# Decision trees Bagging

**Vladislav Goncharenko**

ML Teamlead, DZEN

girafe ai

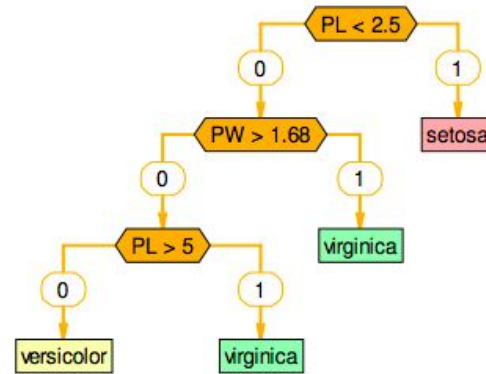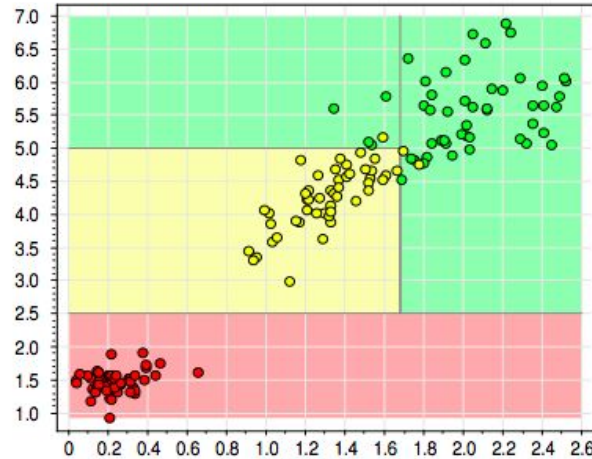MSU, spring 2024

# Outline

1. Intuition
2. Construction procedure
3. Information criteria
4. Special highlights
   - Dealing with missing data
   - Binarization
   - Decision tree as linear model
   - Standards
   - Hyperparameters
5. Bootstrap and Bagging
6. Random Forest

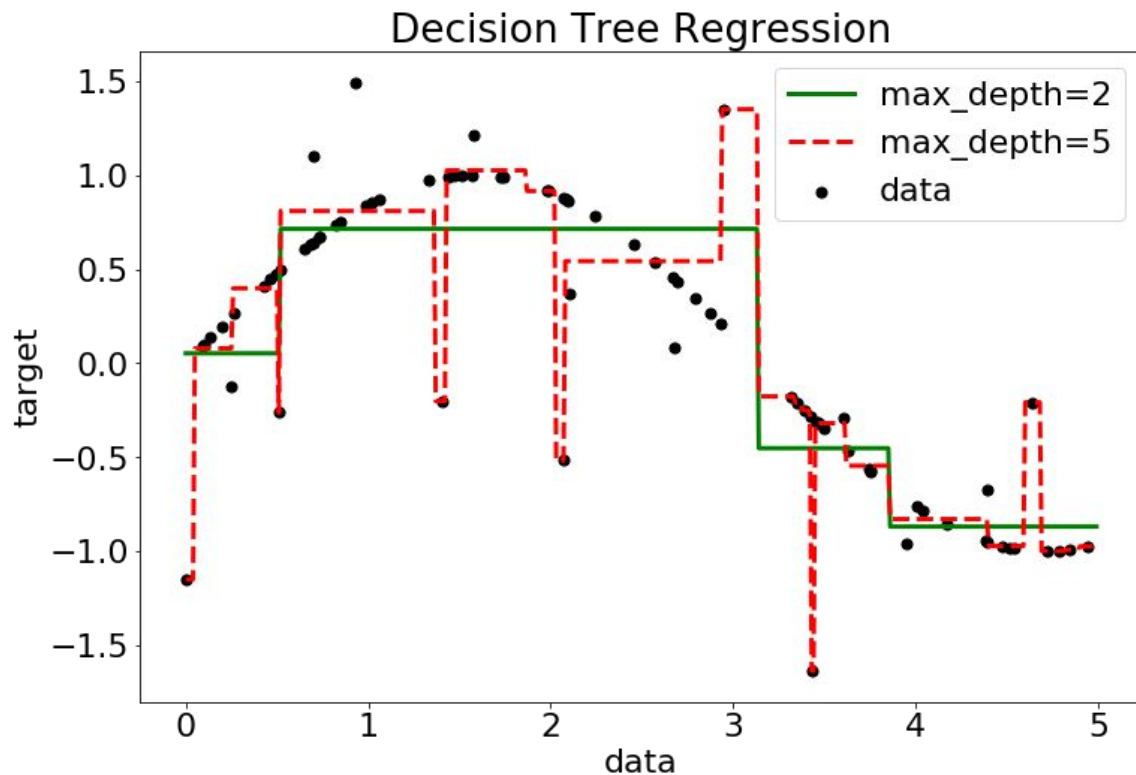# Decision Tree: intuition

girafe
ai

01

# Decision tree for Iris data set



| | |
|---|---|
| setosa | $r_1(x) = \big[PL \leqslant 2.5\big]$ |
| virginica | $r_2(x) = \big[PL > 2.5\big] \wedge \big[PW > 1.68\big]$ |
| virginica | $r_3(x) = \big[PL > 5\big] \wedge \big[PW \leqslant 1.68\big]$ |
| versicolor | $r_4(x) = \big[PL > 2.5\big] \wedge \big[PL \leqslant 5\big] \wedge \big[PW < 1.68\big]$ |

# Decision tree in regression


Decision Tree Regression

Green - decision tree of depth 2

Red - decision tree of depth 5

Every leaf corresponds to some constant.

# Construction procedure
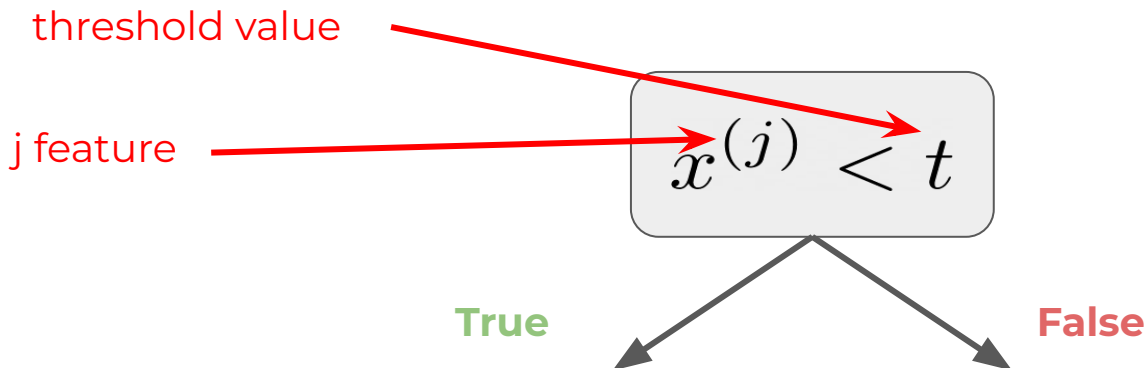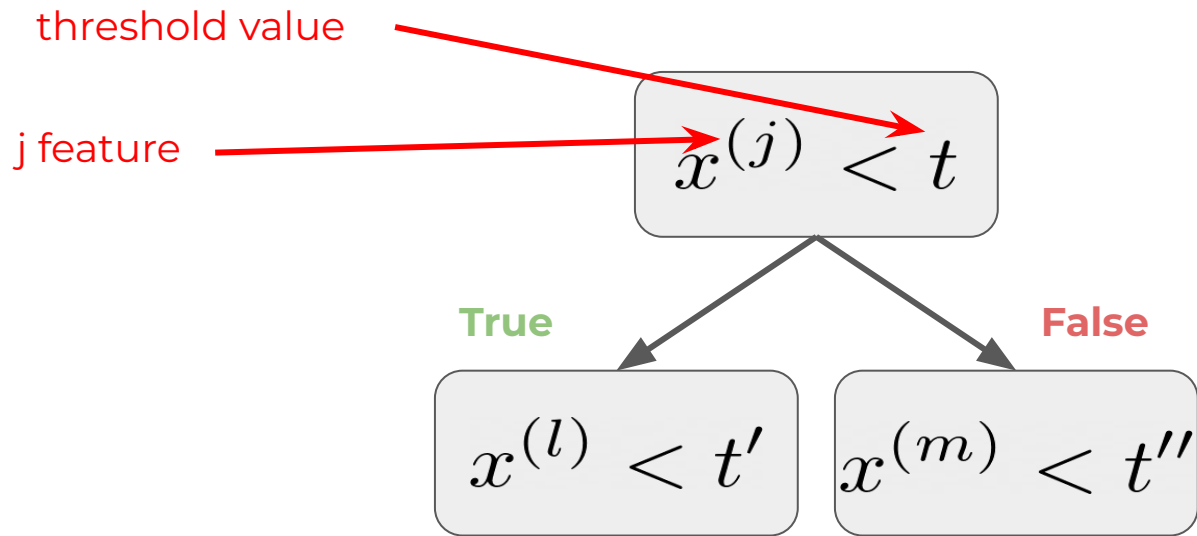
girafe
ai

02

# Constructing decision trees

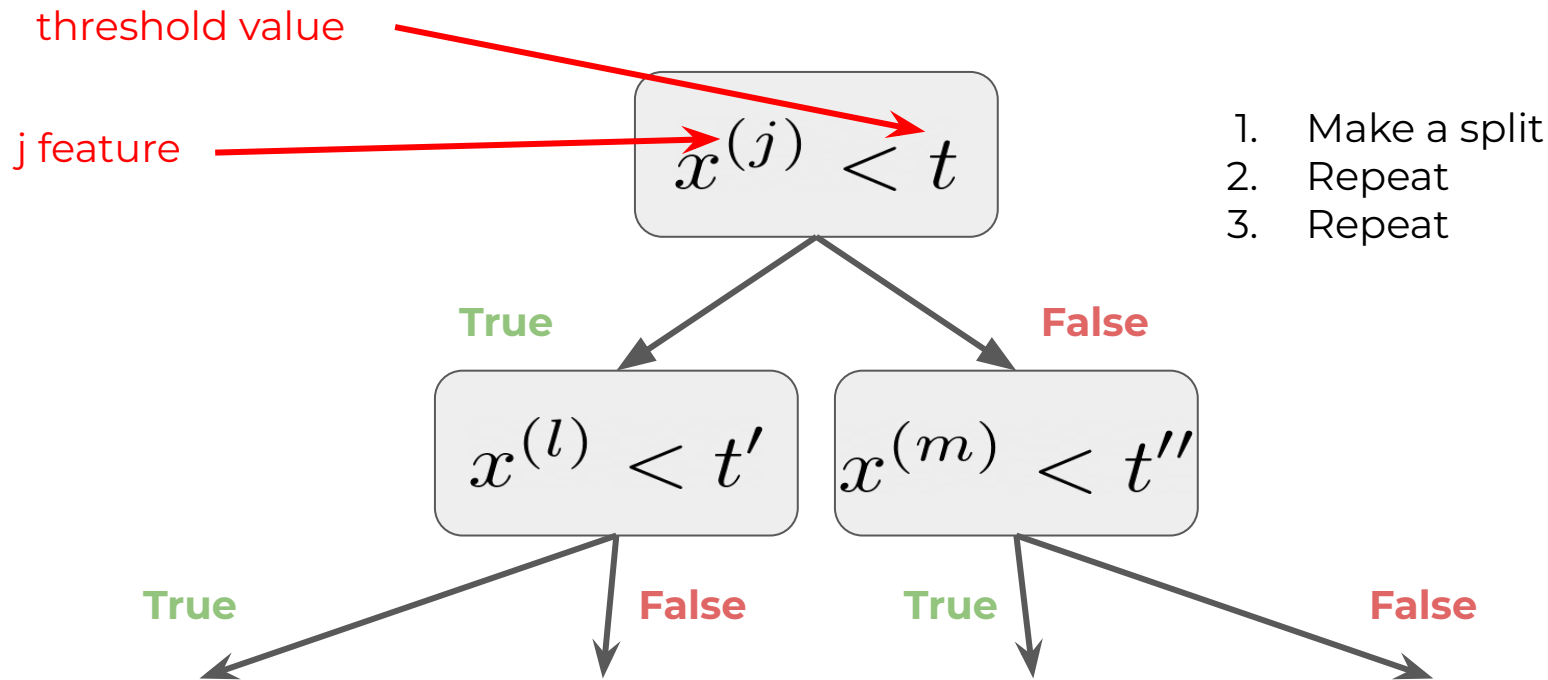threshold value

j feature

$$x^{(j)} < t$$

**True**

**False**

1. Make a split

# Constructing decision trees

threshold value

j feature

$$x^{(j)} < t$$

**True**

**False**

$$x^{(l)} < t'$$

$$x^{(m)} < t''$$

1. Make a split
2. Repeat

# Constructing decision trees

threshold value

j feature

$$x^{(j)} < t$$

1. Make a split
2. Repeat
3. Repeat

**True**

**False**

$$x^{(l)} < t'$$

$$x^{(m)} < t''$$

**True**

**False**

**True**

**False**

# Constructing decision trees

threshold value

j feature

$$x^{(j)} < t$$

1. Make a split
2. Repeat
3. Repeat
4. ...

**True**

**False**

$$x^{(l)} < t'$$

$$x^{(m)} < t''$$

**True**

**False**

**True**

**False**

threshold value

j feature

True

True

# How to split data properly?

$$Q$$

$$x^{(j)} < t$$

$$L \qquad R$$

What is H?

$$\frac{|L|}{|Q|} H(L) + \frac{|R|}{|Q|} H(R) \longrightarrow \min_{j,t}$$
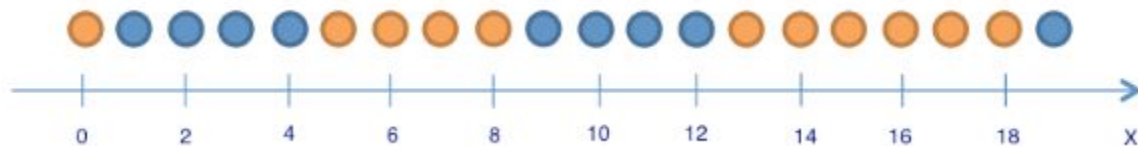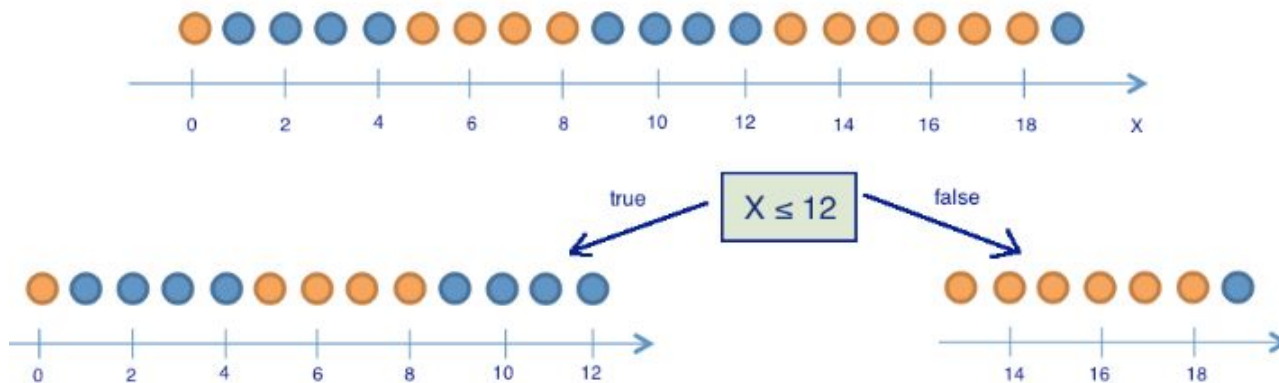
# Information criteria

girafe
ai

03

# Information criteria

H(R) is measure of "heterogeneity" of our data.

Consider binary classification problem:

# Information criteria

H(R) is measure of "heterogeneity" of our data.

Consider binary classification problem:

# Information criteria

H(R) is measure of "heterogeneity" of our data.

Consider binary classification problem:

Obvious way:
Misclassification criteria:
$$H(R) = 1 - \max\{p_0, p_1\}$$

1. Entropy criteria:
$$H(R) = -p_0 \log p_0 - p_1 \log p_1$$

2. Gini impurity:
$$H(R) = 1 - p_0^2 - p_1^2 = 1 - 2p_0 p_1$$

16

# Information criteria

H(R) is measure of "heterogeneity" of our data.

Consider multiclass classification problem:

Obvious way:
Misclassification criteria:

$$H(R) = 1 - \max_k \{p_k\}$$

1. Entropy criteria:

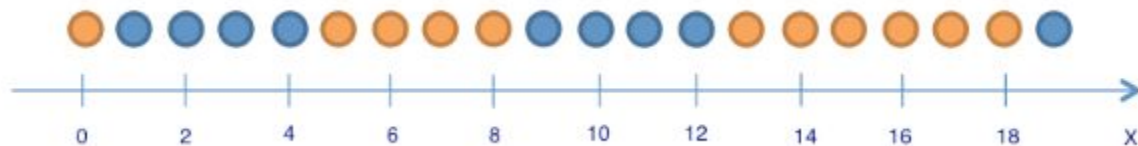$$H(R) = -\sum_{k=0}^{K} p_k \log p_k$$

2. Gini impurity:

$$H(R) = 1 - \sum_k (p_k)^2$$

# Information criteria

H(R) is measure of "heterogeneity" of our data.
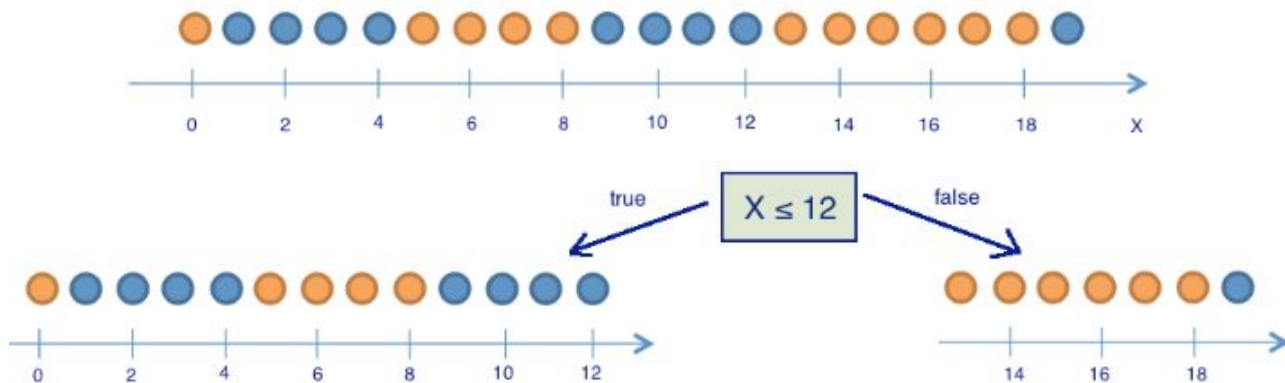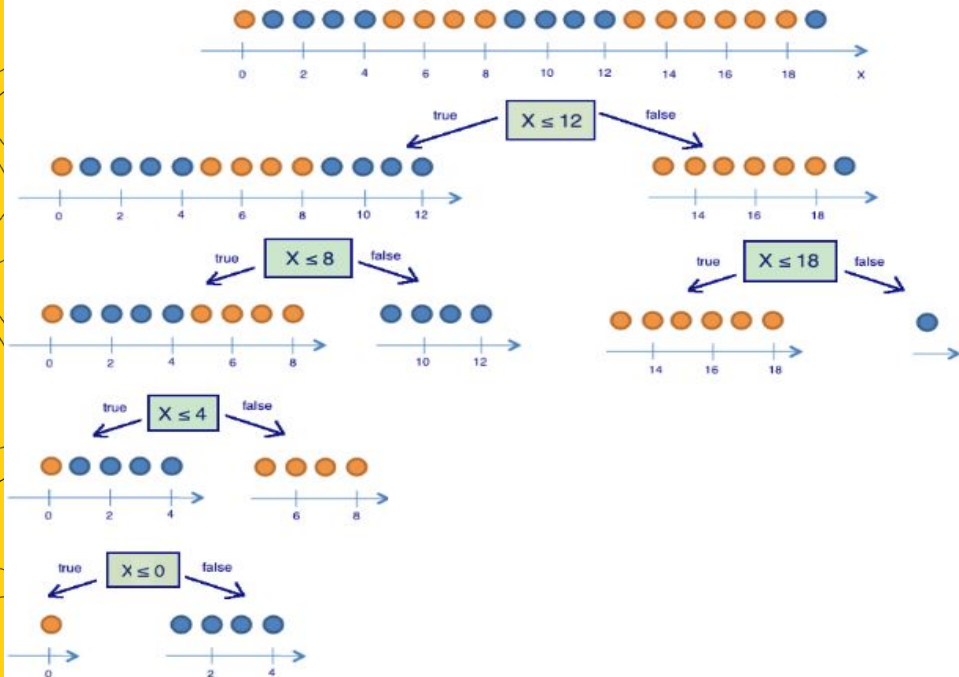
Consider binary classification problem:

# Information criteria

H(R) is measure of "heterogeneity" of our data.

Consider binary classification problem:

# Information criteria: Entropy



$$S = -M \sum_{k=0}^{K} p_k \log p_k$$

In binary case N = 2

$$S = -p_+ \log_2 p_+ - p_- \log_2 p_- = -p_+ \log_2 p_+ - (1 - p_+) \log_2 (1 - p_+)$$

source: https://habr.com/ru/company/ods/blog/322534/

# Information criteria: Gini impurity

$$G = 1 - \sum_k (p_k)^2$$

In binary case N = 2

$$G = 1 - p_+^2 - p_-^2 = 1 - p_+^2 - (1 - p_+)^2 = 2p_+(1 - p_+)$$

source: https://habr.com/ru/company/ods/blog/322534/

# Information criteria
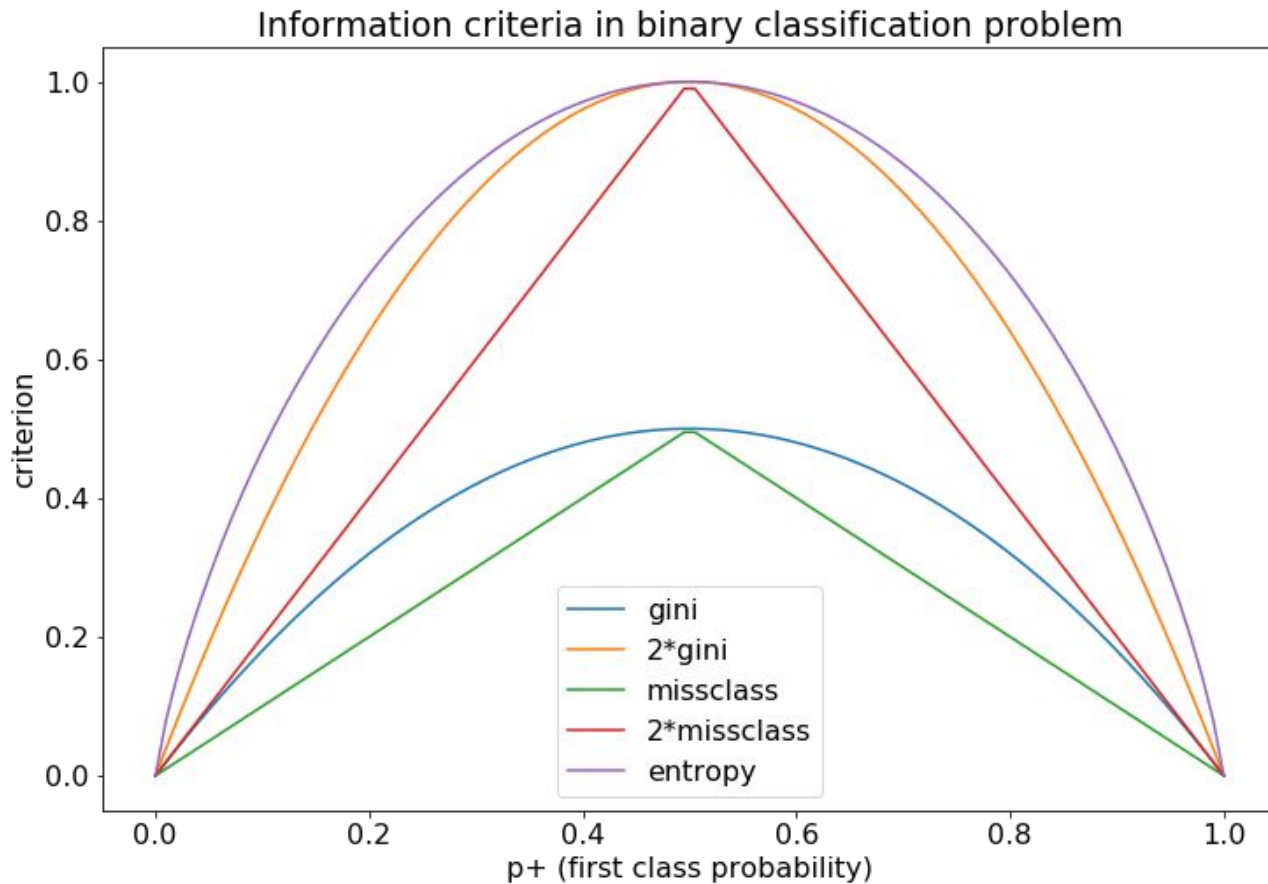
H(R) is measure of "heterogeneity" of our data.

Consider multiclass classification problem:

Obvious way: Misclassification criteria:

$$H(R) = 1 - \max_k \{p_k\}$$

1. Entropy criteria:

$$H(R) = -\sum_k p_k \log_2 p_k$$

2. Gini impurity:

$$H(R) = 1 - \sum_k (p_k)^2$$

# Information criteria



Information criteria in binary classification problem

# Information criteria

H(R) is measure of "heterogeneity" of our data.

Consider regression problem:

1. Mean squared error

$$H(R) = \min_c \frac{1}{|R|} \sum_{(x_i, y_i) \in R} (y_i - c)^2$$

What is the constant?

$$c^* = \frac{1}{|R|} \sum_{y_i \in R} y_i$$

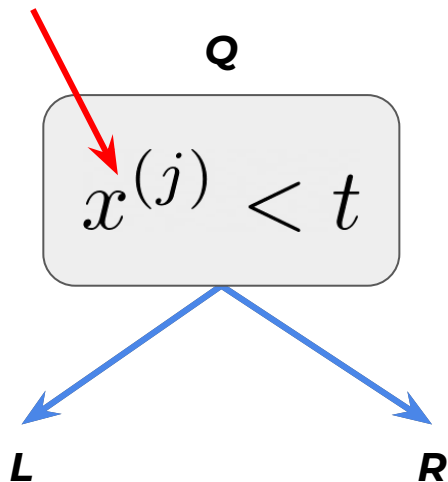# Special highlights

girafe
ai

04

# Missing values in Decision Trees

If the value is missing, one might use both sub-trees and average their predictions.

But this will negatively affect model computational performance.

Missing value



$$\hat{y} = \frac{|L|}{|Q|}\hat{y}_L + \frac{|R|}{|Q|}\hat{y}_R$$

26

# Missing values in Catboost

**Forbidden**: Missing values are not supported, their presence is interpreted as an error

**Min**: Missing values are processed as the minimum value (less than all other values) for the feature. It is guaranteed that a split that separates missing values from all other values is considered when selecting trees.

**Max**: Missing values are processed as the maximum value (greater than all other values) for the feature. It is guaranteed that a split that separates missing values from all other values is considered when selecting trees.
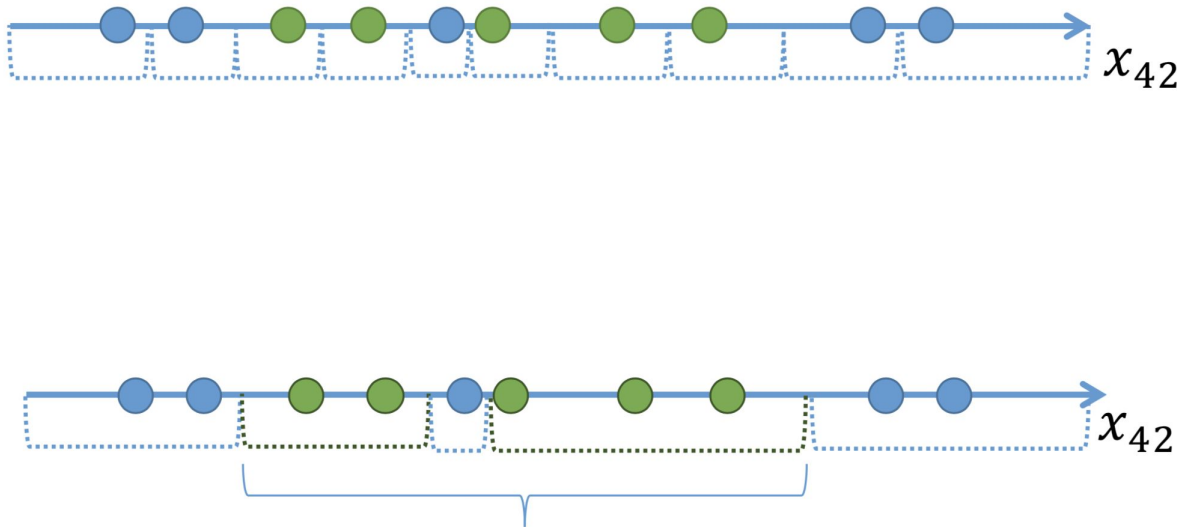
The **default** processing mode **is Min**

[Documentation](Documentation)

27

# Binarization

Idea: instead selecting one threshold define several for one feature.



e.g. Border count hyperparameter in Catboost (defaults to 254)

# Decision Trees as Linear models

Let J be the subspace of the original feature space, corresponding to the leaf of the tree.

Prediction takes form

$$\hat{y} = \sum_j w_j [x \in J_j]$$

# Standards

- ID-3
  - Entropy criteria; Stops when no more gain available
- C4.5
  - Normalised entropy criteria; Stops depending on leaf size; Incorporates pruning
- C5.0
  - Some updates on C4.5
- CART
  - Gini criteria; Cost-complexity Pruning; Surrogate predicates for missing data;
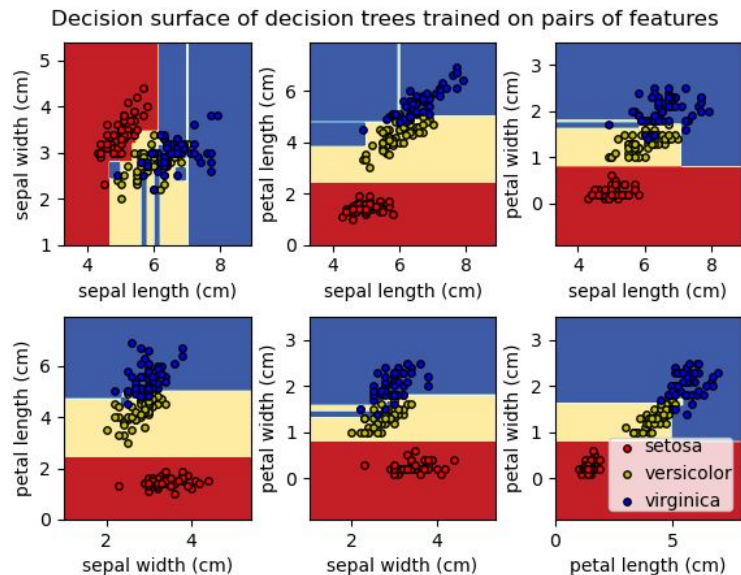- etc.

# Hyperparameters

- max_depth: min 1
- min_samples_split: min 2
- min_samples_leaf: min 1
- min_impurity_decrease

Minor

- criterion:
  - gini, entropy, log_loss for classification
  - MSE or MAE for regression
- splitter: best, random
- max_features: sqrt, log2

As of sklearn implementation



Decision surface of decision trees trained on pairs of features

# Bootstrap and Bagging

girafe
ai

06

# Bootstrap

Consider dataset X containing m objects.

Pick m objects with return from X and repeat in N times to get N datasets.

Error of model trained on Xj: $\varepsilon_j(x) = b_j(x) - y(x), \qquad j = 1, \ldots, N,$

Then $\mathbb{E}_x(b_j(x) - y(x))^2 = \mathbb{E}_x \varepsilon_j^2(x).$

The mean error of N models: $E_1 = \dfrac{1}{N} \sum_{j=1}^{N} \mathbb{E}_x \varepsilon_j^2(x).$

# Bootstrap

Consider the errors unbiased and uncorrelated:

$$\mathbb{E}_x \varepsilon_j(x) = 0;$$
$$\mathbb{E}_x \varepsilon_i(x) \varepsilon_j(x) = 0, \quad i \neq j.$$

$$E_N = \mathbb{E}_x \left( \frac{1}{N} \sum_{j=1}^{n} b_j(x) - y(x) \right)^2 =$$

$$= \mathbb{E}_x \left( \frac{1}{N} \sum_{j=1}^{N} \varepsilon_j(x) \right)^2 =$$

The final model averages all predictions:

$$a(x) = \frac{1}{N} \sum_{j=1}^{N} b_j(x).$$

$$= \frac{1}{N^2} \mathbb{E}_x \left( \sum_{j=1}^{N} \varepsilon_j^2(x) + \underbrace{\sum_{i \neq j} \varepsilon_i(x) \varepsilon_j(x)}_{=0} \right) =$$

Error decreased by N times!

$$= \frac{1}{N} E_1.$$

# Bootstrap

Consider the errors ~~unbiased and uncorrelated~~:

This is a lie

$$\mathbb{E}_x \varepsilon_j(x) = 0;$$

$$\mathbb{E}_x \varepsilon_i(x) \varepsilon_j(x) = 0, \quad i \neq j.$$

$$E_N = \mathbb{E}_x \left( \frac{1}{N} \sum_{j=1}^{n} b_j(x) - y(x) \right)^2 =$$

$$= \mathbb{E}_x \left( \frac{1}{N} \sum_{j=1}^{N} \varepsilon_j(x) \right)^2 =$$

The final model averages all predictions:

$$a(x) = \frac{1}{N} \sum_{j=1}^{N} b_j(x).$$

$$= \frac{1}{N^2} \mathbb{E}_x \left( \sum_{j=1}^{N} \varepsilon_j^2(x) + \underbrace{\sum_{i \neq j} \varepsilon_i(x) \varepsilon_j(x)}_{=0} \right) =$$
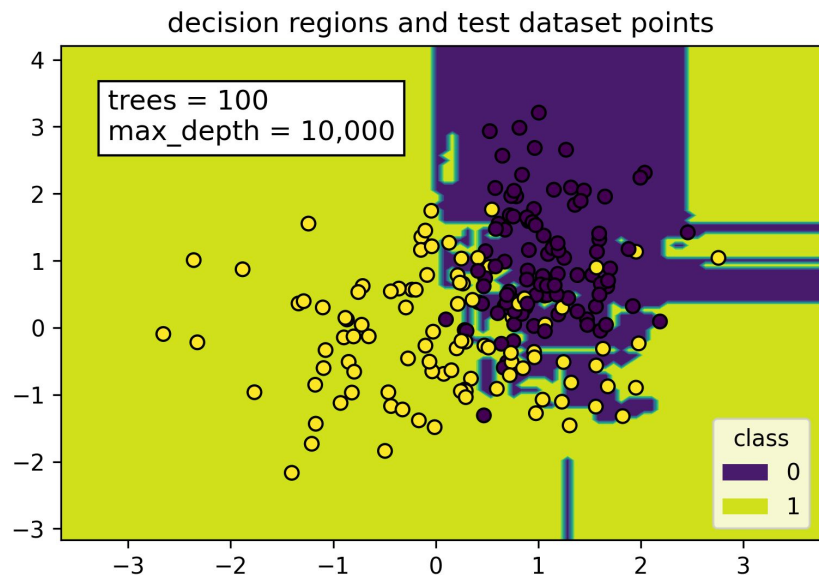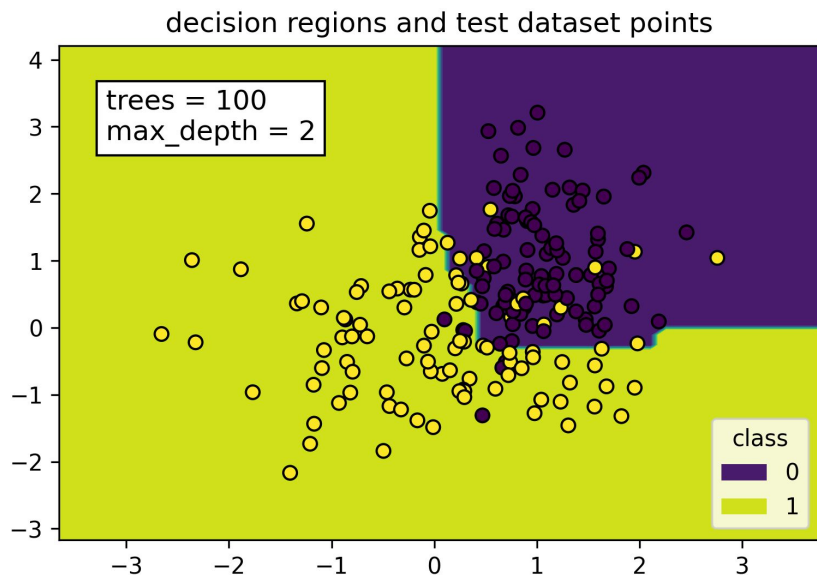
Error decreased by N times!

$$= \frac{1}{N} E_1.$$

# Bagging = Bootstrap aggregating

Decreases the variance if the basic algorithms are not correlated

# Bagging overfitting


decision regions and test dataset points

trees = 100
max_depth = 2

class
0
1

decision regions and test dataset points

trees = 100
max_depth = 10,000

class
0
1

https://www.tommykhoo.com/machine-learning/rf_overfit/
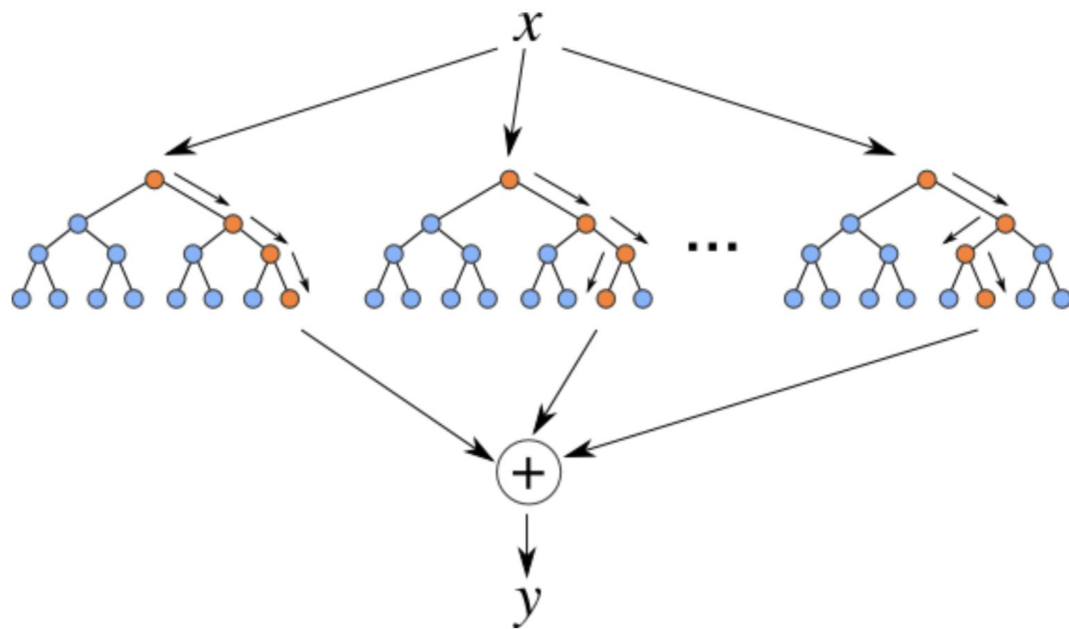
# Random Forest

girafe
ai

# RSM - Random Subspace Method

Same approach, but with features

# Random Forest
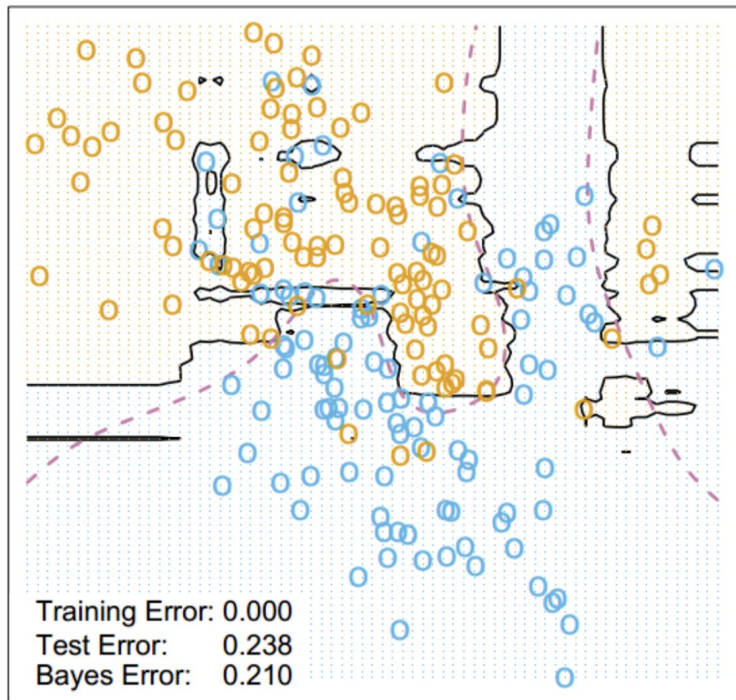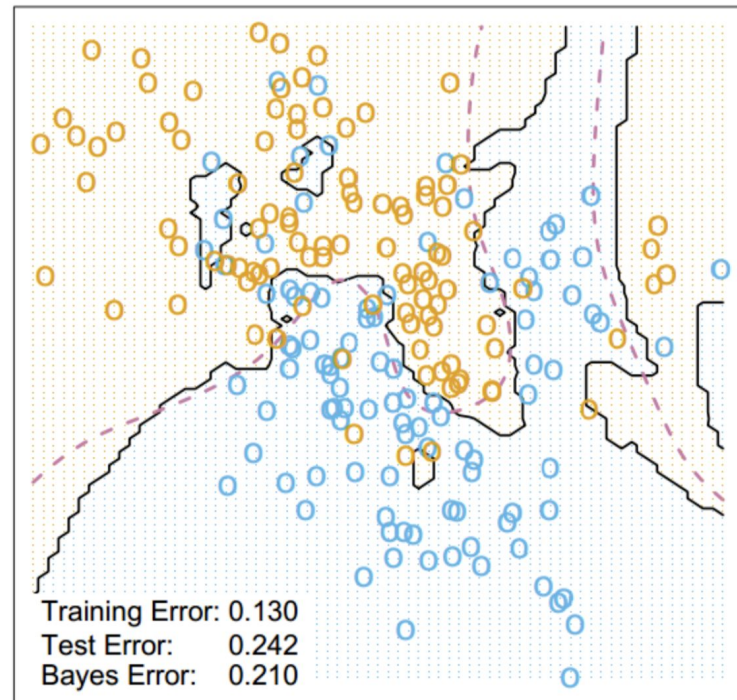
Bagging + RSM = Random Forest

# Random Forest

- One of the greatest "universal" models
- There are some modifications: Extremely Randomized Trees, Isolation Forest, etc

## Random Forest Classifier

Training Error: 0.000
Test Error: 0.238
Bayes Error: 0.210

## 3−Nearest Neighbors

Training Error: 0.130
Test Error: 0.242
Bayes Error: 0.210

# Revise

1. Intuition
2. Construction procedure
3. Information criteria
4. Decision trees special highlights
   - Decision tree as linear model
   - Dealing with missing data
   - Categorical features
5. Bootstrap and Bagging
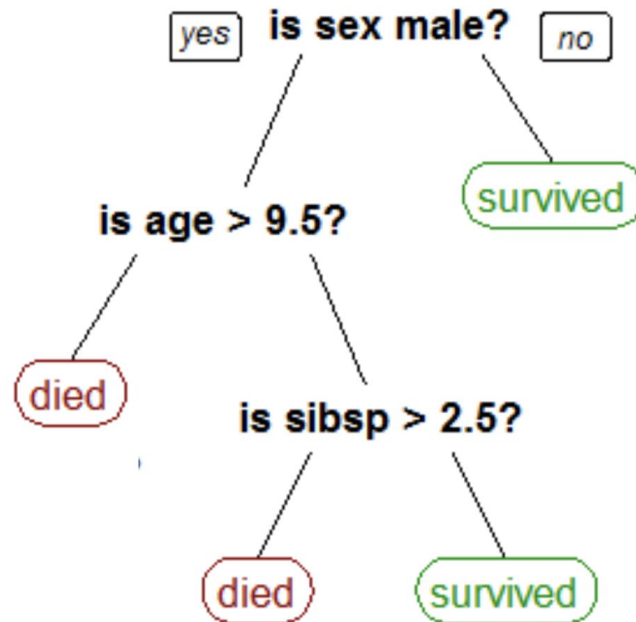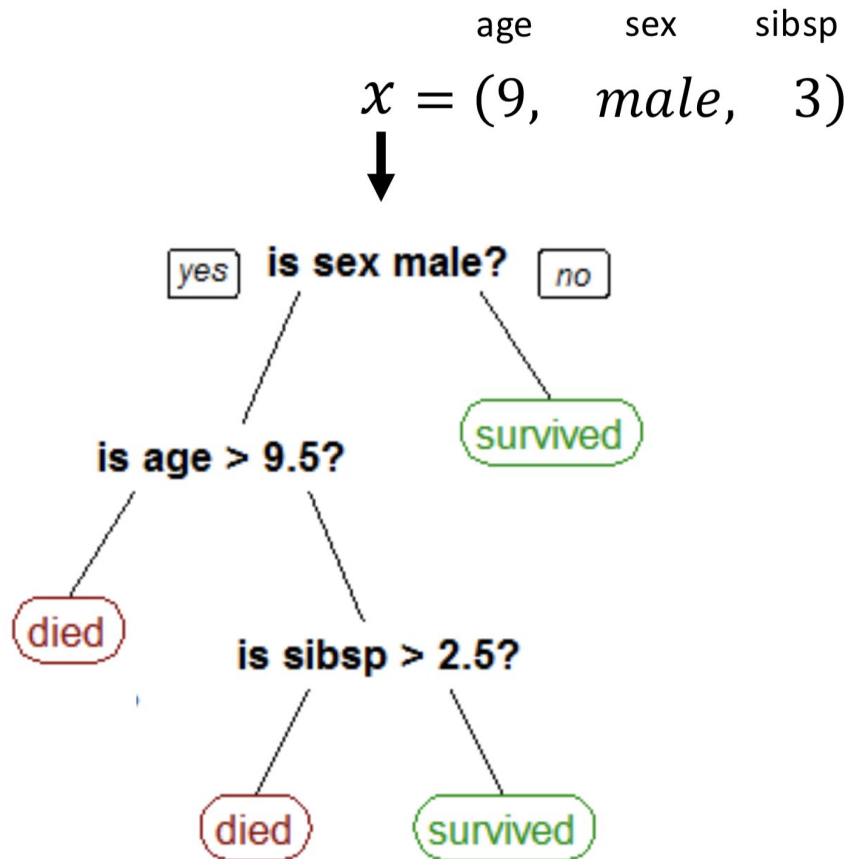6. Random Forest

# Thanks for attention!
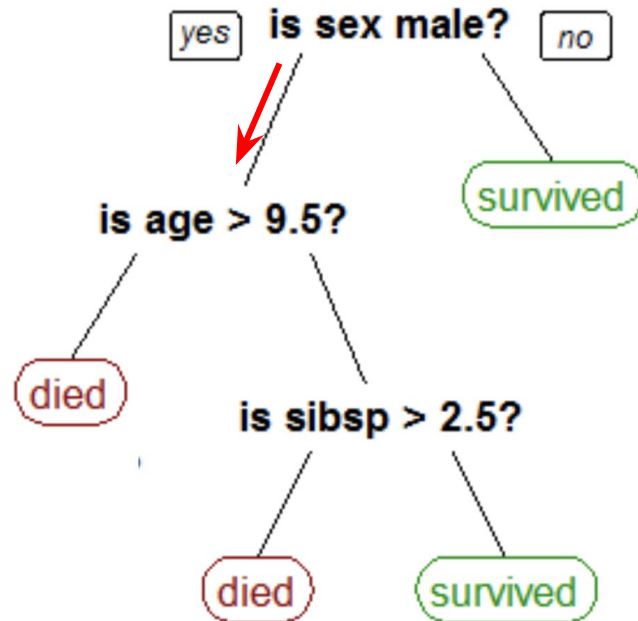
Questions?

girafe
ai

# Decision tree

$$x = (9, \quad male, \quad 3)$$



is sex male?

yes / no

is age > 9.5?

survived

died

is sibsp > 2.5?

died

survived

# Decision tree

$$x = (\underset{age}{9}, \quad \underset{sex}{male}, \quad \underset{sibsp}{3})$$



**is sex male?** yes / no

yes → **is age > 9.5?**

no → survived

is age > 9.5? → died / **is sibsp > 2.5?**

is sibsp > 2.5? → died / survived

# Decision tree

# Decision tree

$$x = (\underset{\text{age}}{9}, \quad \underset{\text{sex}}{male}, \quad \underset{\text{sibsp}}{3})$$

is sex male?

yes    no

survived

is age > 9.5?

died

is sibsp > 2.5?

died    survived

# Decision tree



age      sex      sibsp

$$x = (9, \quad male, \quad 3)$$

is sex male?

yes    no

is age > 9.5?

survived

died

is sibsp > 2.5?

died      survived

49

# Decision tree



$$x = (9, \quad male, \quad 3)$$

age     sex     sibsp

**is sex male?**

yes    no

**is age > 9.5?**    survived

died    **is sibsp > 2.5?**

died    survived

$$y = died$$

# Decision tree in classification

# Decision tree in classification



is sex male?

yes · no

is age > 9.5?

survived
0.73  36%

died
0.17  61%

is sibsp > 2.5?

died
0.05  2%

survived
0.89  2%

Part of survived passengers
= part of objects of 1st class

# Decision tree in classification



is sex male?
yes / no

is age > 9.5?

survived
0.73  36%

died
0.17  61%

is sibsp > 2.5?

died
0.05  2%

survived
0.89  2%

Part of dataset in this leaf
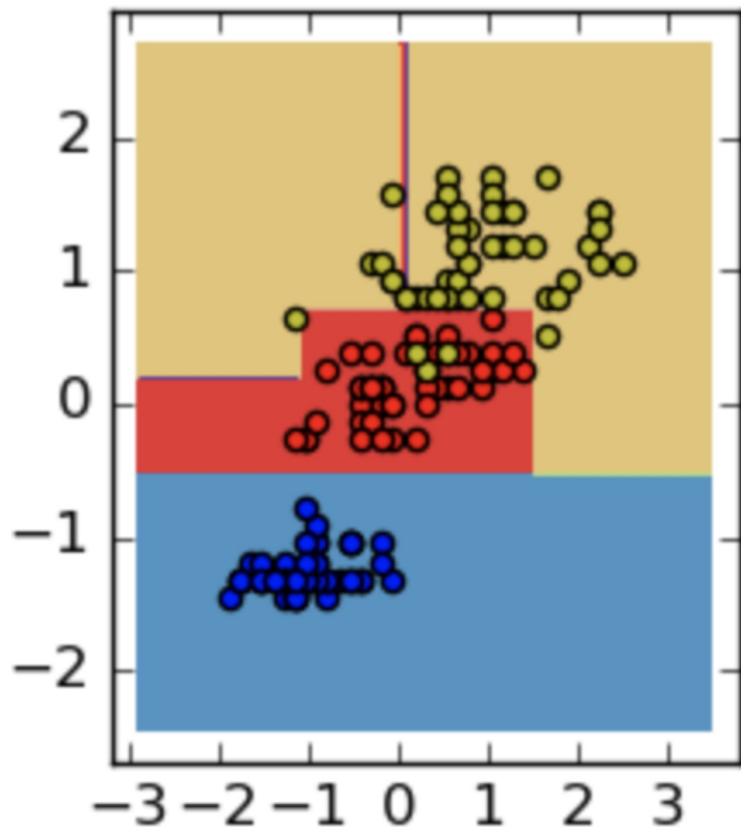
Part of survived passengers = part of objects of 1st class

# Decision tree in classification



Classification problem with 3 classes and 2 features.