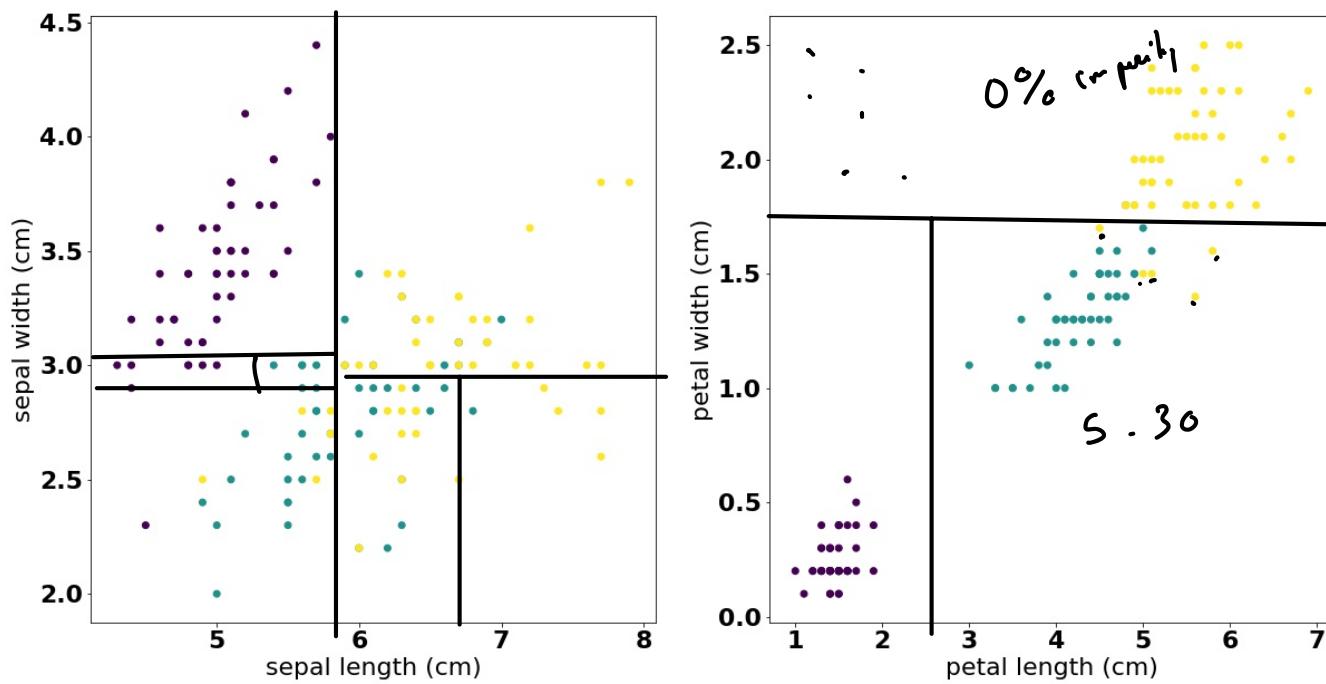
Figure: Iris: setosa, versicolor, virginica

One of the most widely used toy dataset in classification:

- 3 classes : ['setosa', 'versicolor', 'virginica']
- 50 points per class.
- 4 features: ['sepal length', 'sepal width', 'petal length', 'petal width']



Exercise: Guess Classification tree for Iris Dataset



In practice : Scikit Learn

Universal Python library for Machine Learning:

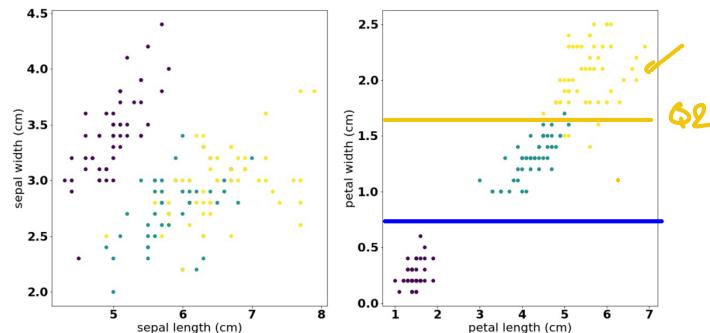
- Very easy to use
- Very well documented
- Open Source

<https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html>

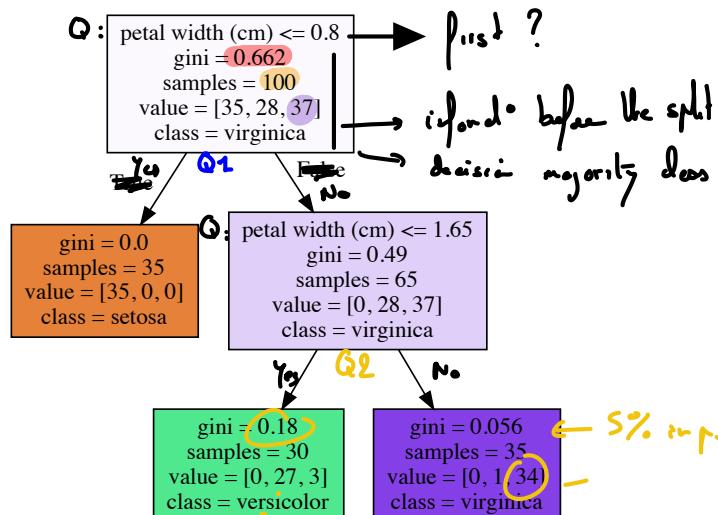
Example: Iris Dataset

What do we need to specify ?

150 pts
lit
100
vd
→ 50



Output:



```
from sklearn.datasets import load_iris
from sklearn.model_selection import cross_val_score
from sklearn.tree import DecisionTreeClassifier, plot_tree
from sklearn.model_selection import train_test_split
from sklearn.model_selection import train_test_split
import matplotlib.pyplot as plt
```

```
# Choose the Learning algorithm
clf = DecisionTreeClassifier(max_depth = 2)
```

max ab of qst^o

```
#Load the dataset
```

```
iris = load_iris()
X_ftr, X_vld, y_ftr, y_vld = train_test_split(
    iris.data, iris.target, test_size=0.33, random_state=42)
```

hyperparam :

```
# Fit the model
clf.fit(X_ftr, y_ftr)
```

```
# Plot tree
```

```
plt.figure(figsize=(5,4))
plot_tree(clf, feature_names = iris.feature_names,
          class_names=iris.target_names,
          filled = True)
plt.show()
```

Example:

```
X_test[0,:], iris.target_names[y_test[0]]
```

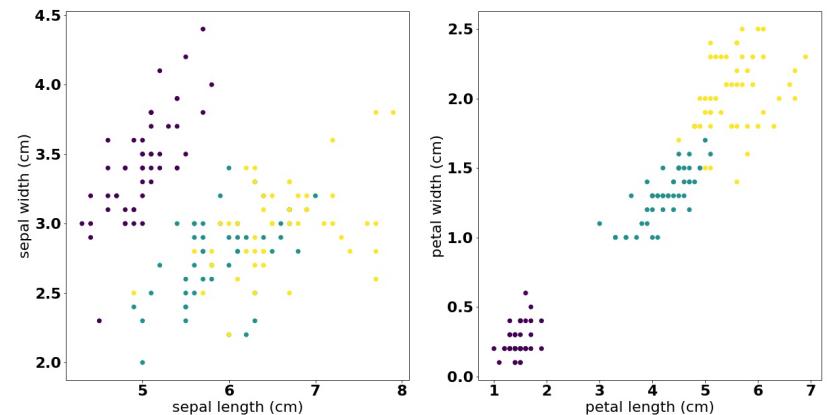
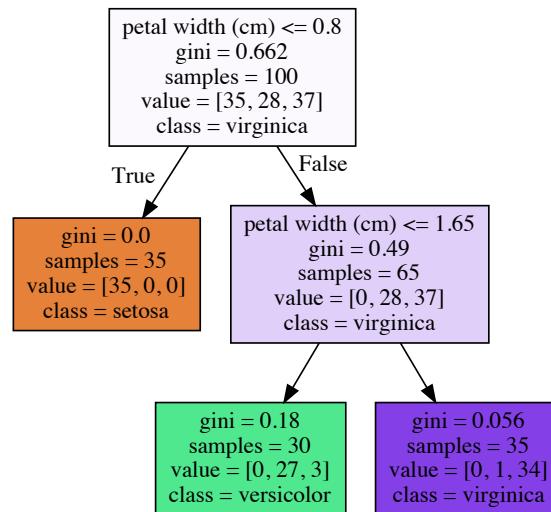
Returns:

```
(array([4.8, 3.1, 1.6, 0.2]), 'setosa')
```

Code: <https://colab.research.google.com/drive/1dXYIjacNAeASgl54vIWjWjMQlfgigSVA?usp=sharing#scrollTo=VeJHkJd7RsDw>

Example: Iris Dataset

Output:



Example:

```
X_test[0,:], iris.target_names[y_test[0]]
```

Returns:

```
(array([4.8, 3.1, 1.6, 0.2]), 'setosa')
```

Code: <https://colab.research.google.com/drive/1dXYIjacNAeASgl54vIWyWjMQlfgigSVA?usp=sharing#scrollTo=VeJHkJd7RsDw>

What if I use a deeper tree ? Or if I do not specify the depth in the model definition ?

Wooclap time!



Code d'événement
TAMXAU

- 1 Allez sur wooclap.com
- 2 Entrez le code d'événement dans le bandeau supérieur

- 1 Envoyez **@TAMXAU** au
06 44 60 96 62
- 2 Vous pouvez participer

Désactiver les réponses par SMS

1

The depth is "learned" in DT : *nope its a hyperparameter*

2

The depth is a hyperparameter

✓ *(yes!)*

3

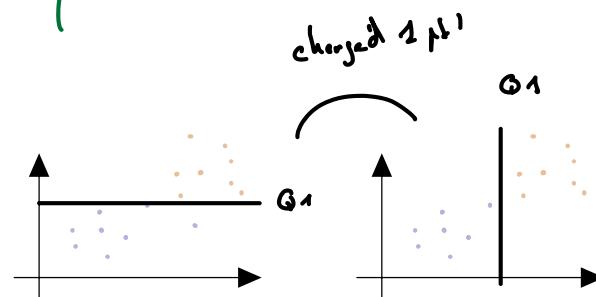
I use the validation set to find hyperparameters

always : not just for DT

4

There is no train set for Decision Trees

yes :



Which of the following statements are true about Decision Trees?

1

They are very stable

2

They provide an interpretable approach

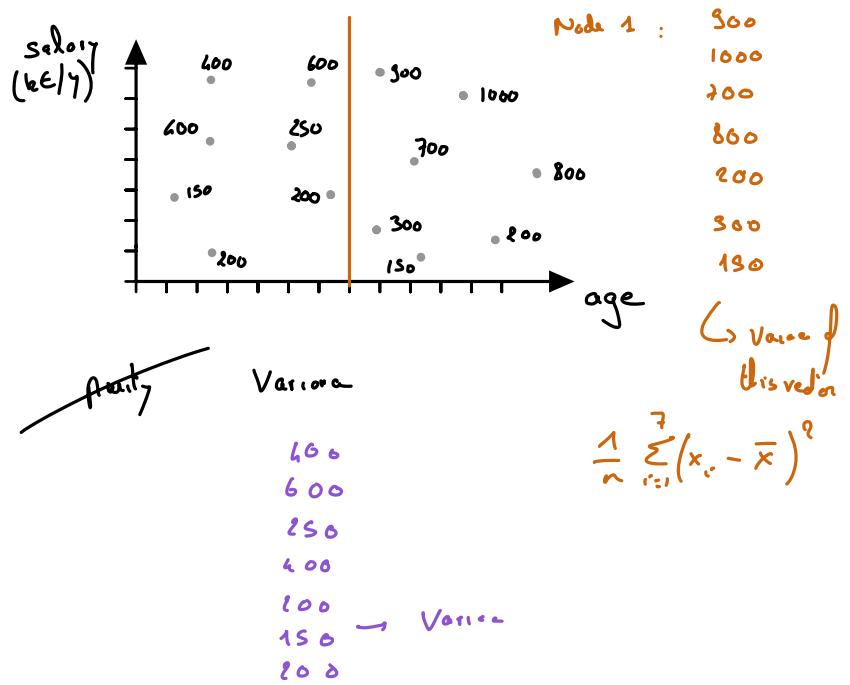
3

They are able to approach the "optimal function" for nearly any data distribution (e.g. beyond linearly separable data)

*4 : Work well for high n ✓
→ complexity per split is proportional*

5 : Do not fit regress°. X

yes they do



Which of the following statements are true about the construction of a decision tree?

1

The same feature cannot be used twice consecutively"



X

you can
but w. ≠ thresholds.

2

The same feature cannot be used twice overall"



X

3

I can use a pair of features if it gives me a better split

→ X

4

I want to find the split that has minimal impurity

5

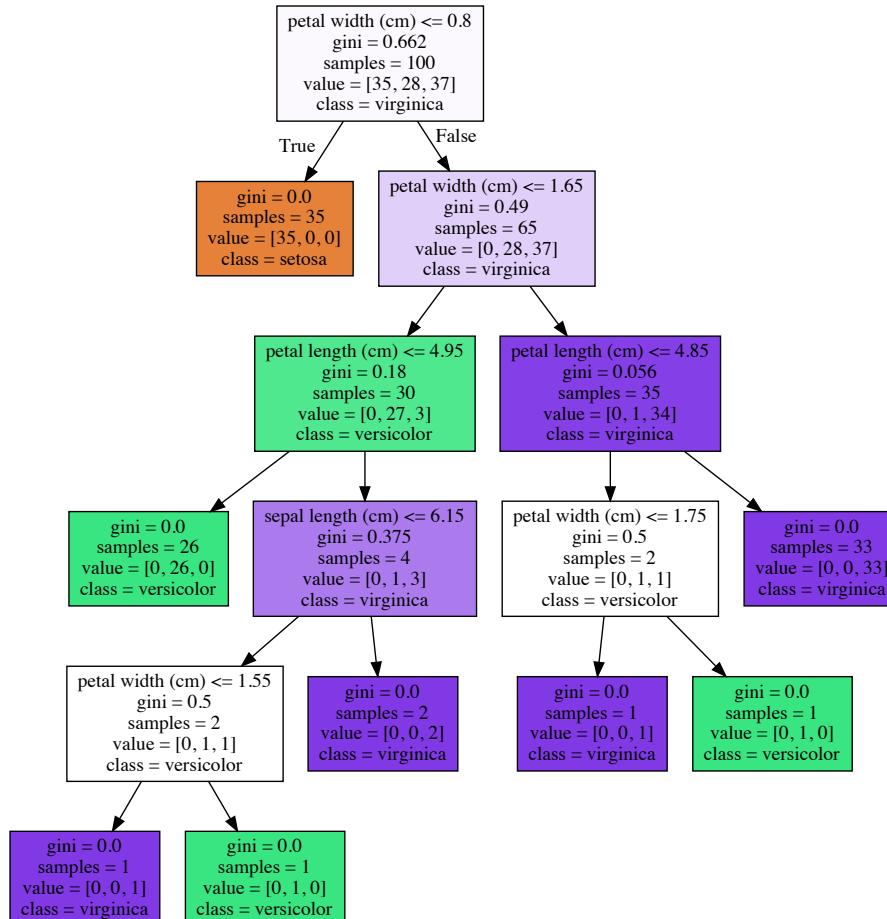
To define a new split, I only use one example

observed

all obs
1 feature

Deeper trees

```
# Choose the Learning algorithm  
clf = DecisionTreeClassifier(random_state=0)
```

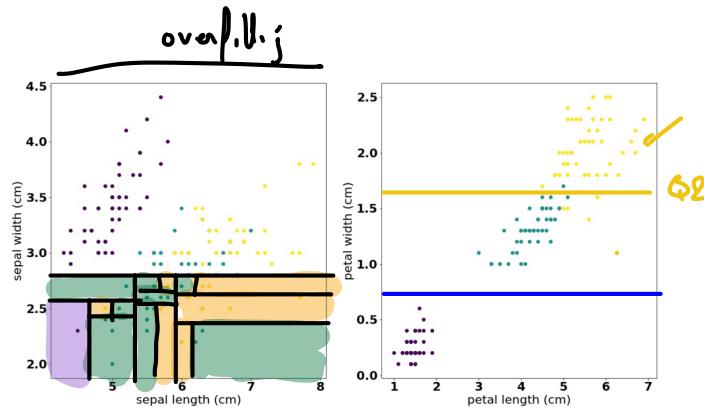


- When does the algorithm stop then ?
- What do you think about it ?

→ overfit

Stopping / Pruning

Trees that achieve perfect classification may suffer from over-fitting. (see example in the lab)



To avoid this problem, we can reduce the complexity of the trees:

① Stopping rules

- ▶ Set a minimum number of samples inside each leaf.
- ▶ Set a maximum depth.

② Pruning

- ▶ Reduced error pruning.
- ▶ Cost complexity pruning.

hyperparameters.
avoid overfitting

Pruning : example of Cost complexity pruning.

We create a sequence of Trees by changing one subtree at each step into a leaf, until we get only the root: the subtree is chosen to minimize the error made.¹ Then the test error is evaluated along the sequence, to choose the optimal pruning.



¹ more details here

Cart: pros and cons

Classif
and Reg

Pros:

- Simple to understand, interpret, visualize
- Implicitly perform features/variable selection
- can handle both categorical and numerical data
- little effort for data preparation
- Non linear relationships do not affect performance
- Can work with large number of observations.

Trees.

+

Cons:

- ① Can create over-complex trees: overfitting
- ② Unstability: small variations of data can generate completely different trees
- ③ Cannot guarantee global optimal tree
- ④ Biased if some classes dominate



Solution : Using Random Forests !

Outline

1 Introduction to Supervised Learning

2 Decision Trees

3 Random Forests

Random Forests

To put it in simple words: « random forests builds multiple decision trees and merges them together to get a more accurate and stable predictions ».

→ key idea: create variable trees + aggregate them

How to create variability ?

ideas: include | errors
 | uncertainty

→ change dataset
Sol 1 → subsampling → bootstrapping
→ check how we cut.

→ sampling the questions I can check.

Sol 2: only "try" a certain nb of features
for each question. e.g., for Q1, only x_1, x_n
can be used

Speed

Interpret

Accuracy

DT is superior

DT - -

RF is superior

Random Forests

To put it in simple words: « random forests builds multiple decision trees and merges them together to get a more accurate and stable predictions ».
→ key idea: create variable trees + aggregate them

How to create variability ?

max_features !

- Instead of the most important features, search the best feature among a random subset (ntry). Sd l
- Work on subsets of data (bootstrap=True). → Sd 1
- More <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestRegressor.html>

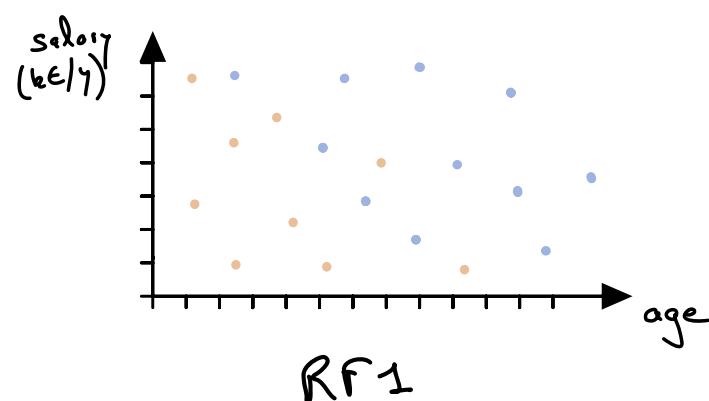
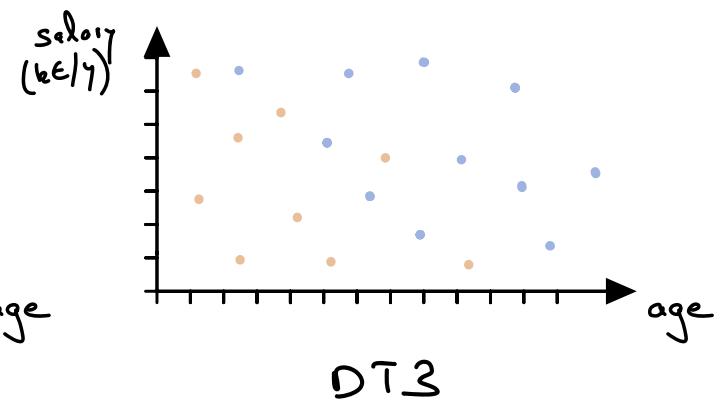
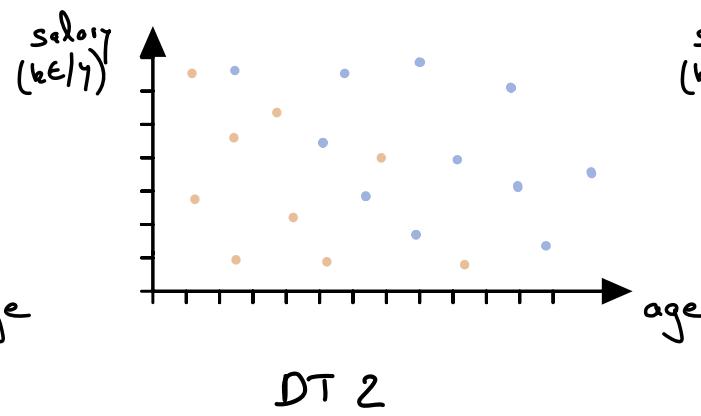
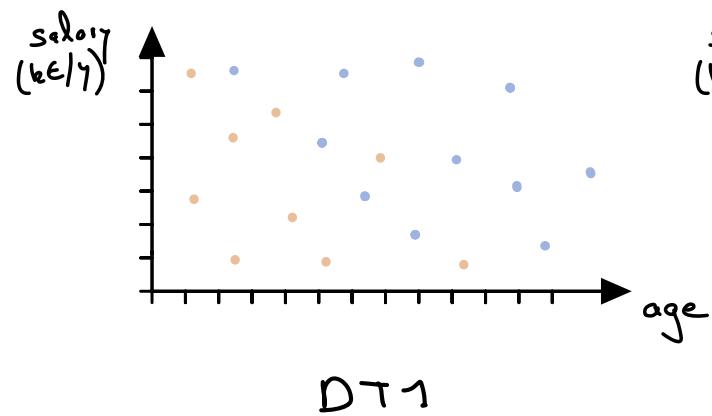
```
# Creating and fitting a Random Forest
from sklearn.ensemble import RandomForestClassifier

plt.figure(figsize=(20,10))
clf = RandomForestClassifier(n_estimators = 100,           ↴ number of
                             max_depth=3, bootstrap = True,   ↴ trees .
                             random_state =43)             ↴ for each tree, the max. depth
clf.fit(X_train,y_train)
```

How to aggregate ? How to interpret a Random Forest ?

Aggregation

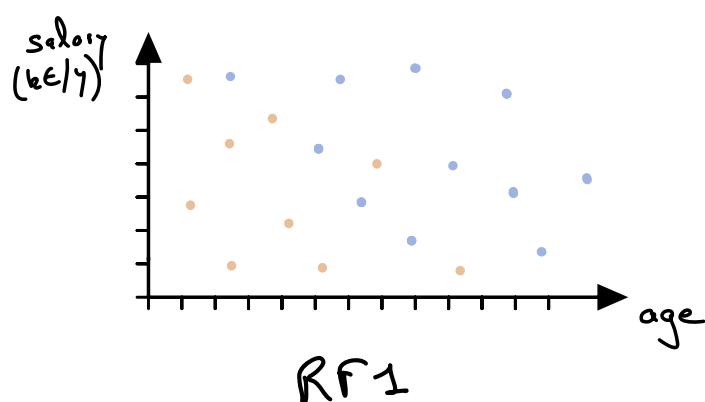
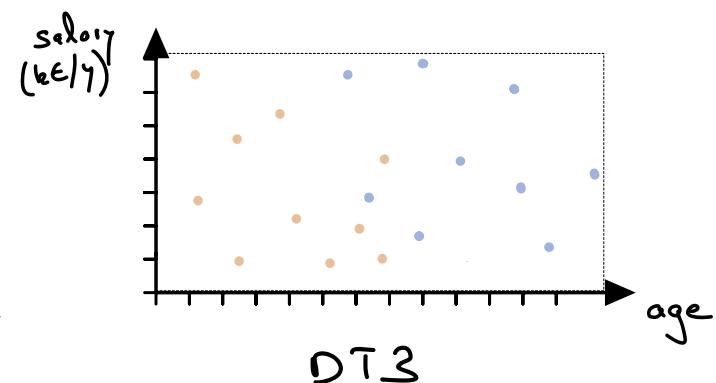
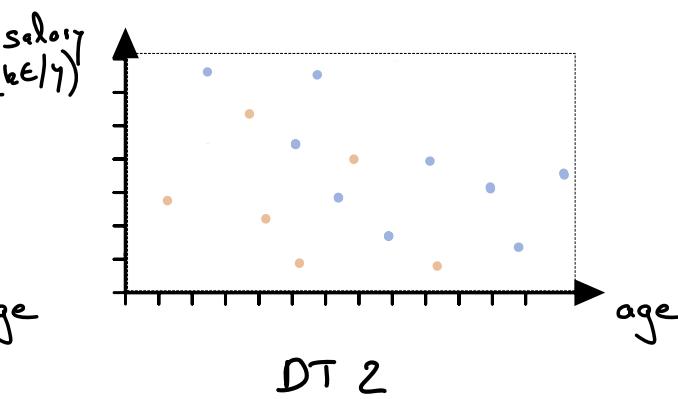
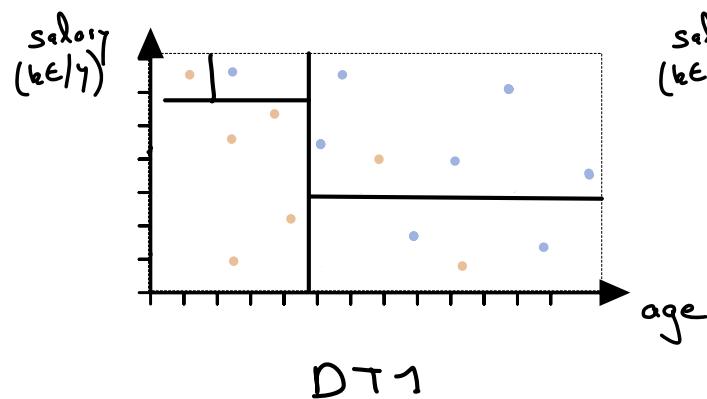
We using a voting or averaging process !



Aggregation

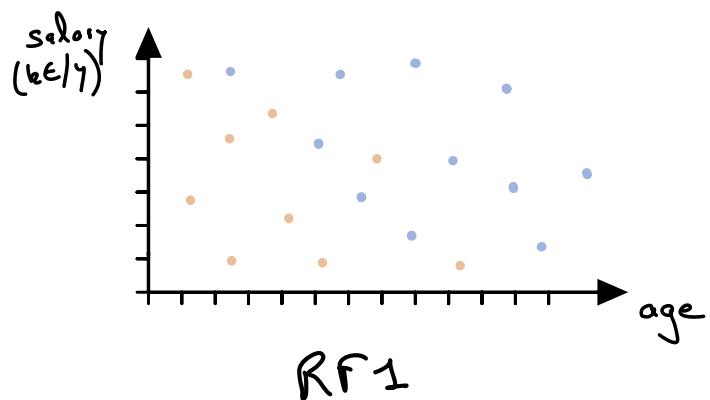
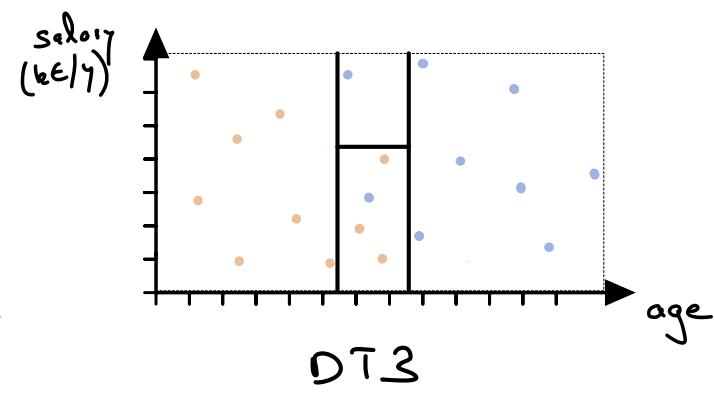
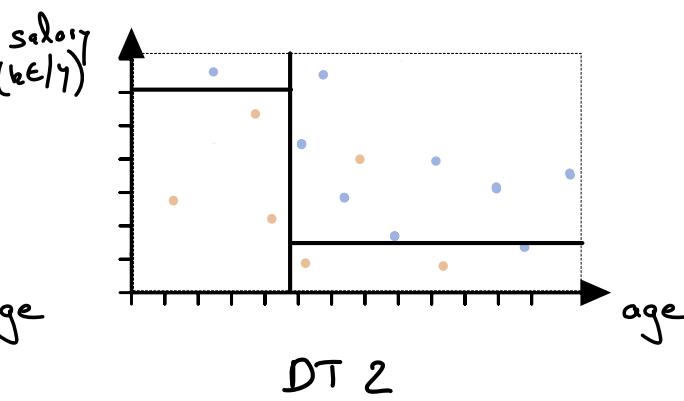
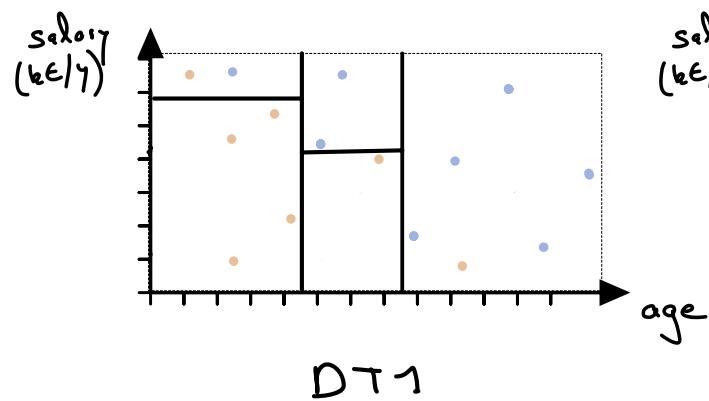
on each, keep 50% of the points

We using a voting or averaging process !



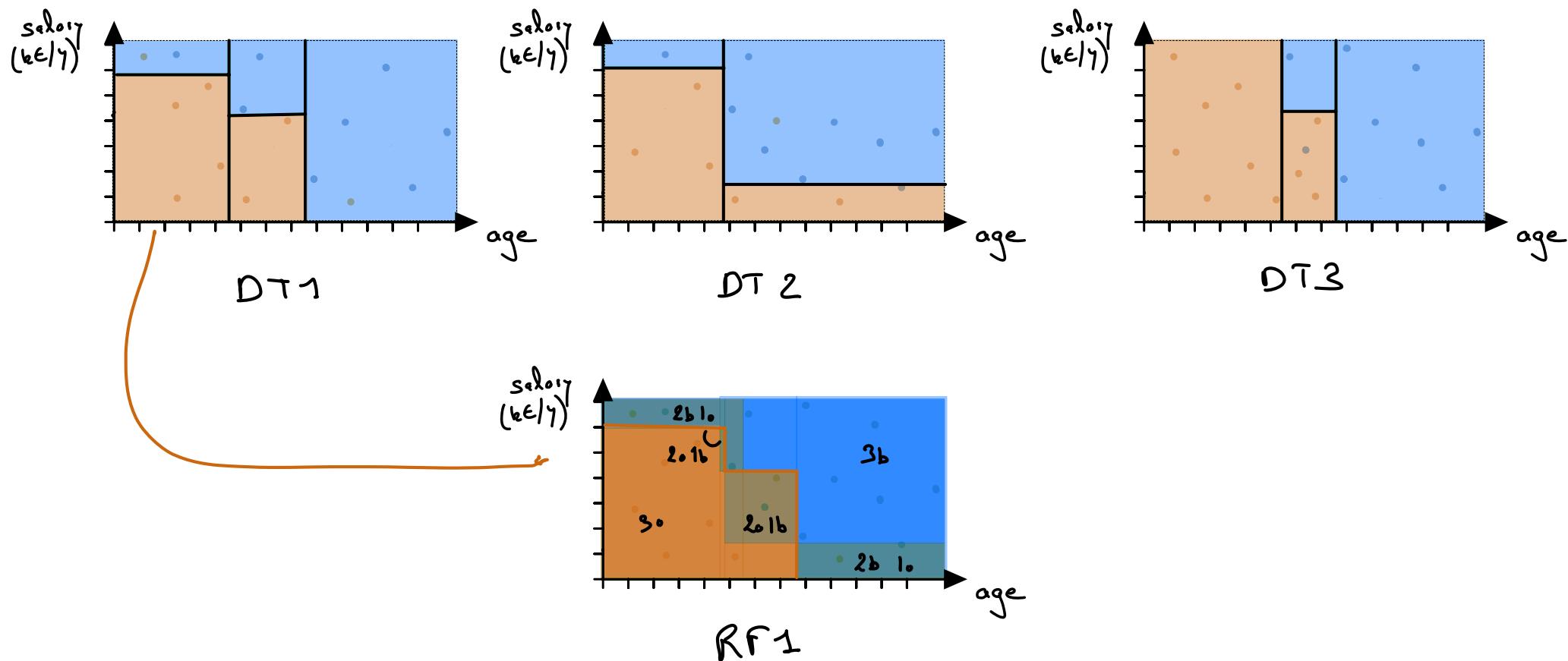
Aggregation

We using a voting or averaging process !



Aggregation

We using a voting or averaging process !



Feature importance

- how much tree nodes using a particular feature reduce impurity along the trees.
- suggests feature selection rule: drop out features with low importance to avoid overfitting.
- Attribute **feature_importance_** in RandomForestRegressor of Sklearn.

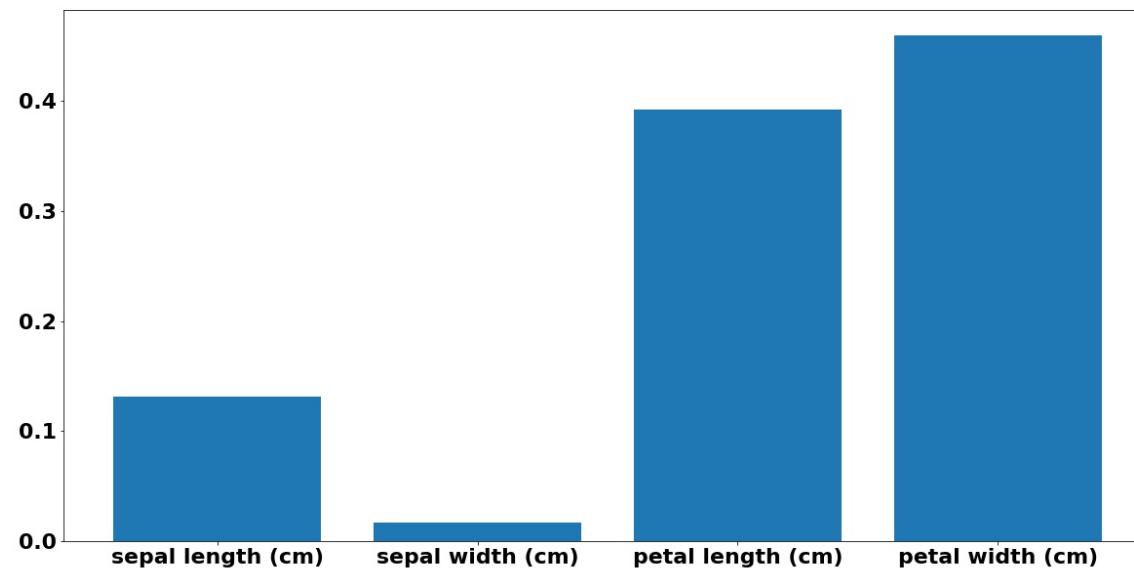


Figure: Feature importance output for a Random Forest Classifier in Python

Comparison with decision trees

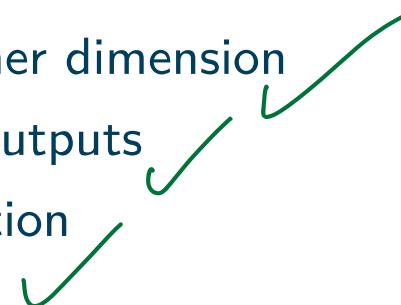
Cons:

- Less interpretable
- Slower to run



Pros:

- Better in higher dimension
- More stable outputs
- Feature selection



Improving predictions - more arguments.

① **n_estimators** : number of trees in the forest

- ▶ improves prediction / stability ..
- ▶ slows down the algorithm

as much as you can

② **max_features** : maximum number of features considered to split a node

③ **min_sample_leaf** : minimum number of samples that should remain in a leaf node.

RF pros and cons

gail Veron

RF Pros:

Why do tree based models outperform DL on tabular data

2022 at. 1690

- ① works for classification and regression
- ② default hyperparameters often produce reasonable prediction -> easy to use.
- ③ avoid overfitting if some good trees in the forest and easy feature selection -> high dimensional pb.
- ④ hard to beat in performance.
- ⑤ Bonus : supports multiple outputs!

RF Cons:

- ① Slow to provide new predictions -> ineffective for « real-time predictions ».
- ② Good for prediction but bad for description..

What about Regression ?

→ At the same

aggregat° is by averaging

Questions?

Boosting

often even better
than RF

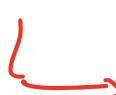


Another import technique (**very powerful !**)

Main idea:

- Give more importance to difficult point iteratively.
- Incrementally building an ensemble by training each new instance to emphasize the training instances previously mis-modeled.

Example: AdaBoost

 slower to train because you cannot train independently.

See: <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.AdaBoostClassifier.html>

Bonus: Multiple Outputs

- **Goal:** predict several outputs simultaneously
- **Solution:**
 - ▶ each leaf contains a value for each output
 - ▶ to chose the splits, we use a (weighted) average of the impurity of each output.

Face completion with multi-output estimators



Bonus: complexity

Total complexity for one decision tree:

Worst case:

$$O(n^2 d)$$

Balanced case:

$$O(nd)$$

Tips

- Decision trees tend to overfit: limit the depth or use `min_samples_leaf` (5 is a good initial pick)
- Visualize the tree with a small depth first
- Balance your dataset: trees tend to be biased towards dominating class.

Conclusion

- ① Trees are one of the simplest method (the most intuitive / human like) 
- ② Random Forests give excellent results in many applications. 

Lab :

- Example on a synthetic dataset of time series
- Application to inflation prediction in Brazil. 