School of Data Analysis and Artificial Intelligence Department of Computer Science

# DATA SCIENCE FOR BUSINESS

Lecture 1. Introduction to Data Science

Moscow, April 10th, 2020.

lzhukov@hse.ru

# COURSE TECHNICALITIES

Lectures:  Friday 18.10 - 19.30 , 10 lectures        ZOOM:  https://zoom.us/j/7723819319
Seminars: Friday 19.40 - 21.00,  10 seminars      ZOOM:  https://zoom.us/j/636910206


Class Website:          http://www.leonidzhukov.net/hse/2020/datascience
Seminar Wiki:           http://wiki.cs.hse.ru/Data_Science_for_Business_2020


Telegram Group:         https://t.me/joinchat/ENzQEhr-hra2WhEjxvgayw


Modeling software:      RapidMiner    https://rapidminer.com

# TEACHING TEAM



**Prof. Leonid Zhukov**
lzhukov@hse.ru



**Anvar Kurmukov**
kurmukovai@gmail.com
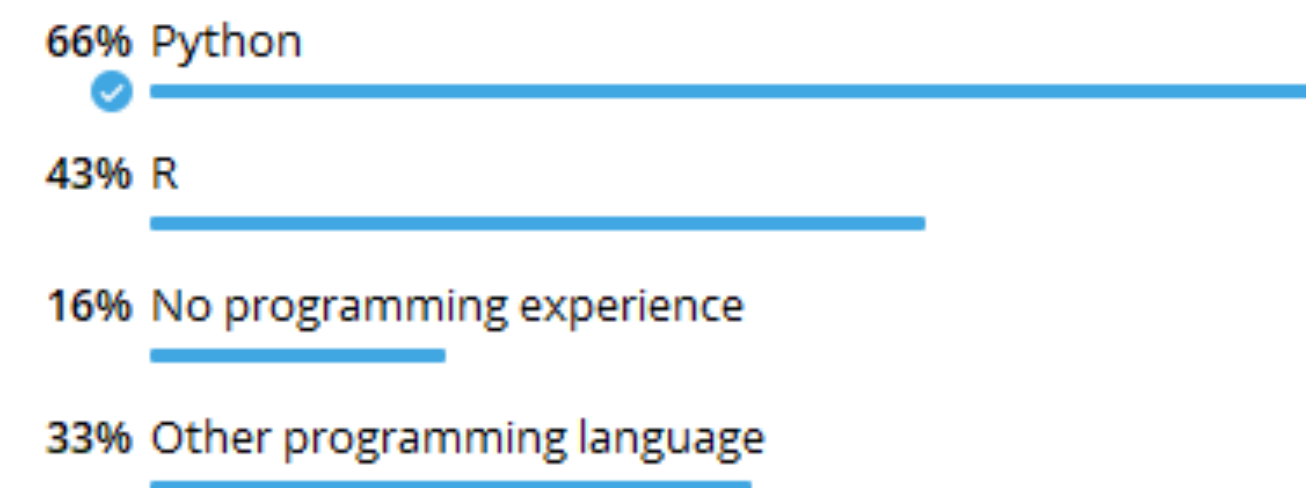


**Ilya Makarov**
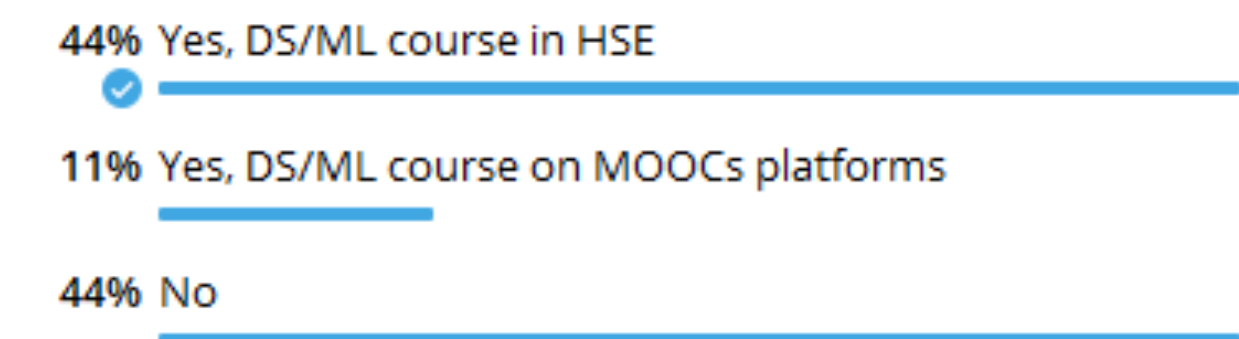iamakrov@hse.ru

# THE CLASS 2020

**What is your faculty?**

Poll

19% Computer Science

33% Business and Management

0% Faculty of Law

0% Faculty of Humanities

5% Faculty of Social Science

2% Faculty of Communication, Media and Design

5% Faculty of World Economy and International Affairs

28% Faculty of Economic Sciences

0% Faculty of Urban and Regional Development

8% Other

**Do you have an experience with any programming language?**

Poll

66% Python

43% R

16% No programming experience

33% Other programming language

**Did you have any Data Science / Machine Learning classes?**

Poll

44% Yes, DS/ML course in HSE

11% Yes, DS/ML course on MOOCs platforms

44% No

**80-85 votes**

# TEXTBOOKS FOR THE COURSE



"A must-read resource for anyone who is serious about embracing the opportunity of big data."
—Craig Vanghan, Global Vice President, SAP

Data Science
for Business

What You Need to Know
About Data Mining and
Data-Analytic Thinking

Foster Provost & Tom Fawcett

Second Edition

Data Science
Concepts and Practice

MK

Vijay Kotu and Bala Deshpande

PREDICTIVE
ANALYTICS

THE POWER TO PREDICT WHO WILL
CLICK, BUY, LIE, OR DIE

ERIC SIEGEL

# RAPIDMINER MODELING SOFTWARE



Do you have RapidMiner installed (with educational license)?

Poll

52% I do

48% I do not

https://rapidminer.com

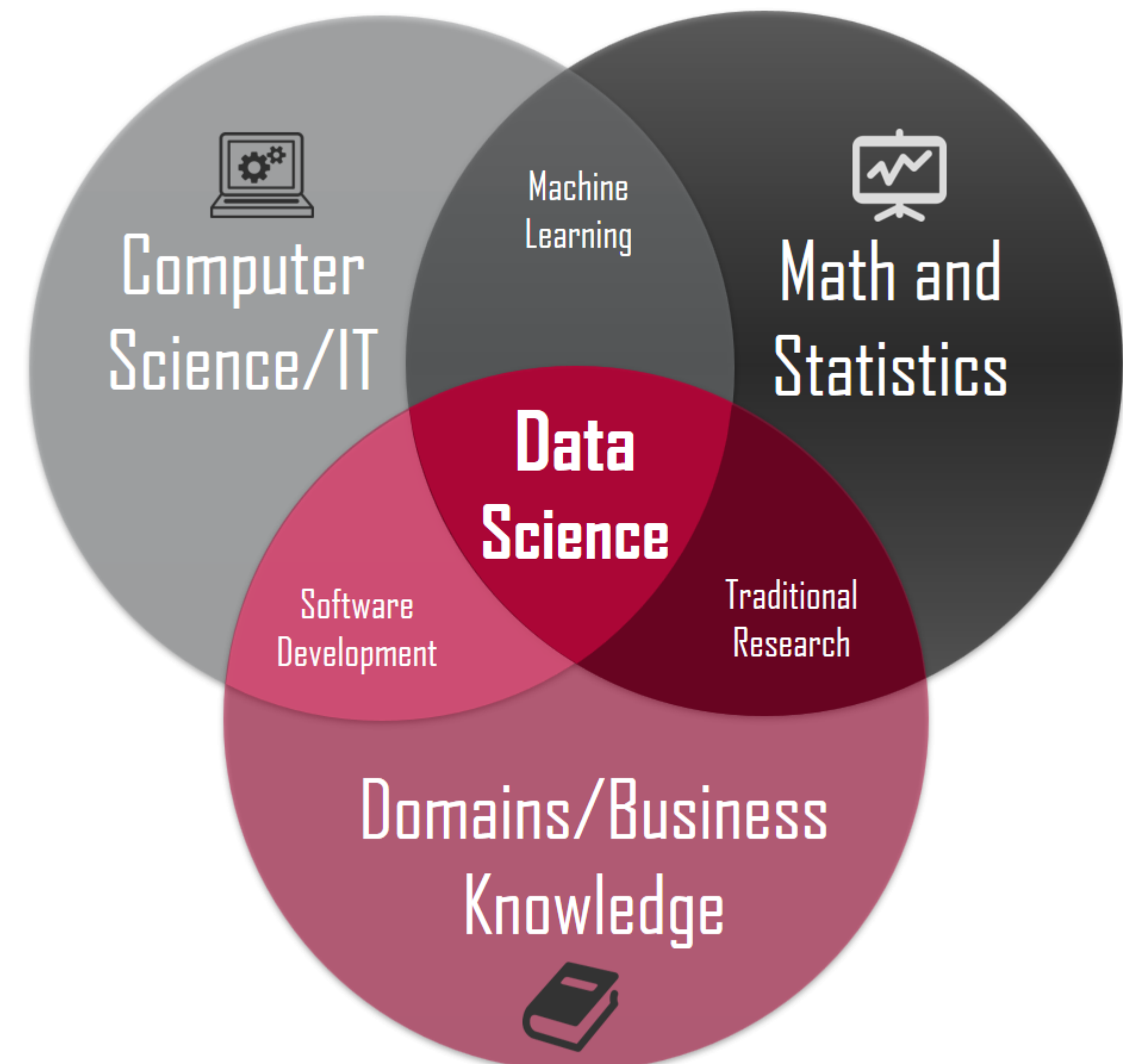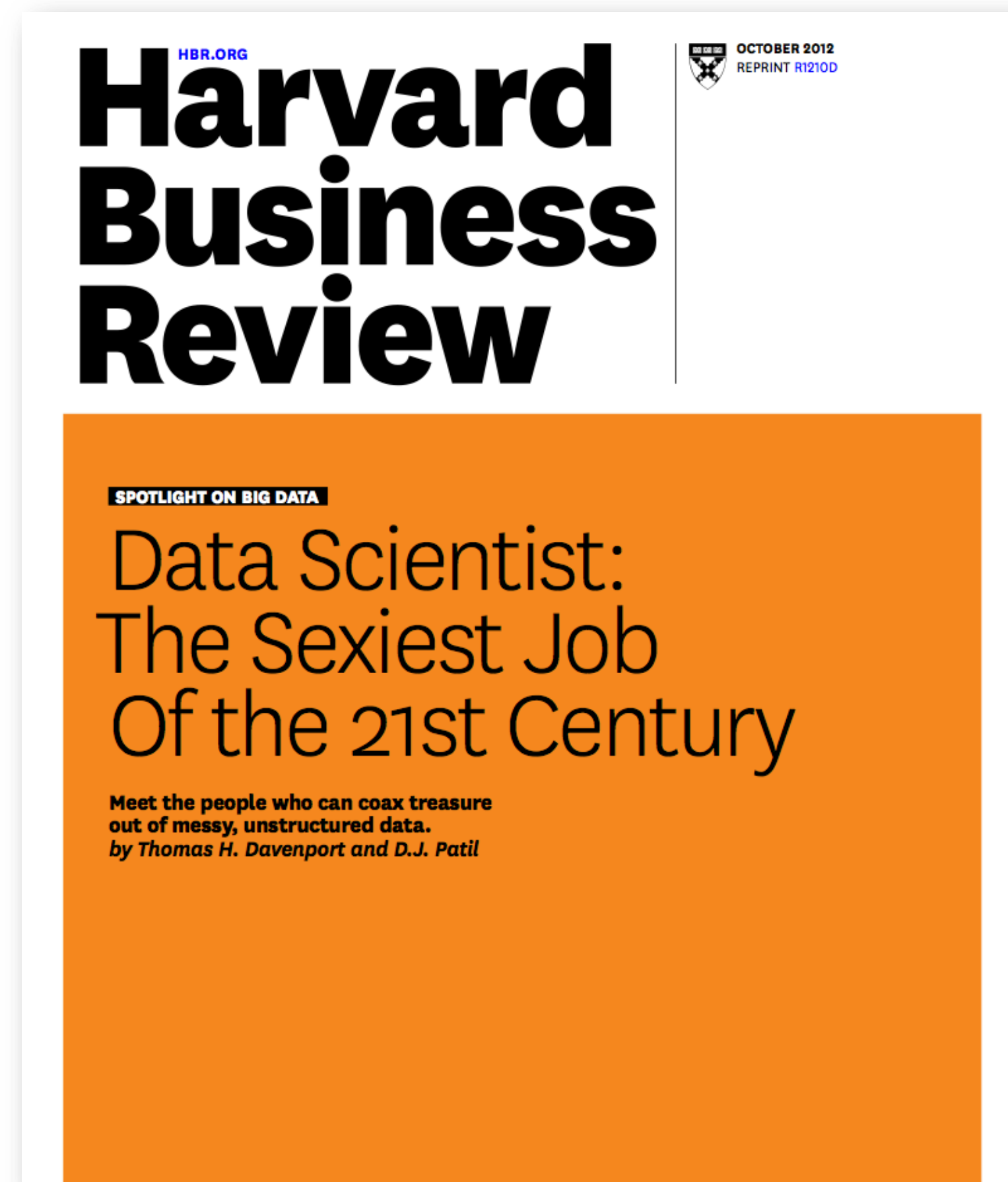# COURSE SCHEDULE

## Lecture topics

1. Introduction to data science.
2. Exploratory data analysis
3. Predictive analytics and machine learning
4. Case study 1: Retail pricing
5. Case study 2: Churn modeling
6. Case study 3: Customer segmentation
7. Case study 4: Personalization
8. Case study 5: Fraud detection
9. Case study 6. Demand forecasting
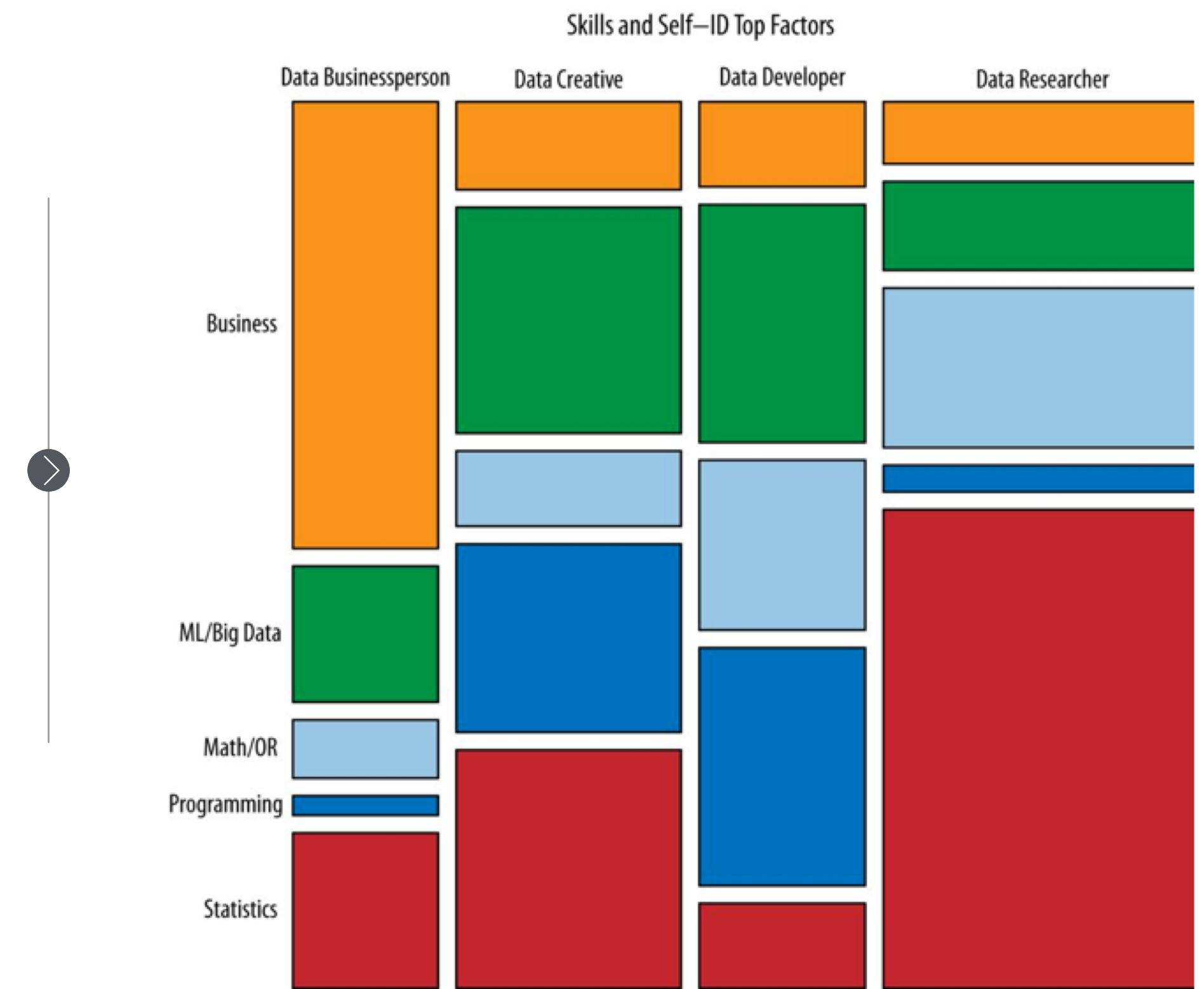10. Impacting the business

      **Exam**

## Seminars - exercises

1. Data flow modeling and RapidMiner
2. Working with data, ETL process, data exploration
3. ML modeling pipeline
4. Regression
5. Classification
6. Clustering
7. Recommender systems
8. Anomaly detection
9. Time series forecasting
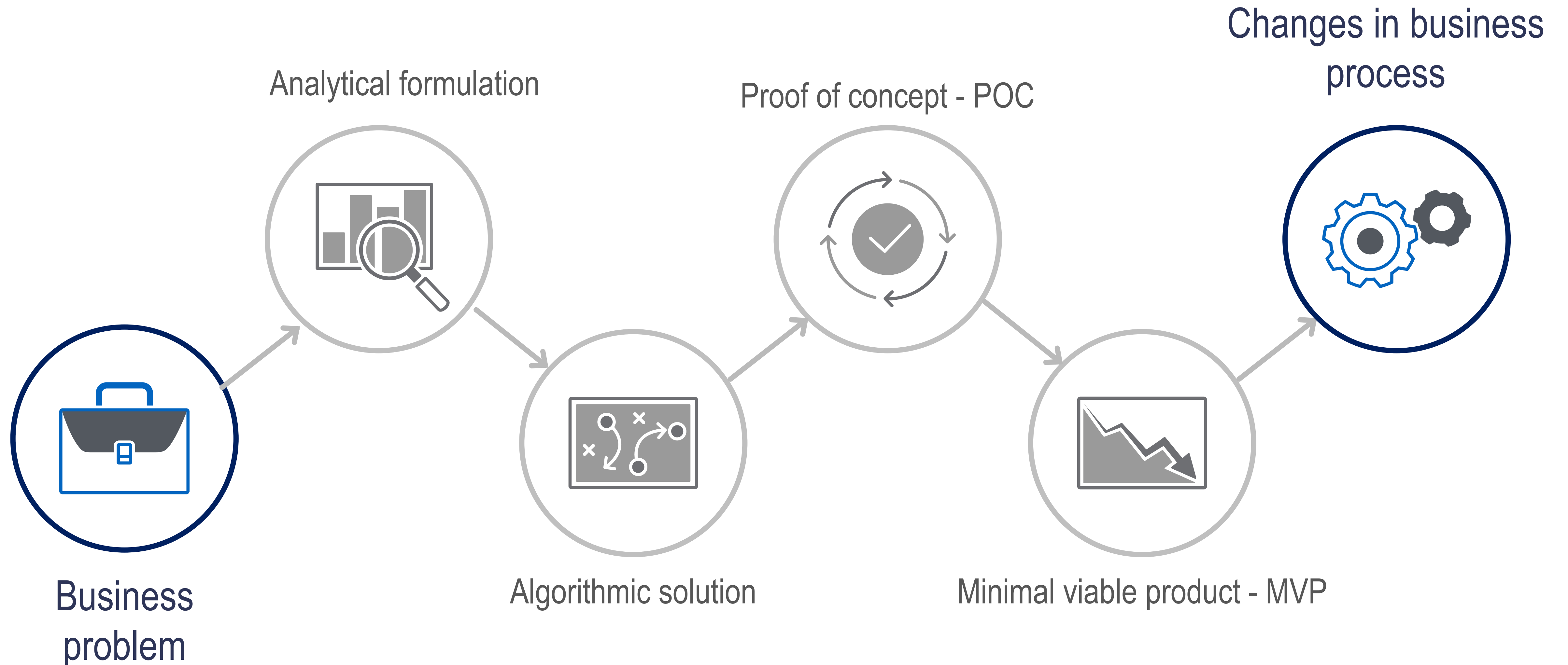10. Problem solving

# DATA SCIENCE

# DATA SCIENTISTS



**Data Scientist**
also known as Data Managers, statisticians.

A data scientist will be able to take data science projects from end to end. They can help store large amounts of data, create predictive modelling processes and present the findings.

*Skills:* Mathematics, Programming, Communication

*Will use programmes such as:*
SQL, Python, R

**Data Engineers**
also known as database administrators and data architects.

They are versatile generalists who use computer science to help process large datasets. They typically focus on coding, cleaning up data sets, and implementing requests that come from data scientists.

*Skills:* Programming, Mathematics, Big data

*Will use programmes such as:*
Hadoop, NoSQL, and Python

**Data Analysts**
also known as business Analysts.

They typically help people from across the company understand specific queries with charts.

*Skills:* Statistics, Communication, Business knowledge

*Will use programmes such as:*
Excel, Tableau, SQL

Skills and Self—ID Top Factors

Data Businessperson | Data Creative | Data Developer | Data Researcher

Business

ML/Big Data

Math/OR

Programming

Statistics

# DATA SCIENCE BUSINESS PROCESS

Changes in business process

Analytical formulation

Proof of concept - POC



Business problem

Algorithmic solution

Minimal viable product - MVP
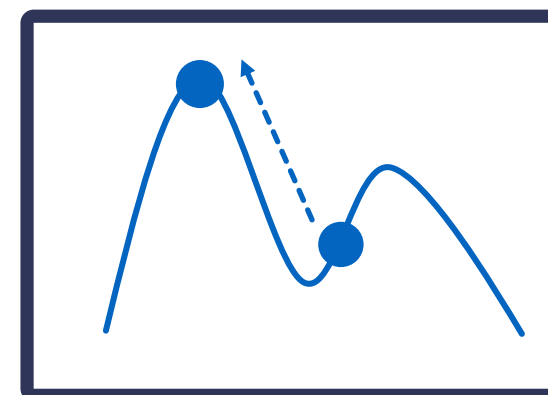
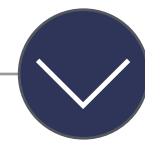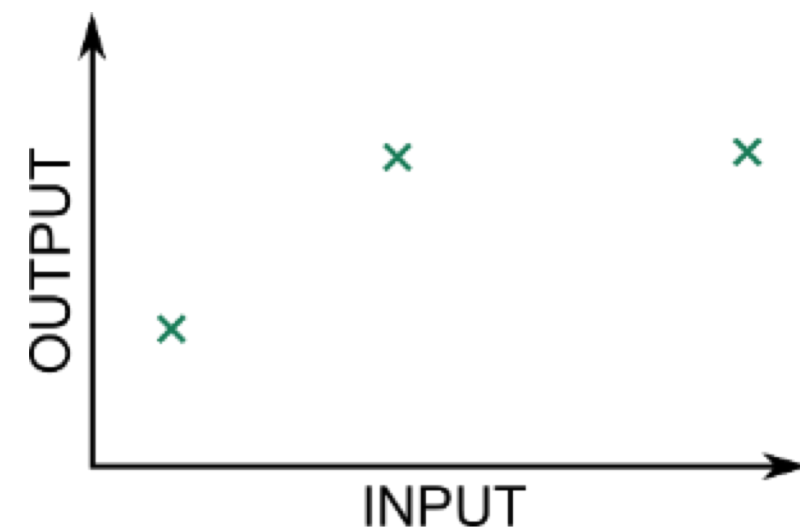# THREE MAIN REASONS TO USE ML IN BUSINESS

Detect

Predict

Explain
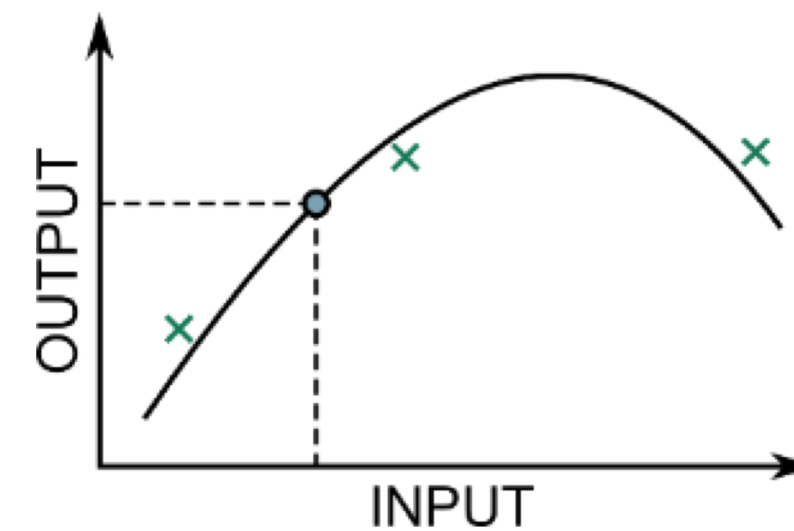
Optimise and improve

# SIMPLE EXAMPLE

**Statistical Analysis**
*"Measure and understand"*



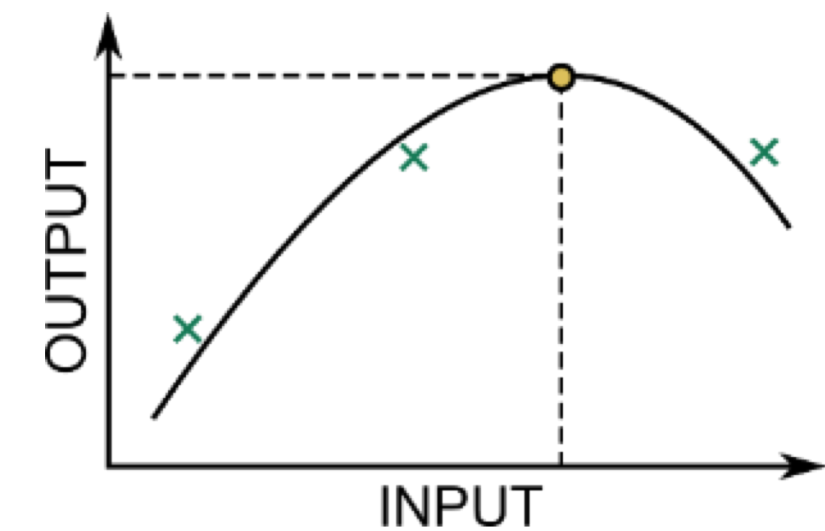Data exploration
Descriptive statistics

**Predictive modeling**
*"ML predict outcomes"*



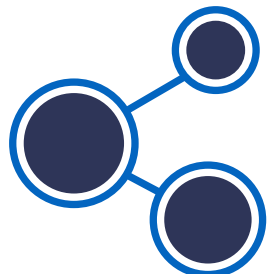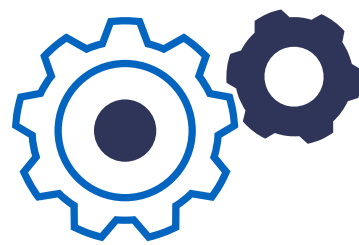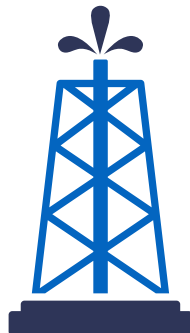Finding  patterns
Predicting outcomes
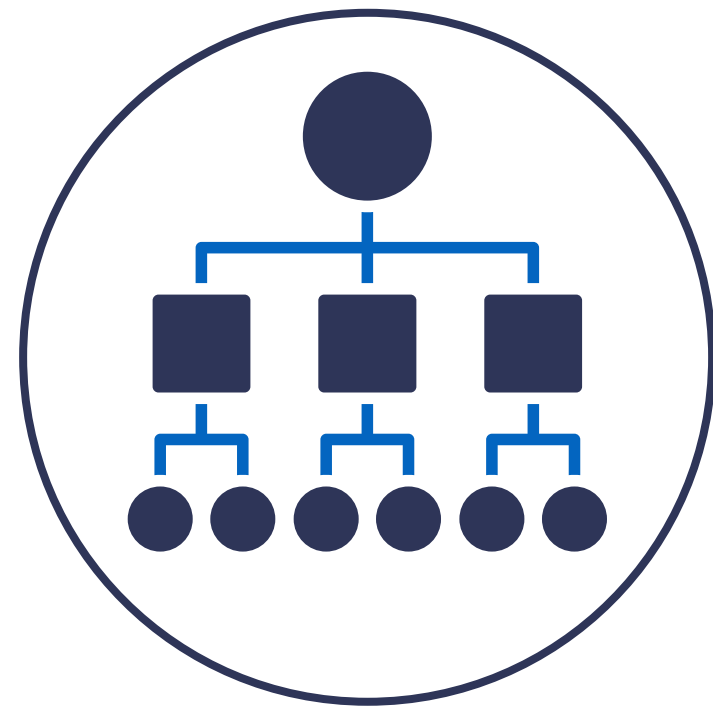
**Optimization**
*"Make optimal decisions"*



Finding optimal values
Finding optimal parameters

# BUSINESS USE CASES

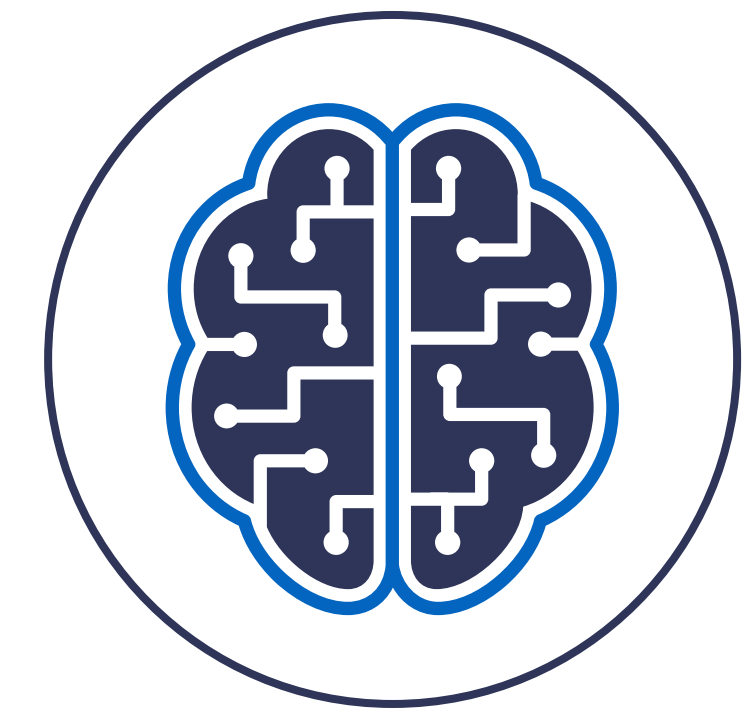| Consumer Goods | Telecoms | Banking/FI | Industrial Goods | Energy | Enterprise |
|---|---|---|---|---|---|
| Demand forecast | Next best offer | Credit risk assessment | Manufacturing process optimization | Production optimization | Back office automation RPA |
| Marketing personalization | Churn and retention modeling | Fraud detection | Predictive maintenance | Predictive maintenance | Performance management |
| Pricing and promo effectiveness | Network optimization | Claim management | Demand and supply forecast | Logistics optimization | Workforce planning |
| Assortment optimization | Infrastructure capacity and utilization | Churn and retention modeling | Operations planning | Project risk management | Scenario simulations |
| Cross sell and upsell | | Next best offer | Energy efficiency | Robotics and automation | |

# THREE TYPES OF MACHINE LEARNING

## Unsupervised Learning

Aim to discover structure:
no target variable known

## Supervised Learning

Aim to predict or model a
known target

## Reinforcement Learning

Optimise actions in a way that
maximises cumulative reward

# Types of algorithms

# Algorithm examples

## Supervised Learning

Algorithms predict class of a new data point from a training set of previously correctly identified observations

## Unsupervised Learning

Algorithms predict results without prior knowledge of the response

**Classification**

**Regression**

**Ranking algorithms**

**Clustering**

**Dimensionality reduction**

**Density estimation**

**Anomaly detection**

Given examples of classes, the model assigns new input data to classes
- **Decision trees,** k-nearest neighbors (kNN), Logistic regression
- **Random Forest,** Support Vector Machines (SVM),Gradient Boosted Decision trees (GBT)
- Neural networks + **Deep Learning**

Given several classes the model assigns input data to classes
- **Linear regression,** Elastic nets
- **Regression trees**

Given ordered pairs examples the model ranks new data

Divide the input data into groups with similar data points assigned to the same group
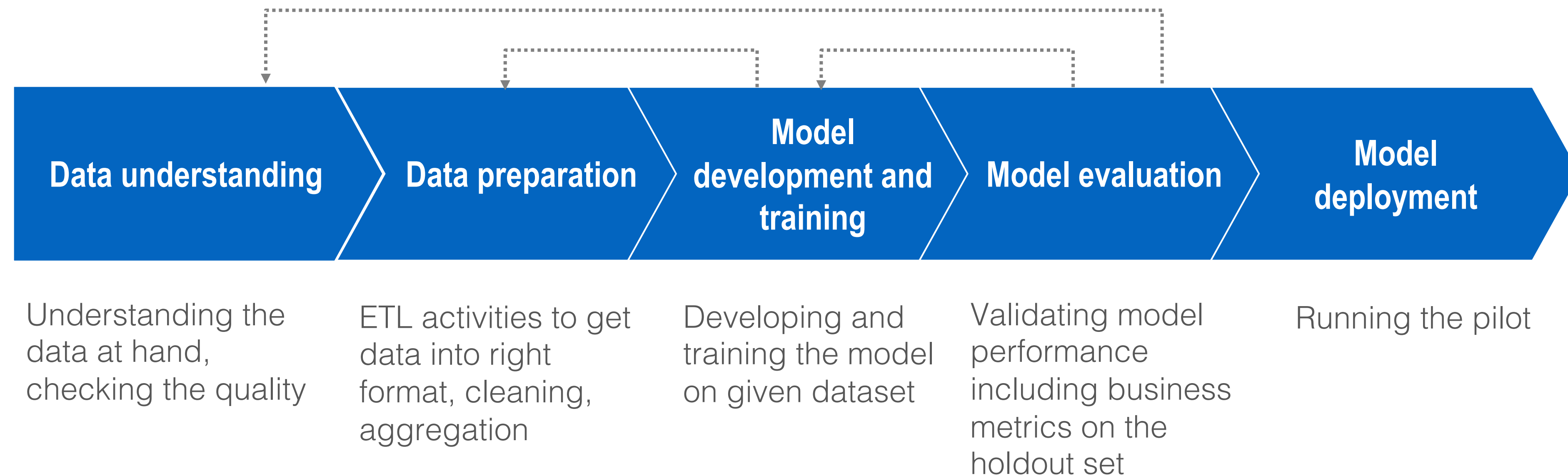- **k-means, spectral**

Mapping input data in lower dimensional space
**PCA**

Estimates the probability distribution of values

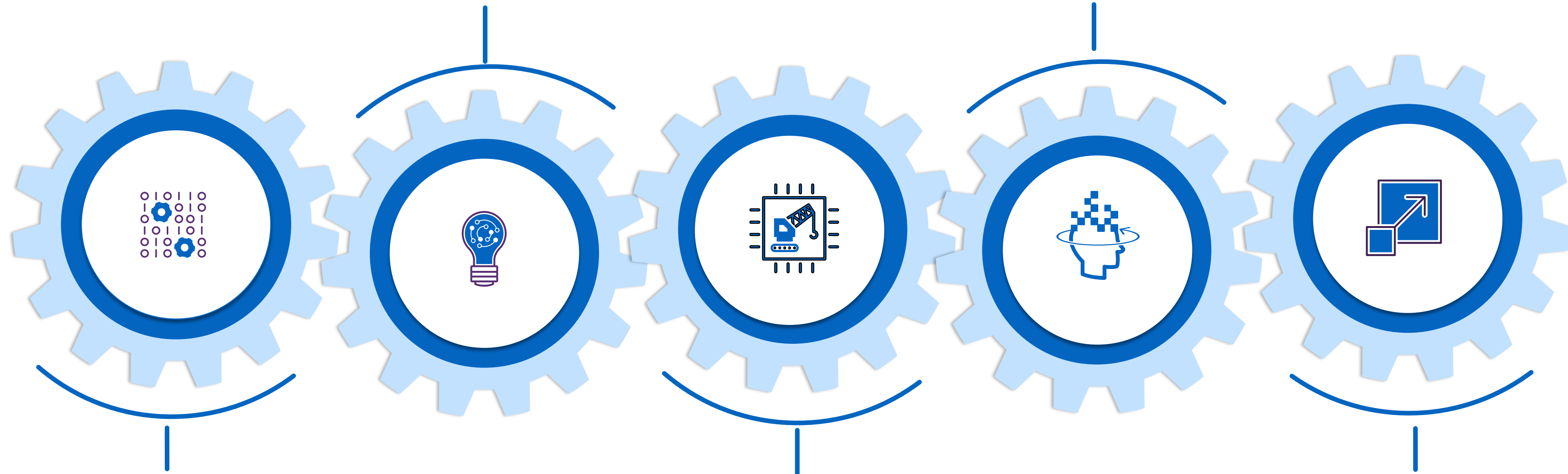Detecting outliers in data

# MODELING PROCESS PIPELINE

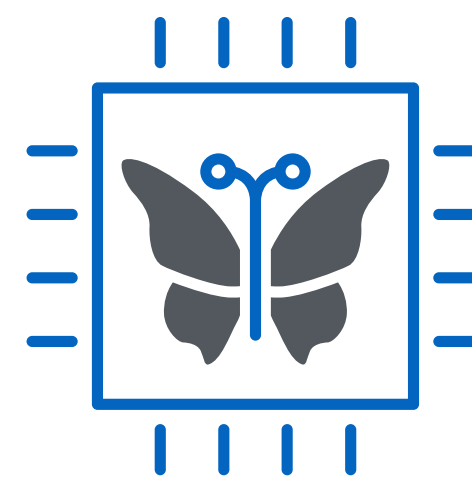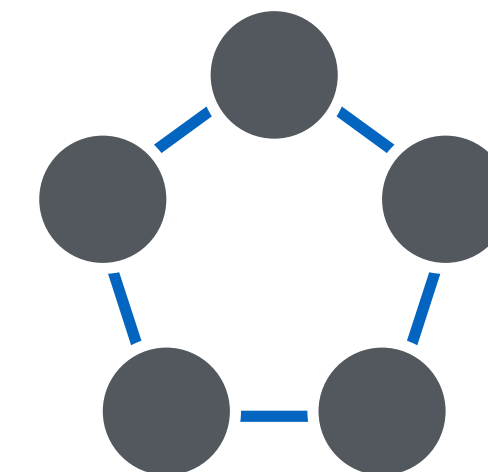| Data understanding | Data preparation | Model development and training | Model evaluation | Model deployment |
|---|---|---|---|---|
| Understanding the data at hand, checking the quality | ETL activities to get data into right format, cleaning, aggregation | Developing and training the model on given dataset | Validating model performance including business metrics on the holdout set | Running the pilot |

rce: BCG

# ANALYTICS PROJECTS SUCCESS FACTORS

Algorithms & Data
- Data analysis
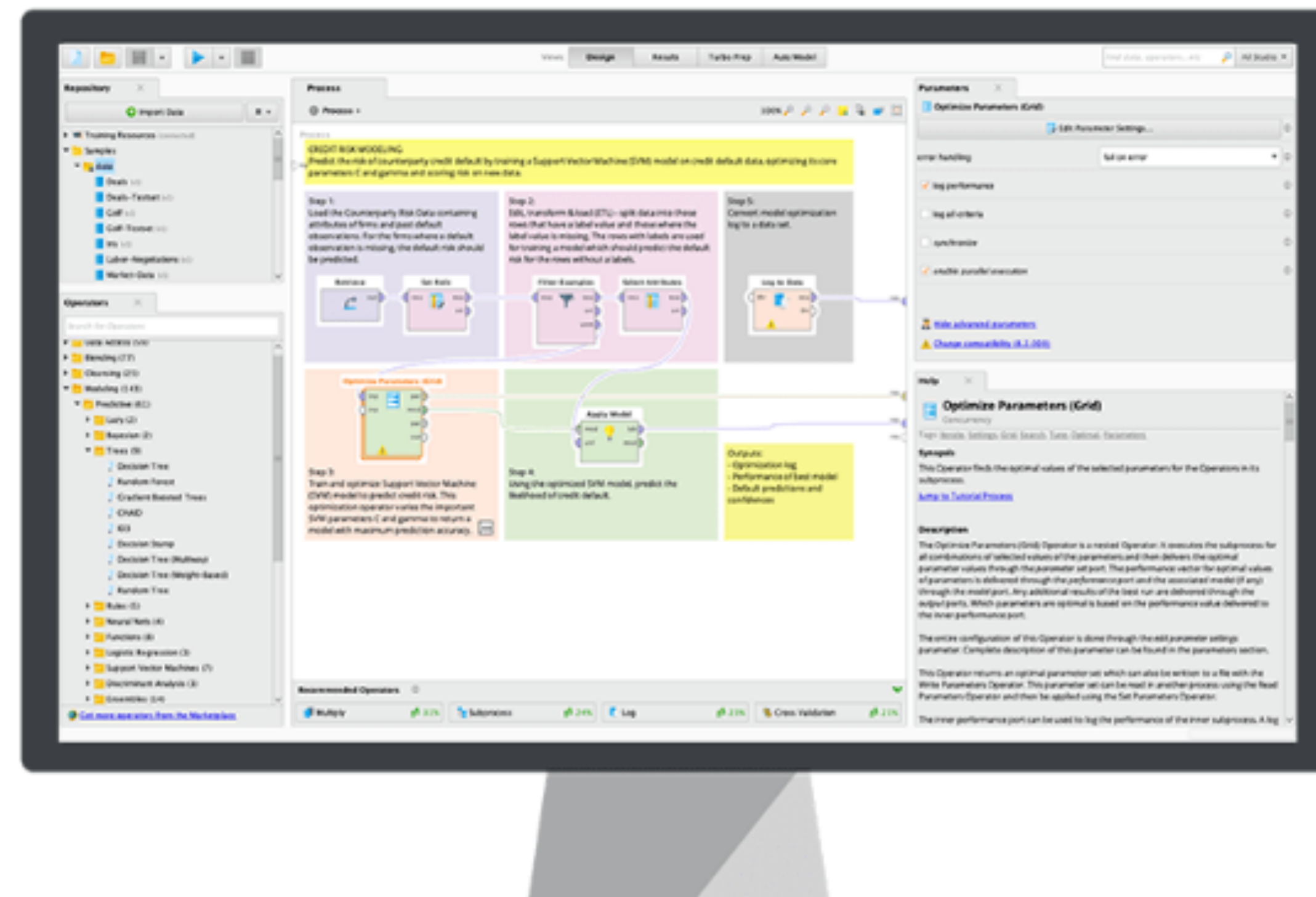- Algorithm development

Technology/IT
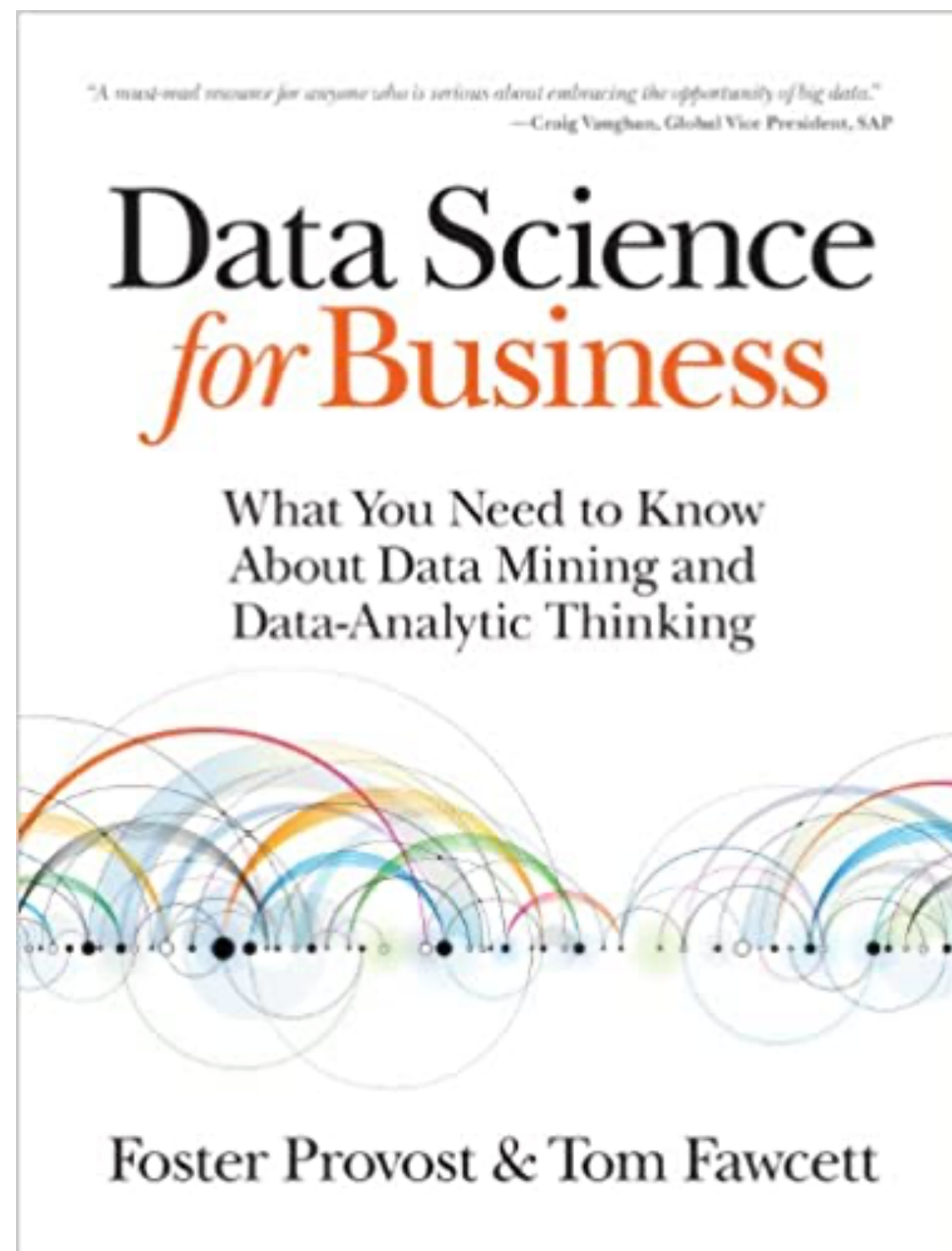- Algorithm industrialization
- Digital platforms development

Business transformation
- Business process redesign
- Enablement
- Change management

# NEXT STEPS