# AUTOMATIZATION
# OF ANALYSES IN RAPID MINER

Petr Berka
University of Economics, Prague
berka@vse.cz

Seminar

Vysoká škola manažmentu v Trenčíne
International Workshop on Knowledge Management

IWKM'2019

**November, 7– 8**
**Bratislava 2019**

# Automatization of analyses in Rapid Miner

Petr Berka
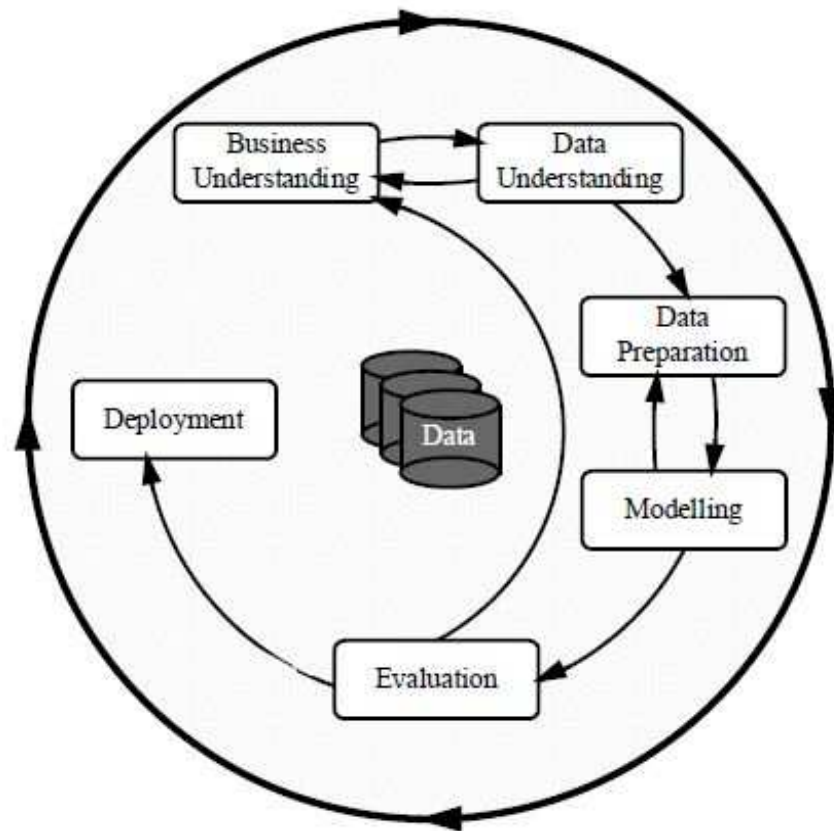
(berka@vse.cz)

University of Economics
Prague

# Automatization of KDD process...

... will allow the domain experts (knowledge workers) to perform data mining tasks without the cooperation with data mining experts, in a "do-it-yourself" way.

As the authors of the DataRobot platform believe: "Automated machine learning creates a new class of citizen data scientists with the power to create advanced machine learning models, all without having to learn to code or understand when and how to apply certain algorithms".

# KDD Process and ...



- **Business Understanding** is the initial phase that focuses on understanding the project objectives
- **Data Understanding** proceeds with activities in order to get familiar with the data
- **Data Preparation** phase covers all activities to construct the final dataset
- In the **Modeling** phase various modeling techniques are applied
- At the **Evaluation** stage the quality of the created model is assessed
- the **Deployment** phase can be as simple as generating a report or as complex as implementing a repeatable data mining process

# ... Possibilities for its Automatization

- **Business Understanding** and **Deployment** steps are closely related to the application domain so its automatization in a general way would be very difficult

- **Data Understanding** is already supported by computing basic characteristics of the data

- data mining automatization is oriented towards supporting the **Data Preparation** (preprocessing) and **Modeling** (learning) steps

# Rapid Miner (rapidminer.com)

- a leading open-source system for knowledge discovery and data mining (www.kdnuggets.com)

- a Leader in 2016 Gartner Magic Quadrant for Advanced Analytics (www.gartner.com)

- the Top 3 Rated Predictive Analytics Software for Enterprise (www.g2crowd.com)

# Rapid Miner Downloads

https://my.rapidminer.com/nexus/account/index.html #downloads

# Rapid Miner Pricing



| | FREE | SMALL | MEDIUM | LARGE |
|---|---|---|---|---|
| | Free | $2,500 Yearly | $5,000 Yearly | $10,000 Yearly |
| # Data Rows | 10,000 | 100,000 | 1,000,000 | Unlimited |
| # Logical Processors | 1 | 2 | 4 | Unlimited |
| Performance Improvements | | 2x | 4x | 10x+ |
| Background Process Execution | | | | ✔ |
| Customer Support | Community | Enterprise | Enterprise | Enterprise |

# Overview of a DM Project

# „Standard" modeling using Rapid Miner

# Automatization no. 1: Parameter Optimization

# Automatization no. 1: Parameter Optimization

# Automatization no. 1: Parameter Optimization

Optimize Parameters (Grid) (2002 rows, 6 columns)

| iteration | k-NN.k | k-NN.weighted_vote | k-NN.nominal_measure | k-NN.numerical_measure | accuracy ↓ |
|-----------|--------|--------------------|-----------------------|--------------------------|------------|
| 1284 | 8 | true | JaccardSimilarity | JaccardSimilarity | 0.738 |
| 1821 | 6 | false | RussellRaoSimilarity | MaxProductSimilarity | 0.737 |
| 973 | 5 | true | JaccardSimilarity | DynamicTimeWarpingDistance | 0.736 |
| 323 | 4 | false | NominalDistance | ChebychevDistance | 0.731 |
| 1215 | 5 | true | SimpleMatchingSimilarity | InnerProductSimilarity | 0.730 |
| 319 | 11 | true | NominalDistance | ChebychevDistance | 0.729 |
| 1338 | 7 | false | RogersTanimotoSimilarity | JaccardSimilarity | 0.728 |
| 1203 | 4 | false | RussellRaoSimilarity | InnerProductSimilarity | 0.728 |
| 1039 | 5 | true | RussellRaoSimilarity | DynamicTimeWarpingDistance | 0.728 |
| 1056 | 11 | false | RussellRaoSimilarity | DynamicTimeWarpingDistance | 0.727 |
| 96 | 8 | true | RogersTanimotoSimilarity | EuclideanDistance | 0.726 |
| 801 | 9 | true | DiceSimilarity | DiceSimilarity | 0.724 |
| 984 | 5 | false | JaccardSimilarity | DynamicTimeWarpingDistance | 0.723 |

# Automatization no. 2: Auto Model

- **Auto Model** finds the best model using multiple machine learning algorithms and hyperparameter optimization

# Auto Model step 1

# Auto Model step 2

# Auto Model step 3

# Auto Model step 4

# Auto Model step 5

# Auto Model step 6

# Auto Model (web/cloud) version

# Automatization no. 3: Turbo Prep

- **Turbo Prep** aims at intuitive data preparation. This extension allows to interactively explore and visualize the data, simplifies data cleansing (automatically removes low quality and correlated data columns) and merges multiple datasets together by automatically identifying matching columns to merge.

# Turbo Prep step 1

# Turbo Prep step 2

# Turbo Prep step 3

# Turbo Prep step 4

# Turbo Prep step 5

# Automatization of analyses in Rapid Miner

- Optimized parameters node always available

- Auto Model and Turbo Prep available in paid version and within free academic license

# Resources

1. Get Started with Fully Automated Data Science, 2019. *Rapid Miner.* [online] Available at: <httphttps://rapidminer.com/products/#automated> [Accessed 1 October 2019].

2. HOFMANN, M. and KLINKENBERG, R., 2013. RapidMiner: Data Mining Use Cases and Business Analytics Applications. Chapman and Hall/CRC.