# trees & forest

agenda:

1. decision trees: introduction

2. decision trees: how does it work?

3. random forest: introduction

4. rondom forest: how does it work?

5. hands on

## Decision Tree (DT)

metaphoric basis:
a standard tree... well... upside down

with all the bells and whistles:
- root
- branching (nodes)
- branches (edges)
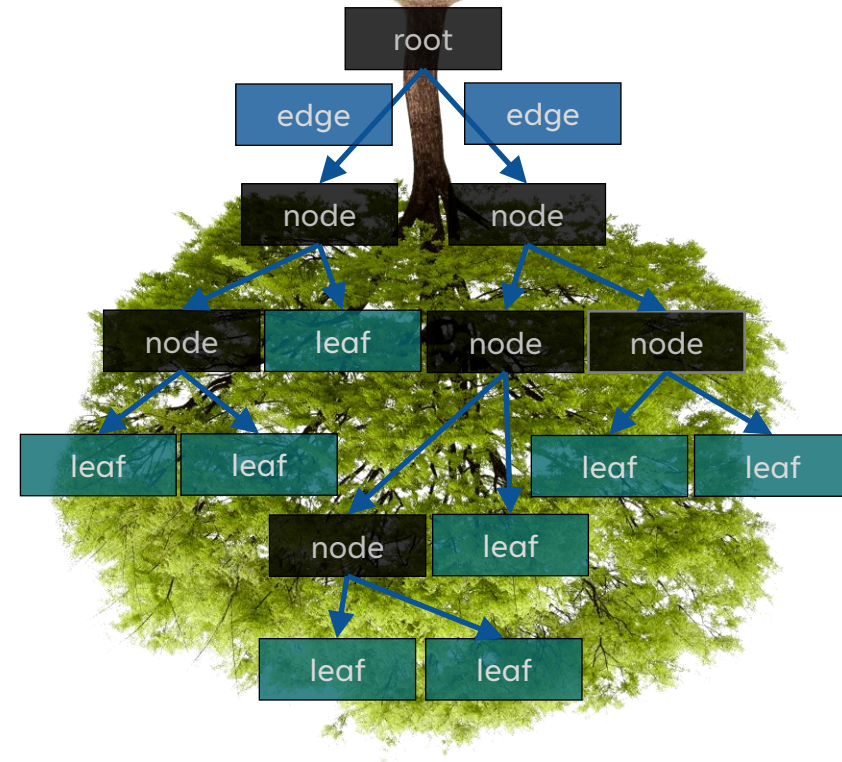- leafs, terminal nodes

## the main principle

„question-and-answer"

root/node:    interrogate the data
edge:         if „yes": go left, else: go right
leaf:         contains final decision / statement

## Decision Tree (DT)

metaphoric basis:
a standard tree... well... upside down

with all the bells and whistles:
- root
- branching (nodes)
- branches (edges)
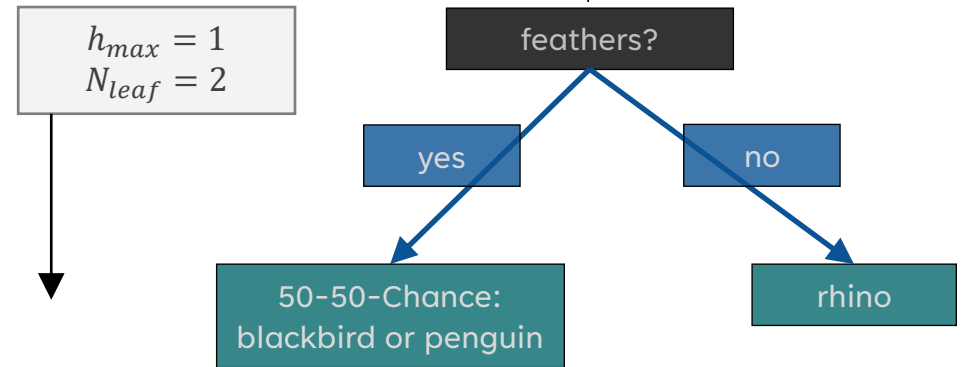- leafs, terminal nodes

## the main principle

„question-and-answer"

root/node:   interrogate the data
edge:           if „yes": go left, else: go right
leaf:            contains final decision / statement

## example

classification of animals

**pruned DT:**
multiple statements in terminal node („impurity")



$$h_{max} = 1$$
$$N_{leaf} = 2$$

feathers?

yes

no

50-50-Chance:
blackbird or penguin

rhino

## Decision Tree (DT)

metaphoric basis:
a standard tree... well... upside down

with all the bells and whistles:
- root
- branching (nodes)
- branches (edges)
- leafs, terminal nodes

## the main principle

„question-and-answer"

root/node:    interrogate the data
edge:          if „yes": go left, else: go right
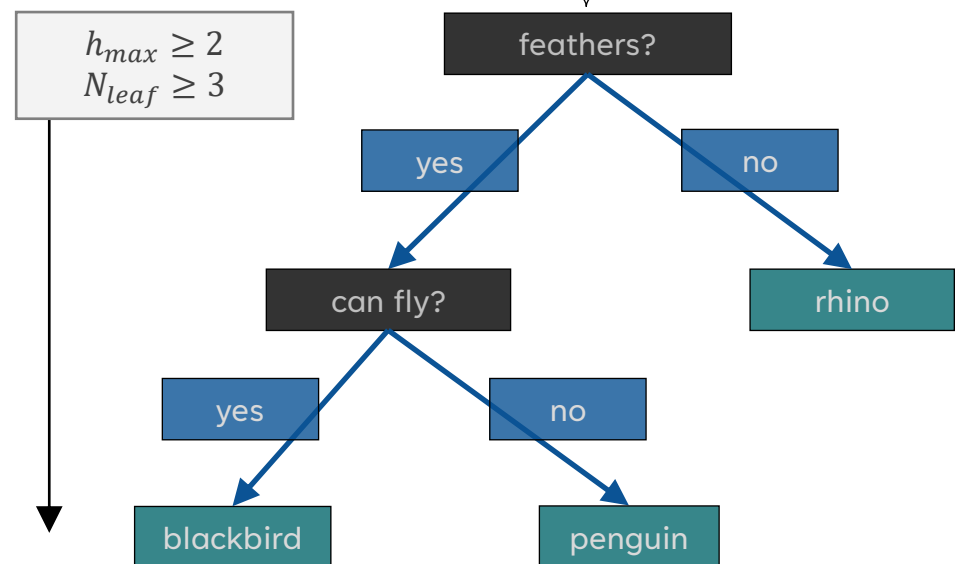leaf:          contains final decision / statement

## example

classification of animals

**unpruned DT:**
single statement per leaf („purity")

➜ important **hyperparameter**: $h_{max}$ or $N_{leaf}$



$h_{max} \geq 2$
$N_{leaf} \geq 3$

feathers?

yes          no

can fly?          rhino

yes          no

blackbird          penguin

## mathematical formulation of Q&A

regression & classification:
1. a DT splits the domain into $M$ subdomains
2. constant value $c_m$ in every subdomain $R_m$

decision to split a domain:
➔ metric: **node impurity** $Q$

## metric for regression

$N$: data points, $y$: true value, $\hat{y}$: predicted value
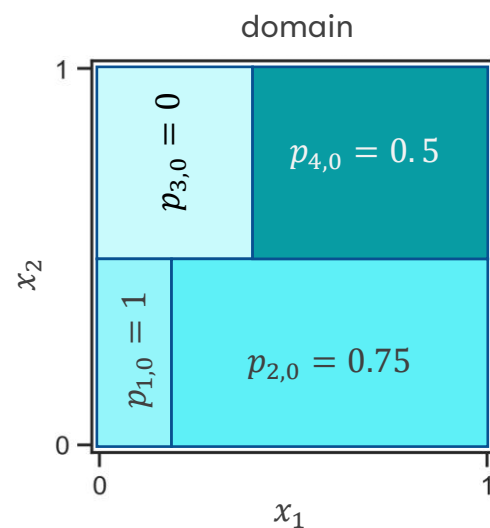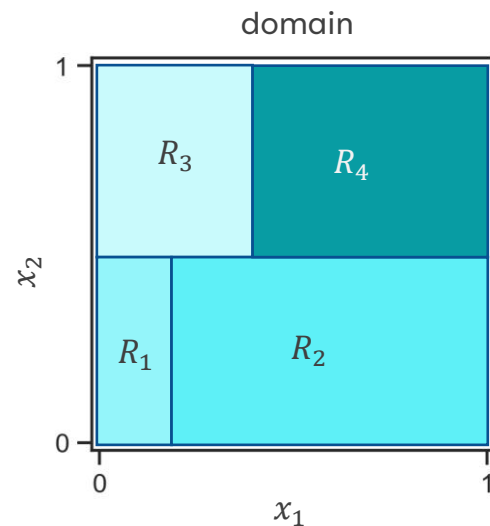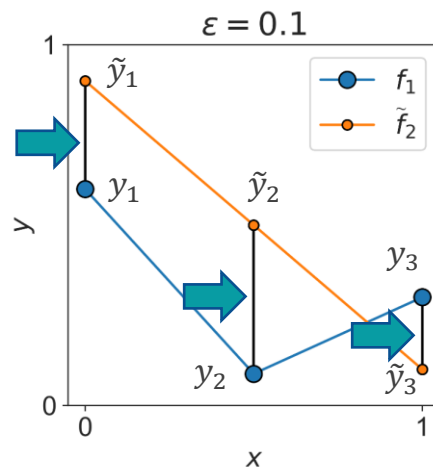**Mean Squared Error** (MSE):

$$Q = \varepsilon_{MSE}(\boldsymbol{y}, \widehat{\boldsymbol{y}}) = \frac{1}{N} \sum_{i=0}^{N-1} (y_i - \hat{y}_i)^2$$

## metric for classification

$p_{m,k}$: probability of occurrence of class $k \in K$
in node $m \in M$
**Gini (impurity) Index**:

$$Q_m = \sum_{k=1}^{K} p_{m,k}(1 - p_{m,k})$$

domain

$x_2$

$R_3$    $R_4$

$R_1$    $R_2$

$x_1$

$\varepsilon = 0.1$

$\tilde{y}_1$

$f_1$
$\tilde{f}_2$

$y_1$    $\tilde{y}_2$

$y_3$

$y_2$    $\tilde{y}_3$

$y$

$x$

domain

$x_2$

$p_{3,0} = 0$    $p_{4,0} = 0.5$

$p_{1,0} = 1$    $p_{2,0} = 0.75$

$x_1$

## mathematical formulation of Q&A

regression & classification:
1. a DT splits the domain into $M$ subdomains
2. constant value $c_m$ in every subdomain $R_m$

decision to split a domain:
➔ metric: **node impurity** $Q$

## metric for regression

$N$: data points, $y$: true value, $\hat{y}$: predicted value
**Mean Squared Error** (MSE):

$$Q = \varepsilon_{MSE}(\boldsymbol{y}, \hat{\boldsymbol{y}}) = \frac{1}{N} \sum_{i=0}^{N-1} (y_i - \hat{y}_i)^2$$

## example

$f: \mathbb{R} \to \mathbb{R}, f(x) = x^2 = y$
mean is „best guess" for each subdomain
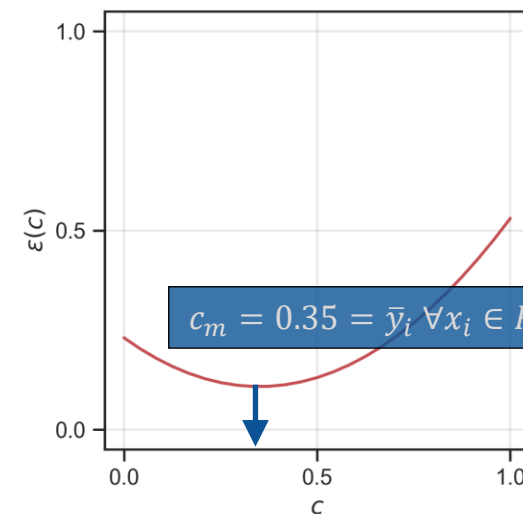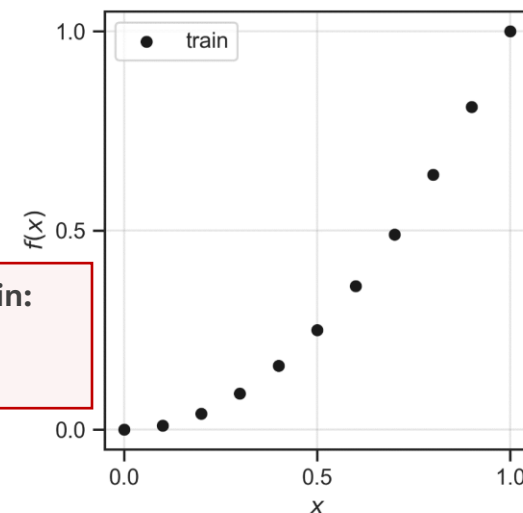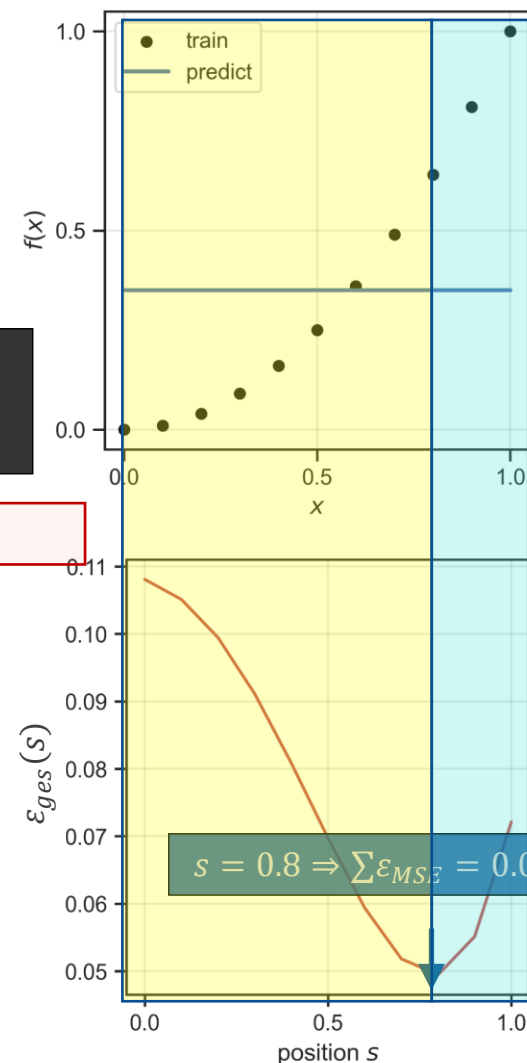$\hat{y}_i = c_m = \bar{y}_i \; \forall x_i \in R_m$

$$n = 11$$
$$\hat{y} =?$$
$$\varepsilon_{MSE} =?$$

**$R_m$ is currently whole domain:**
what is the best
$c_m = \hat{y}_i \forall x_i \in R_m$



$c_m = 0.35 = \bar{y}_i \; \forall x_i \in R_m$

## mathematical formulation of Q&A

regression & classification:
1. a DT splits the domain into $M$ subdomains
2. constant value $c_m$ in every subdomain $R_m$

decision to split a domain:
➔ metric: **node impurity** $Q$

## metric for regression

$N$: data points, $y$: true value, $\hat{y}$: predicted value
**Mean Squared Error** (MSE):

$$Q = \varepsilon_{MSE}(\boldsymbol{y}, \hat{\boldsymbol{y}}) = \frac{1}{N} \sum_{i=0}^{N-1} (y_i - \hat{y}_i)^2$$

## example
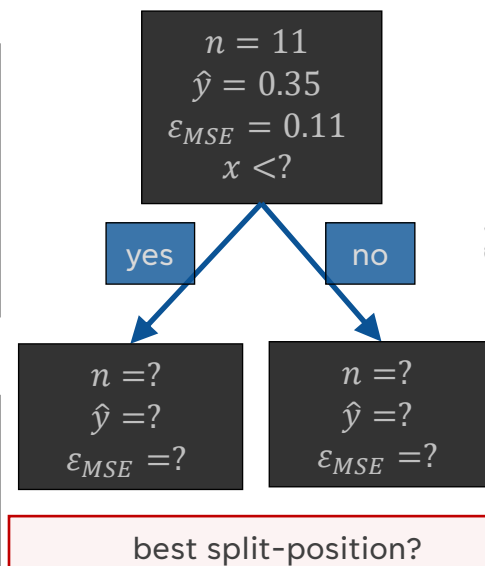
$f: \mathbb{R} \to \mathbb{R}, f(x) = x^2 = y$
mean is „best guess" for each subdomain
$\hat{y}_i = c_m = \bar{y}_i \, \forall x_i \in R_m$

minimize $\varepsilon_{ges}$ by splitting domain at position $s$:

$$\varepsilon_{ges}(s) = \sum_{m=1}^{M=2} \varepsilon_m(s)$$



$n = 11$
$\hat{y} = 0.35$
$\varepsilon_{MSE} = 0.11$
$x <?$

yes          no

$n =?$        $n =?$
$\hat{y} =?$   $\hat{y} =?$
$\varepsilon_{MSE} =?$   $\varepsilon_{MSE} =?$

best split-position?

$s = 0.8 \Rightarrow \sum \varepsilon_{MSE} = 0.05$

## mathematical formulation of Q&A

regression & classification:
1. a DT splits the domain into $M$ subdomains
2. constant value $c_m$ in every subdomain $R_m$

decision to split a domain:
➜ metric: **node impurity** $Q$

## metric for regression

$N$: data points, $y$: true value, $\hat{y}$: predicted value
**Mean Squared Error** (MSE):

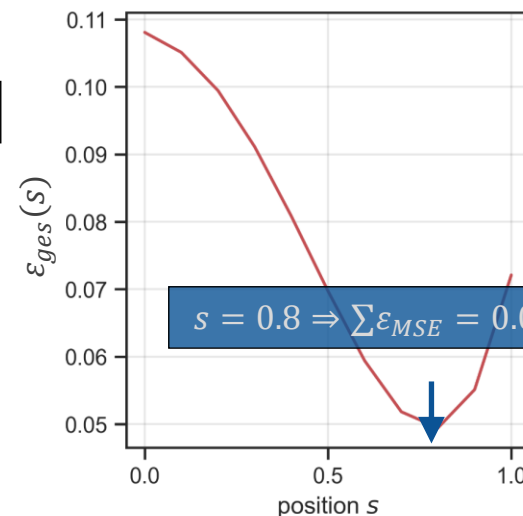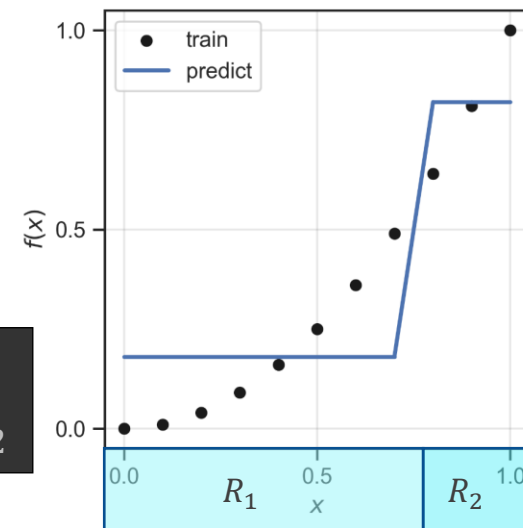$$Q = \varepsilon_{MSE}(\mathbf{y}, \widehat{\mathbf{y}}) = \frac{1}{N} \sum_{i=0}^{N-1} (y_i - \hat{y}_i)^2$$
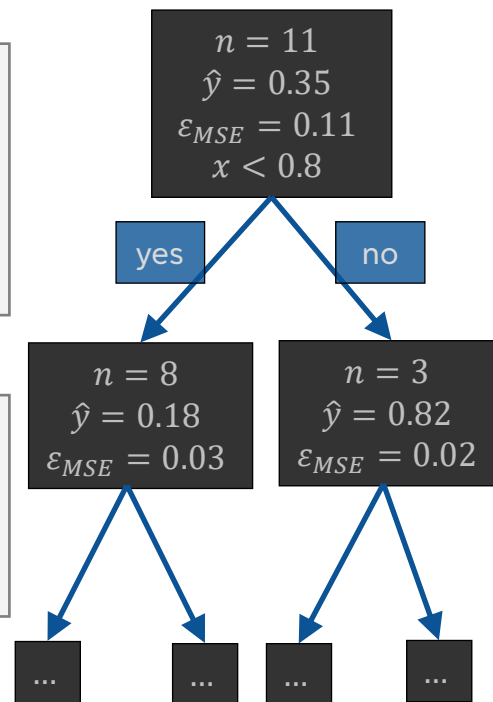
## example

$f: \mathbb{R} \to \mathbb{R},\ f(x) = x^2 = y$
mean ist „best guess" for each subdomain
$\hat{y}_i = c_m = \bar{y}_i\ \forall x_i \in R_m$

minimize $\varepsilon_{ges}$ by splitting domain at position $s$:

$$\varepsilon_{ges}(s) = \sum_{m=1}^{M=2} \varepsilon_m(s)$$

$n = 11$
$\hat{y} = 0.35$
$\varepsilon_{MSE} = 0.11$
$x < 0.8$

yes      no

$n = 8$
$\hat{y} = 0.18$
$\varepsilon_{MSE} = 0.03$

$n = 3$
$\hat{y} = 0.82$
$\varepsilon_{MSE} = 0.02$

...  ...  ...  ...

$s = 0.8 \Rightarrow \sum \varepsilon_{MSE} = 0.05$

## mathematical formulation of Q&A

regression & classification:
1. a DT splits the domain into $M$ subdomains
2. constant value $c_m$ in every subdomain $R_m$

decision to split a domain:
➜ metric: **node impurity** $Q$

## metric for classification

$p_{m,k}$: probability of occurrence of class $k \in K$ in node $m \in M$
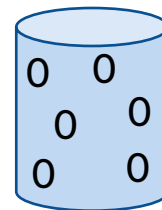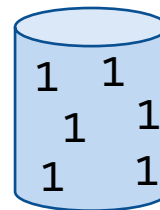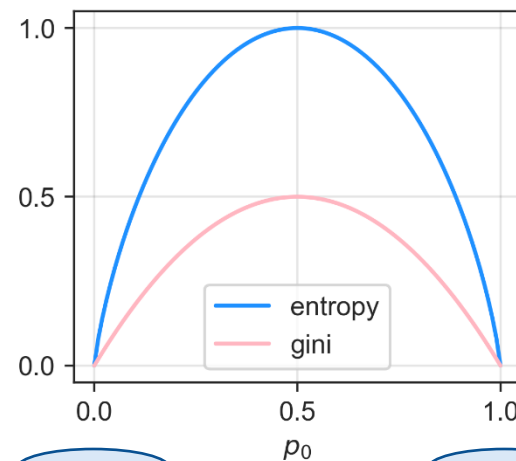
**Gini (impurity) Index**:

$$Q_m = \sum_{k=1}^{K} p_{m,k}(1 - p_{m,k})$$

**Entropy** $H$:

$$H = -\sum_{k=1}^{K} p_{m,k} \ln(p_{m.k})$$

**Gini Index** vs. **Entropy** ($K$=2)

## mathematical formulation of Q&A

regression & classification:
1. a DT splits the domain into $M$ subdomains
2. constant value $c_m$ in every subdomain $R_m$
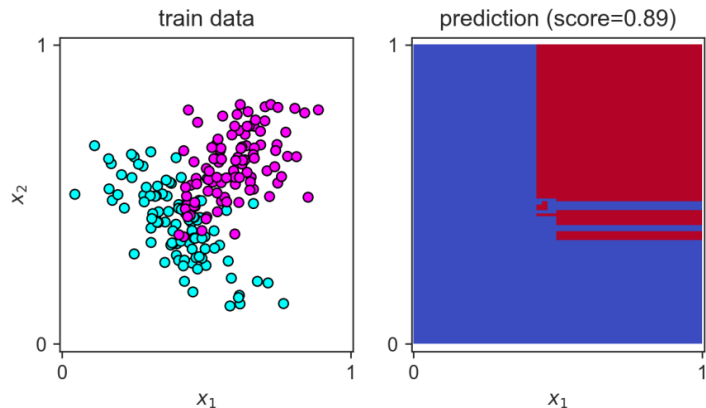
decision to split a domain:
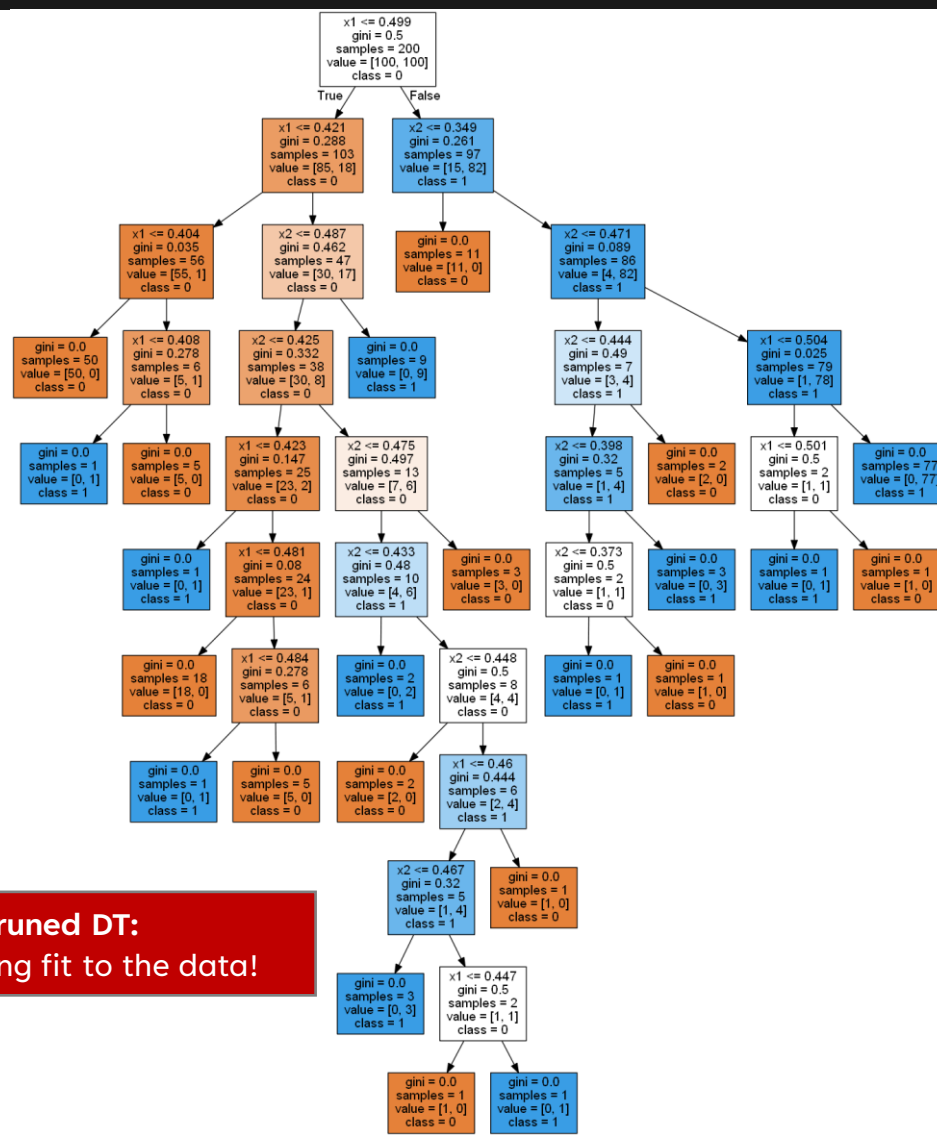➔ metric: **node impurity** $Q$

## metric for classification

$p_{m,k}$: probability of occurrence of class $k \in K$ in node $m \in M$

**Gini (impurity) Index**:

$$Q_m = \sum_{k=1}^{K} p_{m,k}(1 - p_{m,k})$$



train data

prediction (score=0.89)

**unpruned DT:**
strong fit to the data!
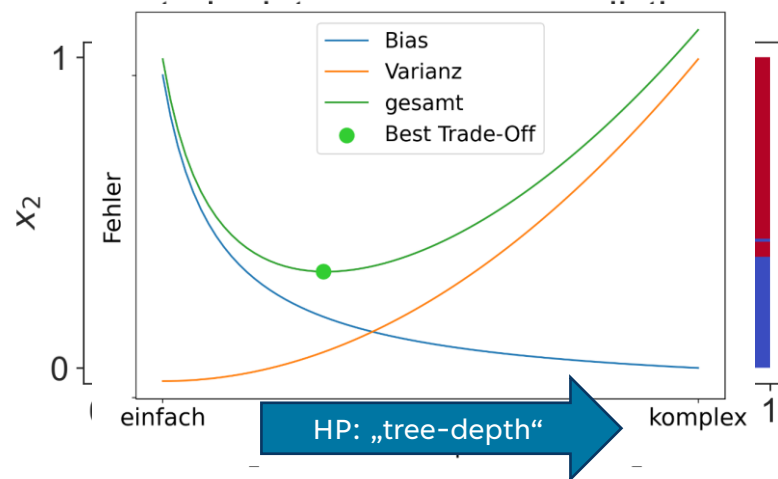
## bias & variance

general definition of a function for approximation:
$$\hat{f}: \mathbb{R}^{n,m} \rightarrow \mathbb{R}^m, \hat{f}(X) = \hat{y} = y + \varepsilon$$

the **approximation error** $\varepsilon$ contains of:
- unknown influences
- **model-bias**:
  simplifications, high when underfitting
- **model-variance**
  high komplexity, high when overfitting
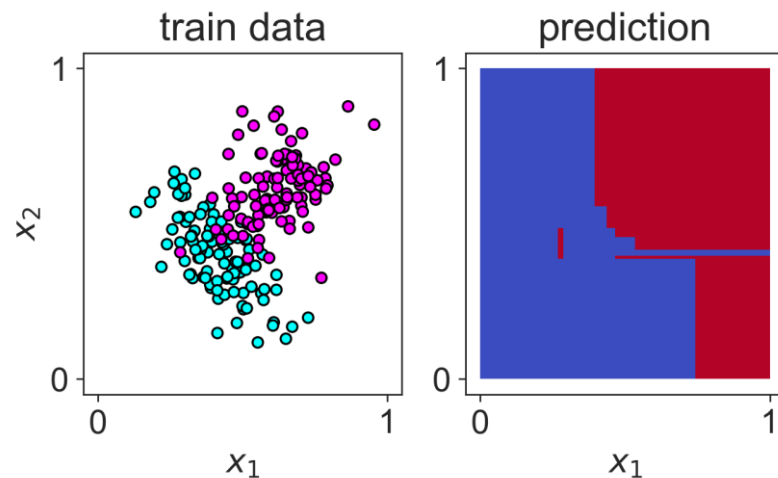- ➔ **IMPORTANT: bias-variance trade-off**



HP: „tree-depth"

## unpruned DT: pros & cons

pros:
- can handle mixed and redundant variables
- **small Bias**
- ...

cons:
- **high varianz**
- **usually prediction is not very good**
- ...

## bias & variance

general definition of function for approximation:
$$\hat{f}\colon \mathbb{R}^{n,m} \to \mathbb{R}^m, \hat{f}(\boldsymbol{X}) = \widehat{\boldsymbol{y}} = \boldsymbol{y} + \boldsymbol{\varepsilon}$$

the **approximation error** $\varepsilon$ contains of:
- unknown influences
- **modell-Bias**:
  simplifications, high when underfitting
- **modell-Variance**
  high komplexity, high when overfitting
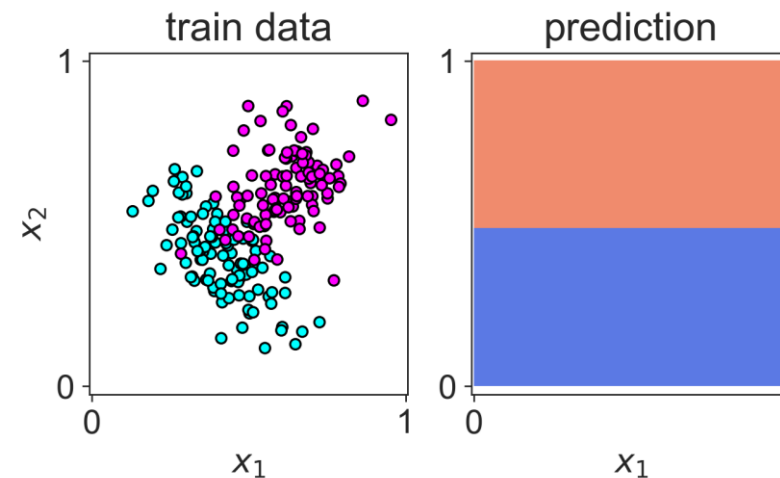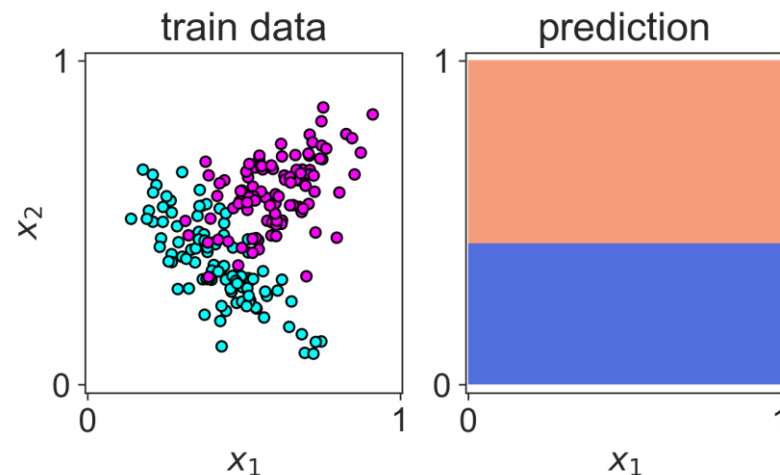- ➔ **IMPORTANT: bias-variance trade-off**



## pruned DT: pros & cons

pros:
- can handle mixed and redundant variables
- **small variance**
- ...

cons:
- **high bias**
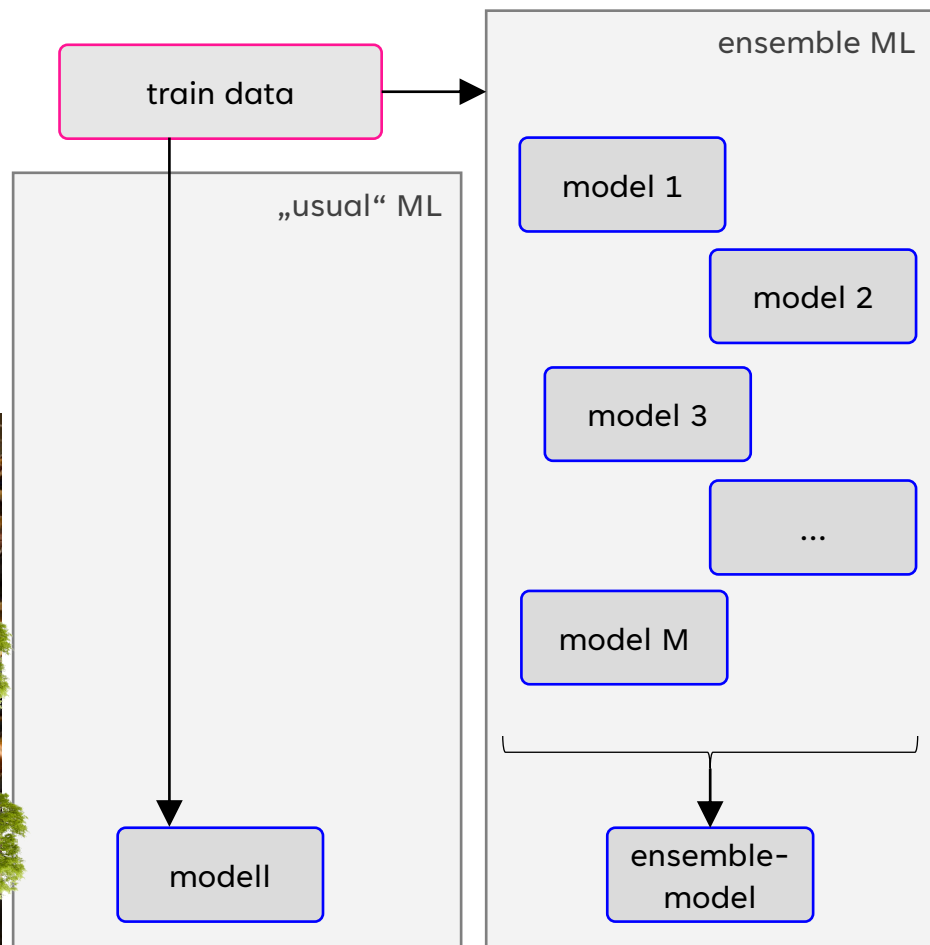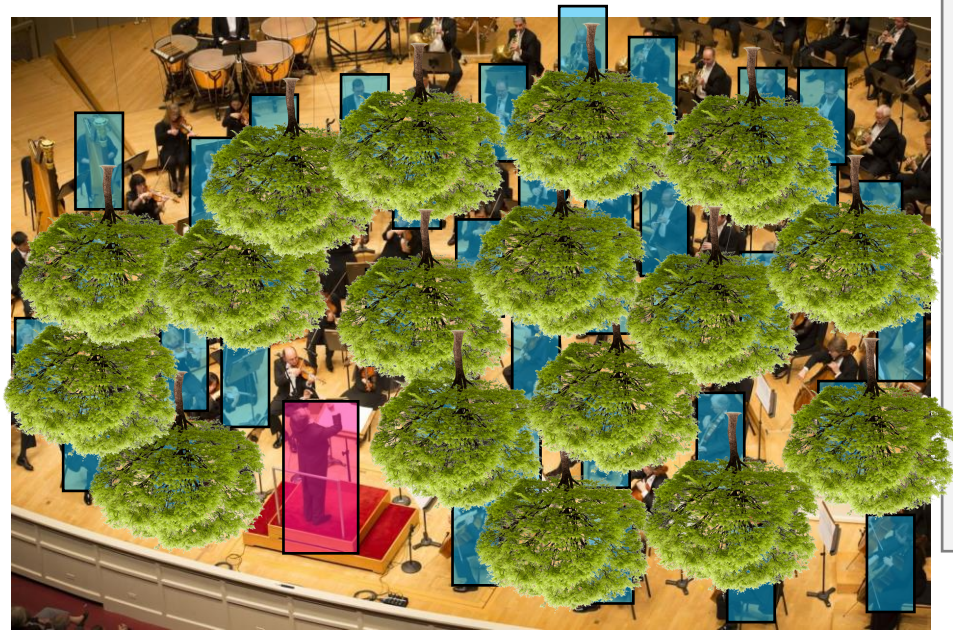- **usually prediction is not very good**
- ...

## what is a random forest?

it's an ensemble

## what is an ensemble?

it's an aggregation...
- of multiple models (usually DTs)
- to exploit pros
- to avoid cons



train data

"usual" ML

modell

ensemble ML

model 1

model 2

model 3

...

model M

ensemble-model

## what is a random forest?

it's an ensemble

## what is an ensemble?

it's an aggregation…
- of multiple models (usually DTs)
- to exploit pros
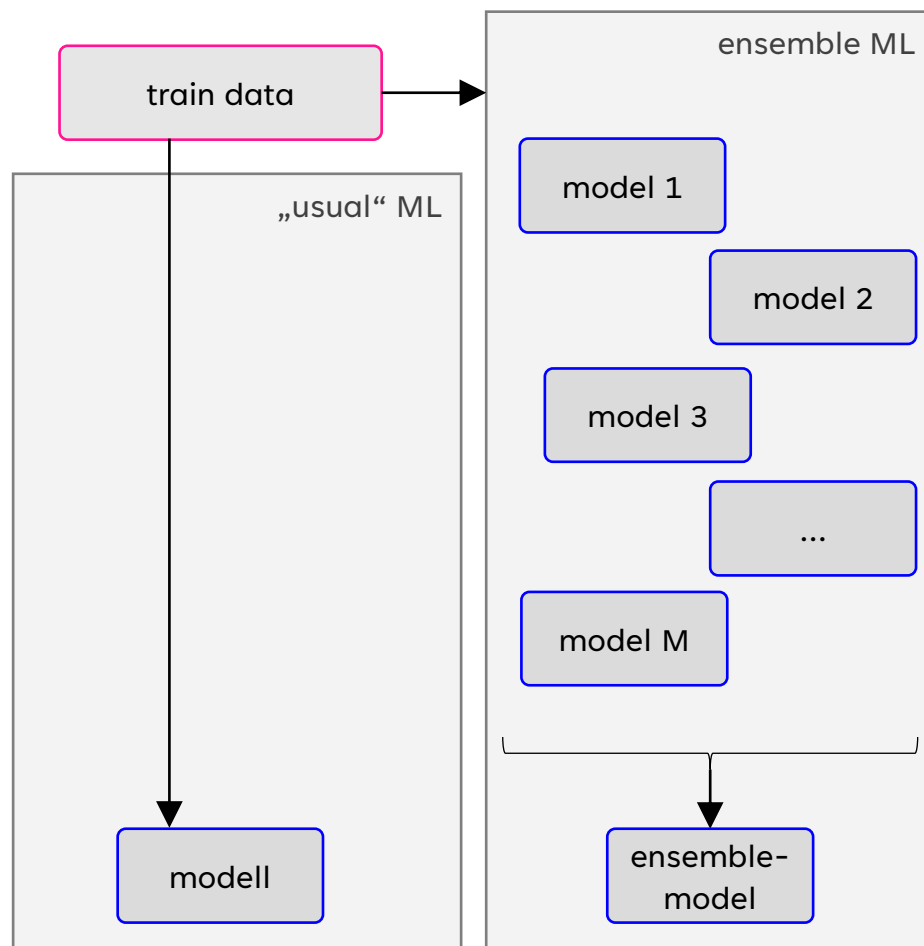- to avoid cons

## different types of ensembles

**Bagging** [Breiman, 1996]
- **Random Forest** [Breiman, 2001]

Boosting:
- AdaBoost [Freund & Shapire, 1996]
- Gradient Boosting [Friedman, 1999]
  (*Extreme Gradient Boosting* [Chen & Guestrin, 2016 ])

Stacking

train data

"usual" ML

ensemble ML

model 1

model 2

model 3

…

model M

modell

ensemble-model

## starting point

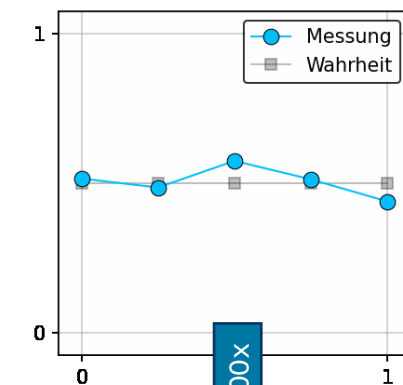unpruned DT*
- pro: small bias
- con: high variance

## bagging

short for „booststrap aggrigation"

idea:

build an ensemble of unpruned DTs
- exploit pro of **small bias** (strong adaption)
- avoid con of high variance by **boostrapping** (averaging **random resamples**)



100x

Box-Whisker-Plot

## starting point

unpruned DT*
- pro: small bias
- con: high variance
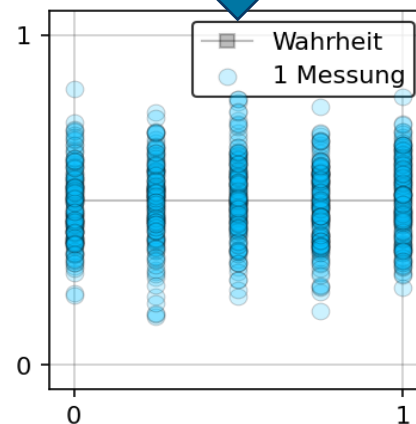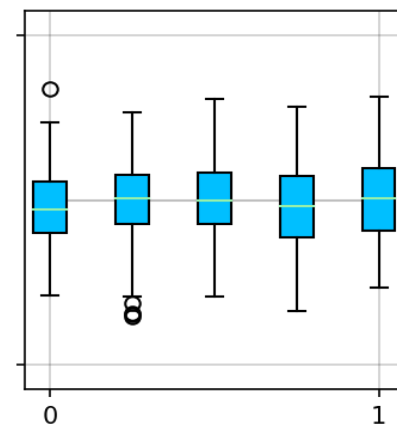
## bagging

short for „booststrap aggrigation"

idea:
build an ensemble of unpruned DTs
- exploit pro of **small bias** (strong adaption)
- avoid con of high variance by **boostrapping** (averaging **random resamples**)

## process

1. resample data $X \in \mathbb{R}^{n,p}$ by bootstrapping $L \in \mathbb{N}$ times:
   $$X \to \{X_{BS,1}, ..., X_{BS,L}\} \text{ with } X_{BS,l} \in \mathbb{R}^{n,p}$$
2. train $L$ models using the $L$ new datasets
   $$\{X_{BS,1}, ..., X_{BS,L}\} \to \{M_1, ..., M_L\}$$
3. answer of the whole ensemble $M_{ges}$ : averaging all ensemble-members
   $$M_{ges} = \frac{1}{L} \sum_{l=1}^{L} M_l$$



bagging

$X \in \mathbb{R}^{n,p}$

model $M_1$ ← $X_{BS,1}$

model $M_2$ ← $X_{BS,2}$

... ←  ...

model $M_L$ ← $X_{BS,L}$

$\frac{1}{L} \sum_{l=1}^{L} M_l$

$M_{ges}$

* bagging works with different types of models too

## starting point

unpruned DT*
- pro: small bias
- con: high variance

## bagging

short for „booststrap aggrigation"

idea:
build an ensemble of unpruned DTs
- exploit pro of **small bias** (strong adaption)
- avoid con of high variance by **boostrapping** (averaging **random resamples**)

## process

1. resample data $X \in \mathbb{R}^{n,p}$
   by bootstrapping $L \in \mathbb{N}$ times:
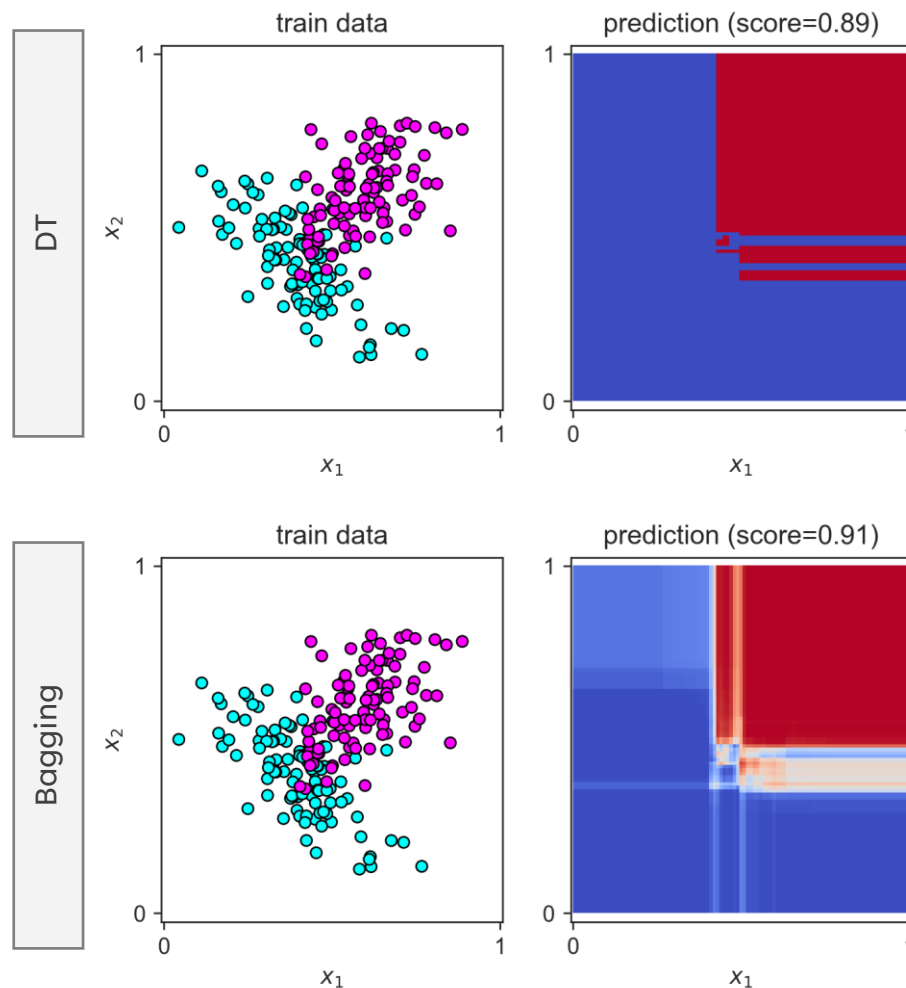   $$X \rightarrow \{X_{BS,1}, ..., X_{BS,L}\} \text{ with } X_{BS,l} \in \mathbb{R}^{n,p}$$
2. train $L$ models using the $L$ new datasets
   $$\{X_{BS,1}, ..., X_{BS,L}\} \rightarrow \{M_1, ..., M_L\}$$
3. answer of the whole ensemble $M_{ges}$ :
   averaging all ensemble-members
   $$M_{ges} = \frac{1}{L} \sum_{l=1}^{L} M_l$$



DT

train data

prediction (score=0.89)

Bagging

train data

prediction (score=0.91)

* bagging works with different types of models too

## extended bagging

decreasing variance even further by...
**uncorrelated DTs**

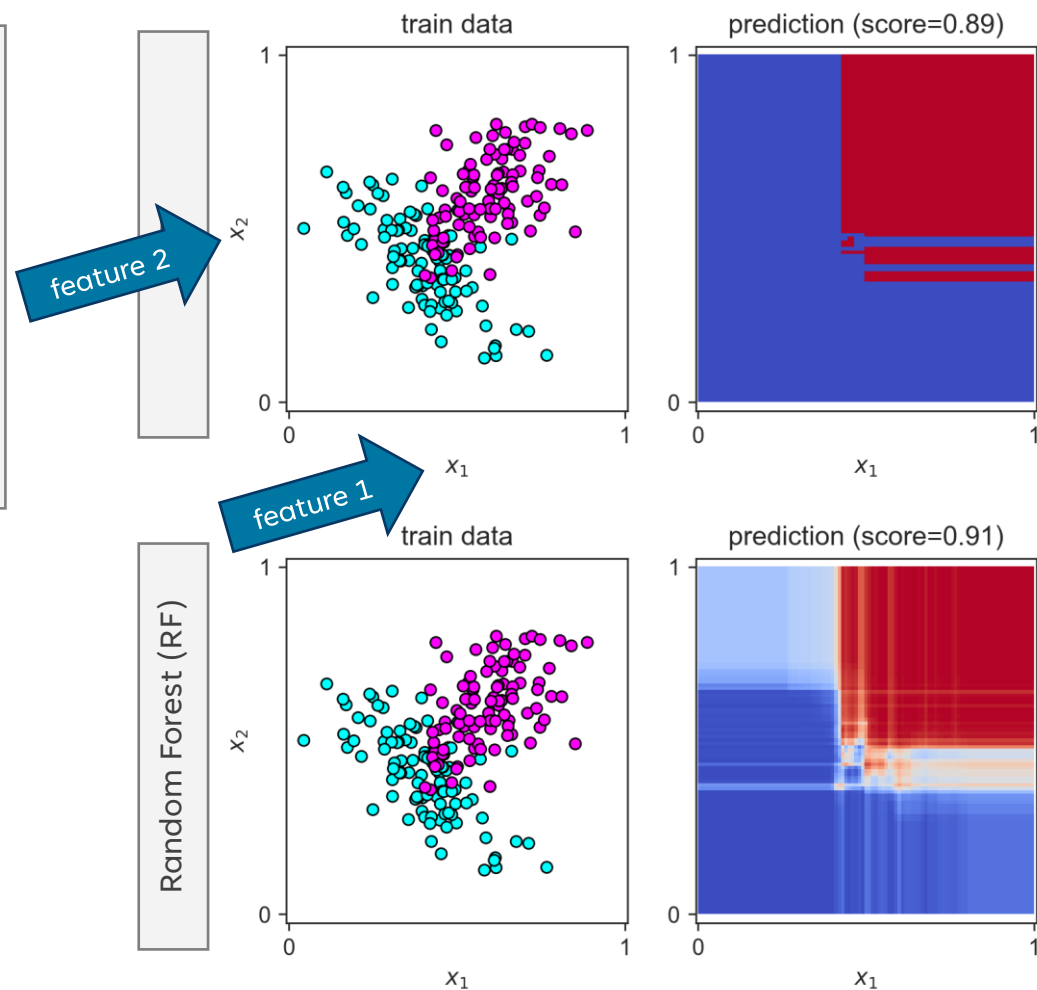$X \in \mathbb{R}^{n,p}$ training data, $n$ samples, $p$ features

<u>random forest:</u>
train DTs with randomly selected $m < p$ features
$X \in \mathbb{R}^{n,m}$

e.g.:
$m = \sqrt{p}$ or $\log(p)$

feature 2

feature 1

Random Forest (RF)

train data

$x_2$

$x_1$

prediction (score=0.89)

$x_1$

train data

$x_2$

$x_1$

prediction (score=0.91)

$x_1$

## code available*

clone or download GitHub-Repository
https://github.com/saifedias/tree_randomForest.git


online Notebook via Binder
https://mybinder.org/v2/gh/saifedias/tree_randomForest.git/HEAD


## would you like to know more? – a short outline

bagging, boosting, stacking
https://towardsdatascience.com/ensemble-methods-bagging-boosting-and-stacking-c9214a10a205

ensemble learning
https://www.kaggle.com/discussions/general/263786

# BTU ML-Group

contact:
marlon.lehmann@b-tu.de