# Introduction to Artificial Intelligence: Methods, Models, Algorithms

**Aleksandr I. Panov and Konstantin S. Yakovlev**

National Research University Higher School of Economics

20 July 2018 – Summer University

apanov@hse.ru

Intro
○●○○○○○○○○

Definition
○○○○○○○○○○○○○

Tree learning
○○○○○○○○○○○○○○○○○○○○○○○○○○○

## Linear models

$$a(x) = w_0 + \sum_{j=1}^{d} w_j x^j$$

Weights can be interpreted if features are scaled

# Example

- The prediction value of the apartment
- Features: area, floor, number of rooms

$$a(x) = 10 \cdot (\text{area}) + 1.1 \cdot (\text{floor}) + 20 \cdot (\text{number of rooms})$$

# Example

- Dependence on the floor is hardly linear
- Quadratic features:

$a(x) = 10 \cdot (\text{area}) + 1.1 \cdot (\text{floor}) + 20 \cdot (\text{number of rooms}) - 0.2 \cdot (\text{area})^2 + 0.5 \cdot (\text{area} \cdot \text{number of rooms}) + \ldots$

# Example

- With cubic features will be even better
- How to interpret the feature
  area $\cdot$ number of rooms$^2$?
- A total of 20 such features

Intro
○○○○●○○○○○

Definition
○○○○○○○○○○○○○

Tree learning
○○○○○○○○○○○○○○○○○○○○○○○○

# Example

- You can binarysoul features: $[x^j > t]$
- $(\text{floor} > 1), (\text{floor} > 2), ..., (\text{floor} > 30)$
- Features will be orders of magnitude more
- Easier to interpret:
  $-2[\text{floor} > 3][\text{area} < 40][\text{number of rooms} < 3]$
- You can use $L_1$-regularization

Intro
0000000000

Definition
000000000000

Tree learning
00000000000000000000000

## Logical rules

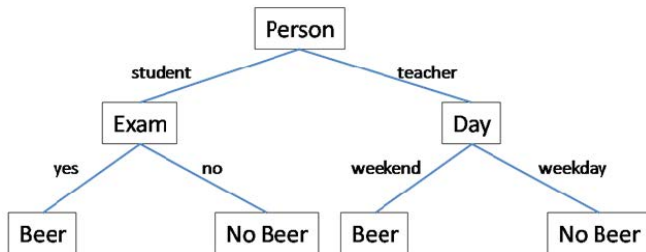[floor $> 3$][area $< 40$][number of rooms $< 3$]

- Easy to explain to the customer (if $\leq 5$ conditions)

- Allow you to extract knowledge from data

- Not the fact that they are optimal in terms of quality

Intro
0000000●000
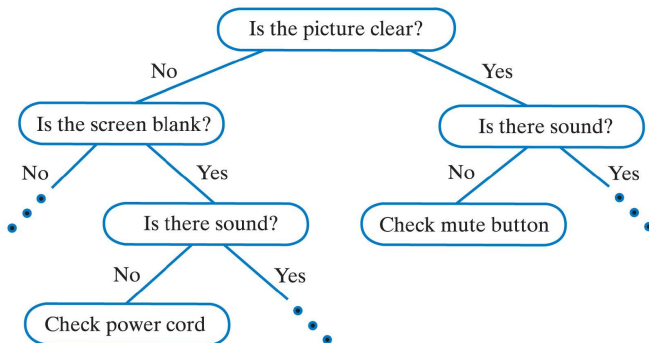
Definition
000000000000

Tree learning
0000000000000000000000000

# Logical rules

- How to construct them?

- Linear model

- Busting, greedy build-up

- Decision trees

Intro
○○○○○○○●○○

Definition
○○○○○○○○○○○○

Tree learning
○○○○○○○○○○○○○○○○○○○○○○○○

# Decision making

Intro
○○○○○○○○○●○

Definition
○○○○○○○○○○○○○

Tree learning
○○○○○○○○○○○○○○○○○○○○○○○○○

# The scheme of dialogue with the client

Intro
○○○○○○○○○●

Definition
○○○○○○○○○○○○○

Tree learning
○○○○○○○○○○○○○○○○○○○○○○○○○○

## The passengers of the Titanic

Intro
ooooooooooo

Definition
●oooooooooooo

Tree learning
oooooooooooooooooooooooo

# Decision tree

- Binary tree
- Each inner node contains a condition
- Each leaf contains prediction (solution)

Intro
○○○○○○○○○○○

Definition
○●○○○○○○○○○○○

Tree learning
○○○○○○○○○○○○○○○○○○○○○○○○○○

# Conditions

- Most popular options:

$$[x^j \leq t] \text{ and } [x^j = t]$$

- Examples:
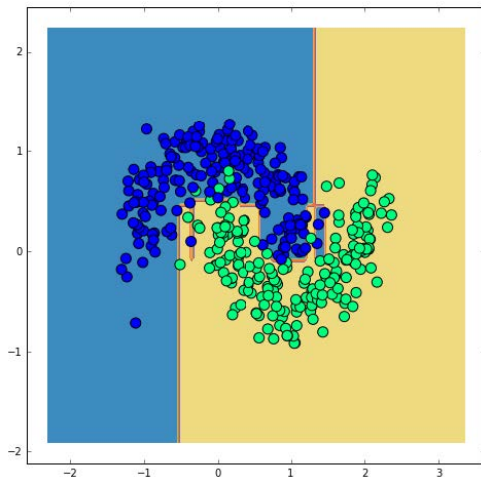
$$[\text{floor} = 5] \text{ or } [\text{area} \leq 30]$$

Intro
0000000000

Definition
000●000000000

Tree learning
0000000000000000000000

# Prediction in the leaf

- Regression: Real number
- Classification: Class or Class probabilities

Intro
○○○○○○○○○○○

Definition
○○○●○○○○○○○○○

Tree learning
○○○○○○○○○○○○○○○○○○○○○○○○○

# Classification

Intro
0000000000

Definition
0000●00000000

Tree learning
00000000000000000000000

# Classification

Intro
0000000000

Definition
000000●0000000

Tree learning
0000000000000000000000000

# Classification

Intro
0000000000

Definition
000000●000000

Tree learning
0000000000000000000000000

# Classification

Intro
○○○○○○○○○○

Definition
○○○○○○○●○○○○○

Tree learning
○○○○○○○○○○○○○○○○○○○○○○○○○○○

# Regression

Intro
○○○○○○○○○○

Definition
○○○○○○○○●○○○○

Tree learning
○○○○○○○○○○○○○○○○○○○○○○○○○○

# Regression

Intro
0000000000
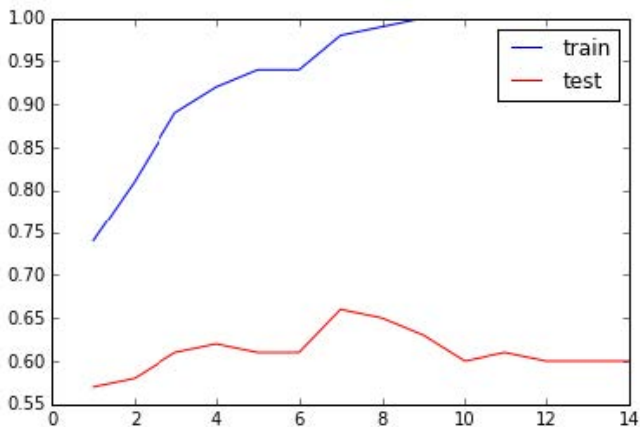
Definition
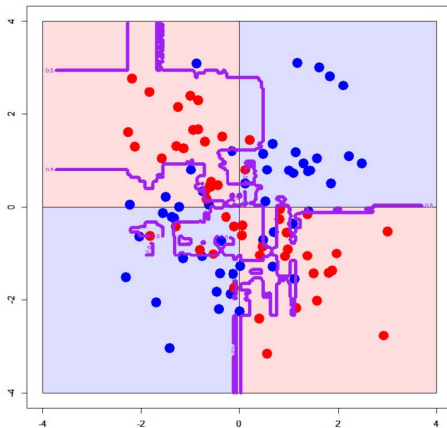000000000●000

Tree learning
0000000000000000000000000

# Decision tree

- Restore complex dependencies
- Can build any complex surface
- The greater the depth the more complex the surface
- Prone to overfitting

# Depth of trees

Intro
0000000000

Definition
000000000000●0

Tree learning
0000000000000000000000000

# Overfitting of trees

Intro
0000000000

Definition
0000000000000●
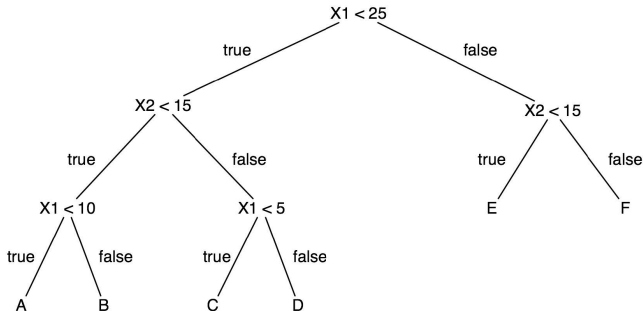
Tree learning
0000000000000000000000

# Overfitting of trees

- The tree can achieve zero error on any sample
- Tackling overfitting: the minimum tree among all with zero error
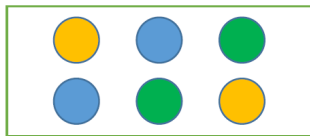- NP-complete task

- *Solution*: greedy building

Intro
0000000000

Definition
0000000000000

Tree learning
●0000000000000000000000000

## Gready formation

- Grow the tree from root to leaves

Intro
○○○○○○○○○○

Definition
○○○○○○○○○○○○○

Tree learning
○●○○○○○○○○○○○○○○○○○○○○○○○○

# Gready formation



How to split the node?

Intro
○○○○○○○○○○○

Definition
○○○○○○○○○○○○○

Tree learning
○○●○○○○○○○○○○○○○○○○○○○○○○○○

# Gready formation

Intro
○○○○○○○○○○

Definition
○○○○○○○○○○○○○

Tree learning
○○○●○○○○○○○○○○○○○○○○○○○○○○

# Gready formation

Intro
0000000000

Definition
0000000000000

Tree learning
0000●0000000000000000000

# Gready formation

Intro
○○○○○○○○○○○

Definition
○○○○○○○○○○○○○

Tree learning
○○○○○●○○○○○○○○○○○○○○○○○○○○

# Gready formation

Intro
0000000000

Definition
0000000000000

Tree learning
0000000●0000000000000000000

# Gready formation

Intro
○○○○○○○○○○

Definition
○○○○○○○○○○○○○

Tree learning
○○○○○○○○●○○○○○○○○○○○○○○○

# How to compare splits?



OR

Intro
○○○○○○○○○○

Definition
○○○○○○○○○○○○○

Tree learning
○○○○○○○○●○○○○○○○○○○○○○○○

# Entropy

- Measure of uncertainty of distribution

Intro
0000000000

Definition
000000000000

Tree learning
00000000000●000000000000000

# Entropy
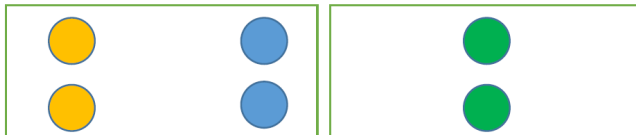
- Discrete distribution
- Accepts $n$ values with probabilities $p_1, \ldots, p_n$
- Entropy:

$$H(p_1, \ldots, p_n) = -\sum_{i=1}^{n} p_i \log p_i$$

Intro
oooooooooo

Definition
ooooooooooooo

Tree learning
oooooooooooo●ooooooooooooo
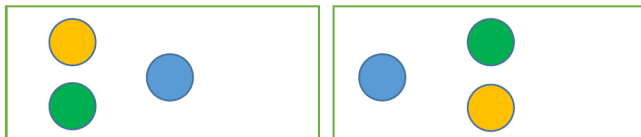
# Entropy

- $(0.2, 0.2, 0.2, 0.2, 0.2) \rightarrow H = 1.60944$
- $(0.9, 0.05, 0.05, 0, 0) \rightarrow H = 0.394398$
- $(0, 0, 0, 1, 0) \rightarrow H = 0$

Intro
○○○○○○○○○○

Definition
○○○○○○○○○○○○○

Tree learning
○○○○○○○○○○○○●○○○○○○○○○○○

# Entropy



- (0.5, 0.5, 0) and (0, 0, 1)
- $H = 0.693 + 0 = 0.693$



- (0.33, 0.33, 0.33) and (0.33, 0.33, 0.33)
- $H = 1.09 + 1.09 = 2.18$

Intro
0000000000
Definition
000000000000
Tree learning
00000000000000●000000000000

# What about regression?

Intro
0000000000

Definition
0000000000000

Tree learning
000000000000000●000000000

# What about regression?

Intro
0000000000

Definition
0000000000000

Tree learning
0000000000000000●00000000

# What about regression?

Intro
0000000000

Definition
0000000000000

Tree learning
00000000000000000●0000000

# What about regression?

- Choose the partition with the least total variance
- The smaller the variance, the less uncertainty

Intro
000000000

Definition
0000000000000

Tree learning
0000000000000000000000000

## Searching the partition

- Let the node $m$ be the sample $X_m$
- $Q(X_m, j, t)$ - is a criteria for the condition error $[x^j \leq t]$
- Search the parameters $t$ and $j$:

$$Q(X_m, t, j) \to \min_{j,t}$$

## Searching the partition

- When we find the partition we split the $X_m$ into two parts:

$$X_l = \{x \in X_m | [x^j \leq t]\}$$
$$X_r = \{x \in X_m | [x^j > t]\}$$

- Repeat the procedure for child nodes

Intro
000000000

Definition
0000000000000

Tree learning
00000000000000000000●0000

# Stop criterion

- At what point should the splitting of nodes be stopped?

- The single item at the node of ?

- Items of the same class at the node?

- Did the depth exceed a threshold?

Intro
0000000000

Definition
0000000000000

Tree learning
00000000000000000000●000

## Prediction in the leaf

- For example, I decided to make a node $m$ leaf
- Which prediction to choose?
- Regression:

$$a_m = \frac{1}{|X_m|} \sum_{i \in X_m} y_i$$

- Classification

$$a_m = \arg\max_{y \in \mathbb{Y}} \sum_{i \in X_m} [y_i = y]$$

## Prediction in the leaf

- For example, I decided to make a node $m$ leaf
- Which prediction to choose?
- Class probabilities:

$$a_{mk} = \frac{1}{|X_m|} \sum_{i \in X_m} [y_i = k]$$

Intro
0000000000

Definition
000000000000

Tree learning
0000000000000000000000●0

# Summary

- Sometimes the model needs to be interpreted

- Decision trees are easy to explain

- Decision trees easily overfits

- Tree construction is a greedy algorithm