

Модели и технологии оперативного анализа данных

Лекция 6 Введение в Data Science

Гедранович Ольга Брониславовна,
старший преподаватель кафедры ИТ, МИУ
volha.b.k@gmail.com



08.02.2017 ● ○ ○ ○ ○ ● ○ ○ ○ ○ ● ○ ○ ○ ○ ● ○ ○ ○ ○ ● ○ ○ ○ ○ ● ○ ○ ○ ○ 2

Вопросы лекции

- Понятие Data Science
- Отличие Data Science от Business Intelligence
- Процесс Data Science
- Технологии и инструментарий Data Science
- Machine Learning



08.02.2017 ● ○ ○ ○ ○ ● ○ ○ ○ ○ ● ○ ○ ○ ○ ● ○ ○ ○ ○ ● ○ ○ ○ ○ 3

Понятие Data Science

- Interdisciplinary field about processes and systems to extract knowledge or insights from data in various forms, either structured or unstructured, which is a continuation of some of the data analysis fields such as statistics, machine learning, data mining, and predictive analytics, similar to Knowledge Discovery in Databases (KDD).
- Раздел информатики, изучающий проблемы анализа, обработки и представления данных в цифровой форме. Объединяет методы по обработке данных в условиях больших объёмов и высокого уровня параллелизма, статистические методы, методы интеллектуального анализа данных и приложения искусственного интеллекта для работы с данными, а также методы проектирования и разработки баз данных.

Нет однозначности

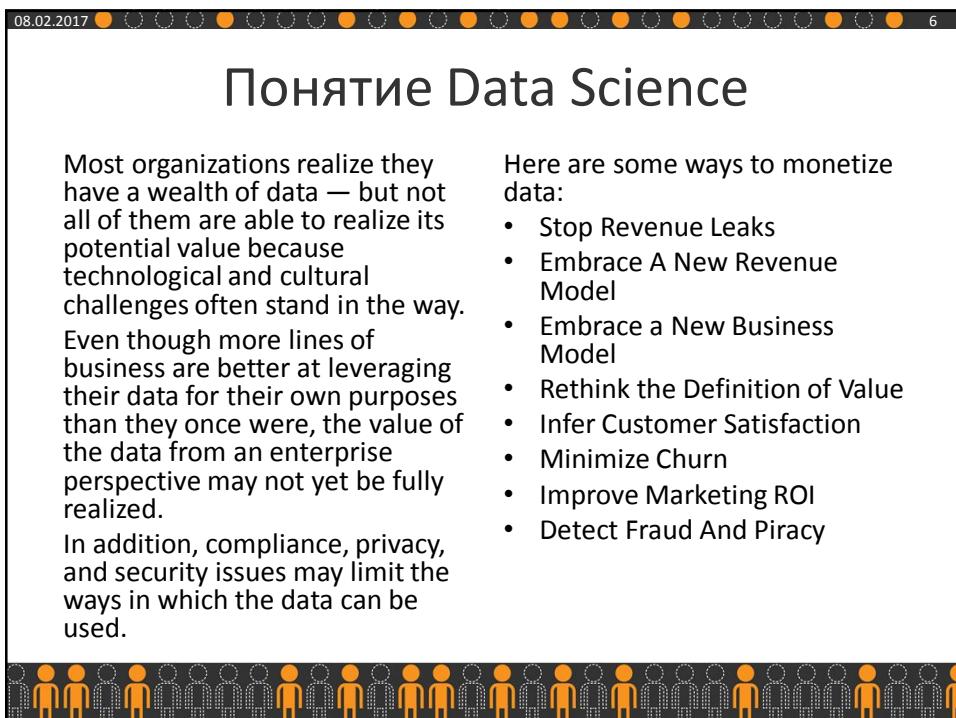
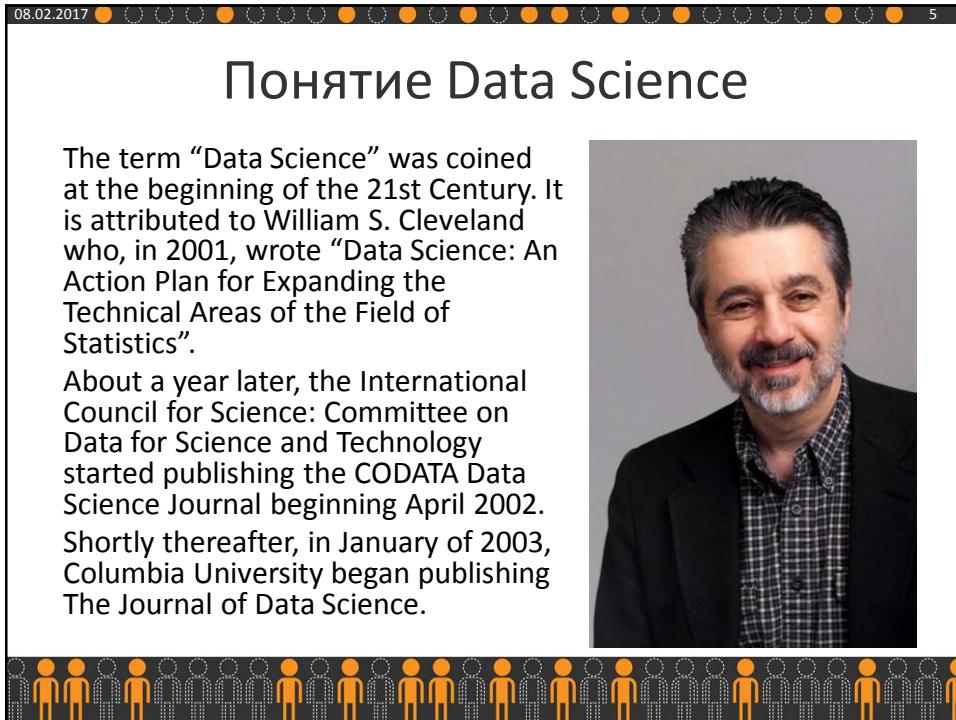


08.02.2017 ● ○ ○ ○ ○ ● ○ ○ ○ ○ ● ○ ○ ○ ○ ● ○ ○ ○ ○ ● ○ ○ ○ ○ ● ○ ○ ○ ○ 4

Понятие Data Science

Data science — as a profession and as an academic discipline unto itself — is new, having been born in the first decade of the 21st century. It is a child born of the mature parental disciplines of scientific methods, data and software engineering, statistics, and visualization.







Понятие Data Science

Какие примеры решений/продуктов, созданных с помощью DS, вы знаете?



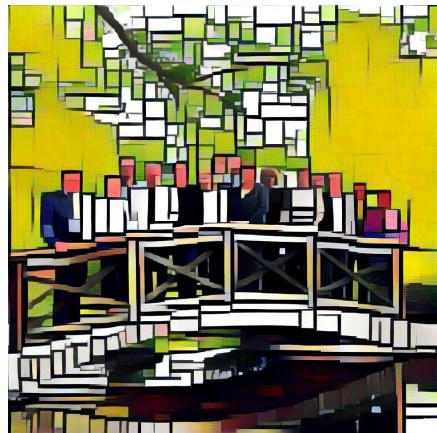
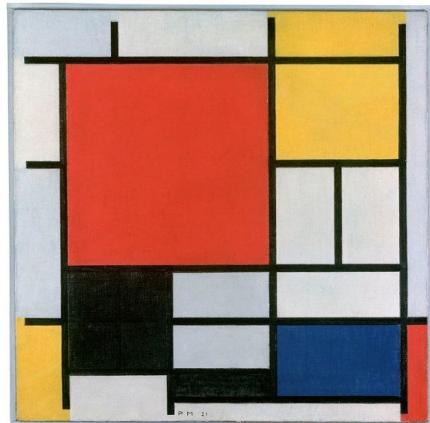
Понятие Data Science

 PRISMA

Prisma transforms your photos and videos into works of art using the styles of famous artists: Van Gogh, Picasso, Levitan, as well as world famous ornaments and patterns. A unique combination of neural networks and artificial intelligence helps you turn memorable moments into timeless art pieces.

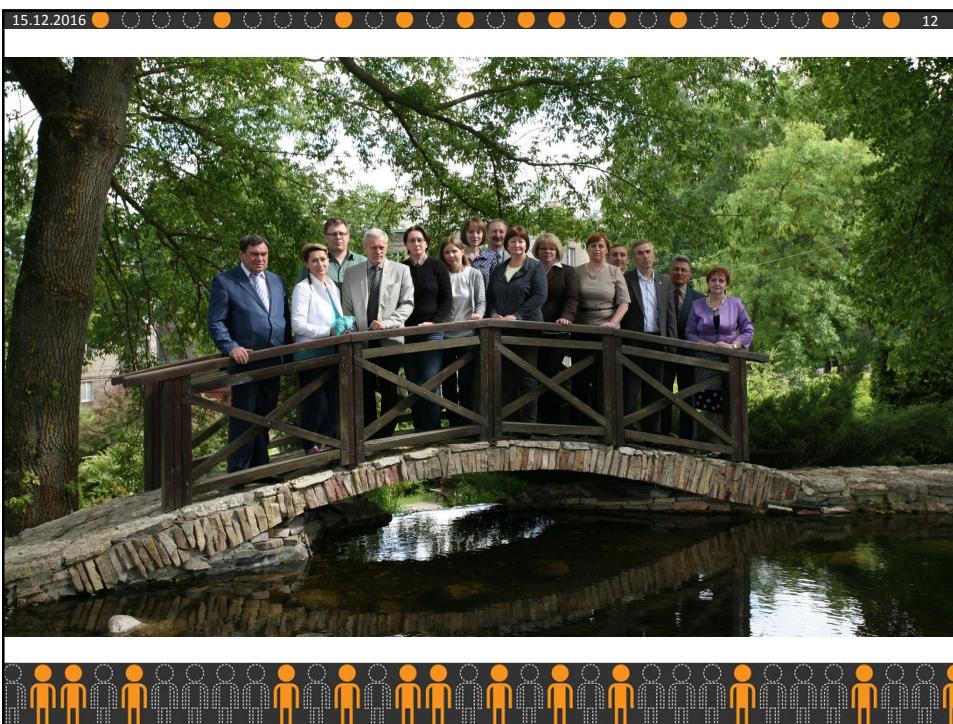


Piet Mondrian “Composition With Red Yellow And Blue”



Edgar Degas “Dancers in blue”





15.12.2016 ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● 13

Понятие Data Science



x.ai – персональный
ассистент, который
планирует встречи



15.12.2016 ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● 14

x.ai

Scenario

You and Mary reconnected during Adweek. She follows up a week later to chat more about a partnership opportunity.

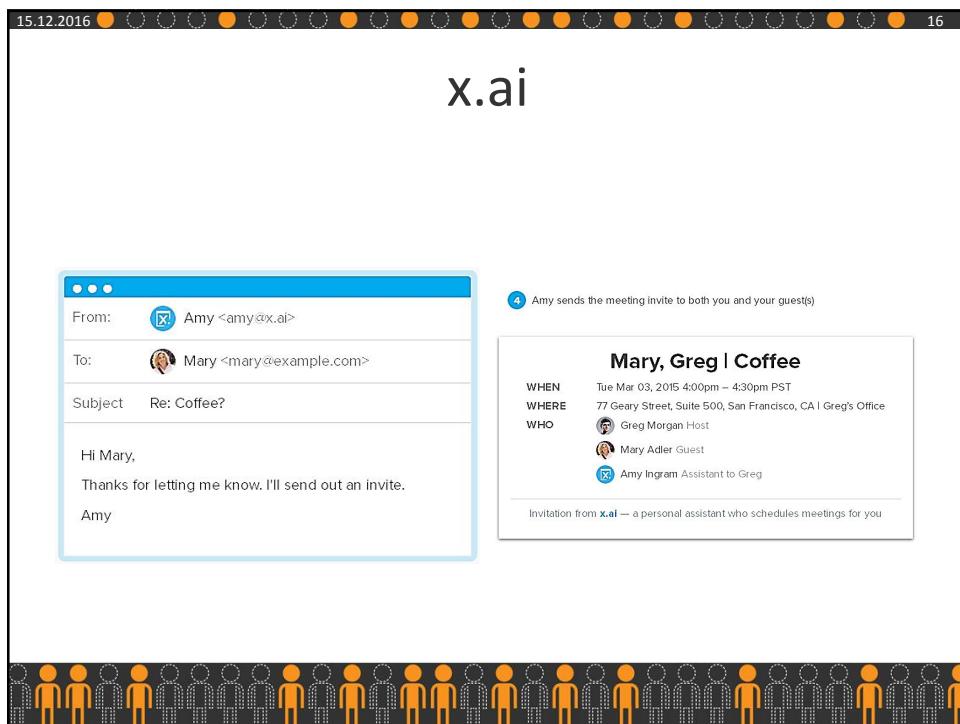
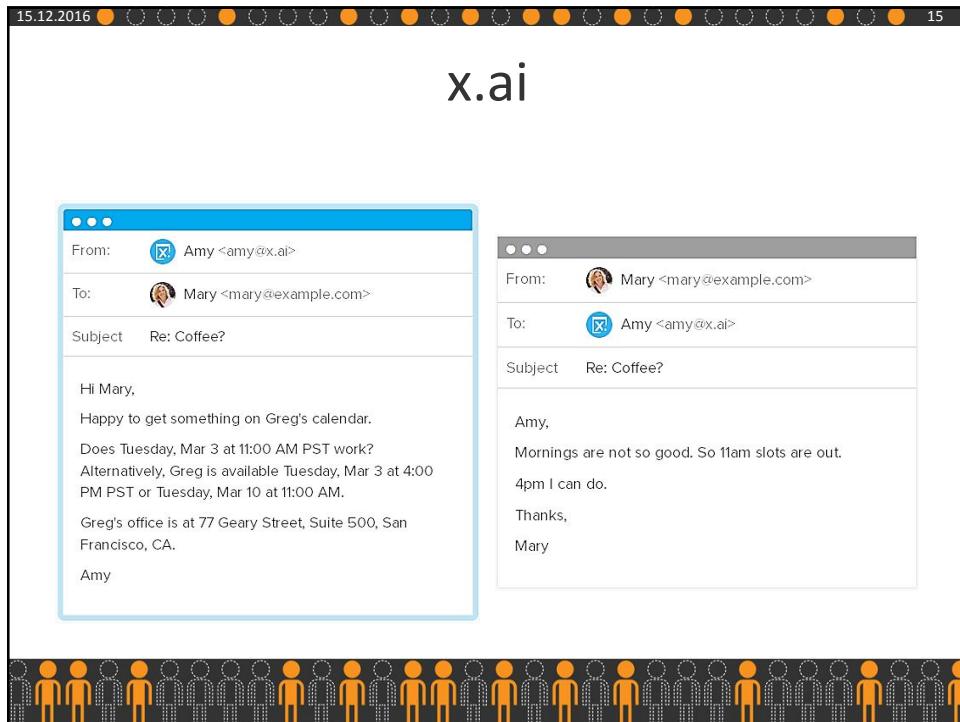
1 A meeting request is made

From:	Mary <mary@example.com>
To:	Greg <greg@example.com>
Subject:	Coffee?
 Hi Greg, It was nice to meet you during Adweek. Do you have a little time for coffee to continue our conversation? I can swing by somewhere close to your office. Thanks, Mary	

2 You reply and cc: Amy

From:	Greg <greg@example.com>
To:	Mary <mary@example.com>
CC:	Amy <amy@x.ai>
Subject:	Re: Coffee?
 Sure, Mary. Amy, can you find 30 minutes for coffee at my office? Cheers, Greg	





15.12.2016 ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● 17

Понятие Data Science



iCarbonX
碳云智能

Optimize your life data
Know more about yourselves

Complete your digital health records by storing genomic data, medical data and daily health data, so as to help you grasp your health condition whenever and wherever you want.



A horizontal row of orange human icons representing a community or population.

08.02.2017 ● ● ● ● ● ● ● ● ● ● ● ● ● ● 18

Основные области применения

BANKING AND SECURITIES	GOVERNMENT
1 Challenges: - Early warning for Securities - Card fraud detection and audit - Enterprise credit risk reporting - Customer behaviour transformation and analytics	2 Challenges: - Integration and use of data The Food and Drug Administration (FDA) is using big data to detect and study patterns of food-related illnesses and diseases, allowing for faster identification and recall of contaminated products.
3 Challenges: - Integrating and analyzing multiple data sources from different providers - Delivering timely and efficient services to customers - Identifying patterns of real-time, media content usage	4 Challenges: - Using Google Maps to deliver better location-based services - Utilizing data from mobile devices to predict where the next matches will be held and what time
5 Challenges: - Increasing the variety of data used in the mining of Natural resources - Integrating data from the manufacturing industry	6 Challenges: - Integration and use of data The Food and Drug Administration (FDA) is using big data to detect and study patterns of food-related illnesses and diseases, allowing for faster identification and recall of contaminated products.
7 Challenges: - Improving pricing, targeted services to customers - Underutilization of data - Hunger for better insight	8 Challenges: - Unbiased Data derived from customer loyalty cards, POS scanners, RFID etc. - Optimized staffing through data from shopping patterns, traffic analysis, weather forecast
9 Challenges: - Increasing the variety of data used in the mining of Natural resources - Integrating data from the manufacturing industry	10 Challenges: - 80% of electric grid assets will become obsolete by 2030 - Wind energy capacity increased by 10% annually - Increased energy management and efficiency - Smart meter readers allow a data collection system to collect data from individual users of the grid

Some applications of big data by governments, private organizations and individuals include:
 - Traffic control and big data
 - Intelligent transport systems by predicting traffic conditions
 - Predictive maintenance of railroads, telecommunications, transportation, energy generation and delivery, and more. For these areas, big data can be used to predict maintenance needs in提前 maintenance
 - Social media structures

A horizontal row of orange human icons representing a community or population.

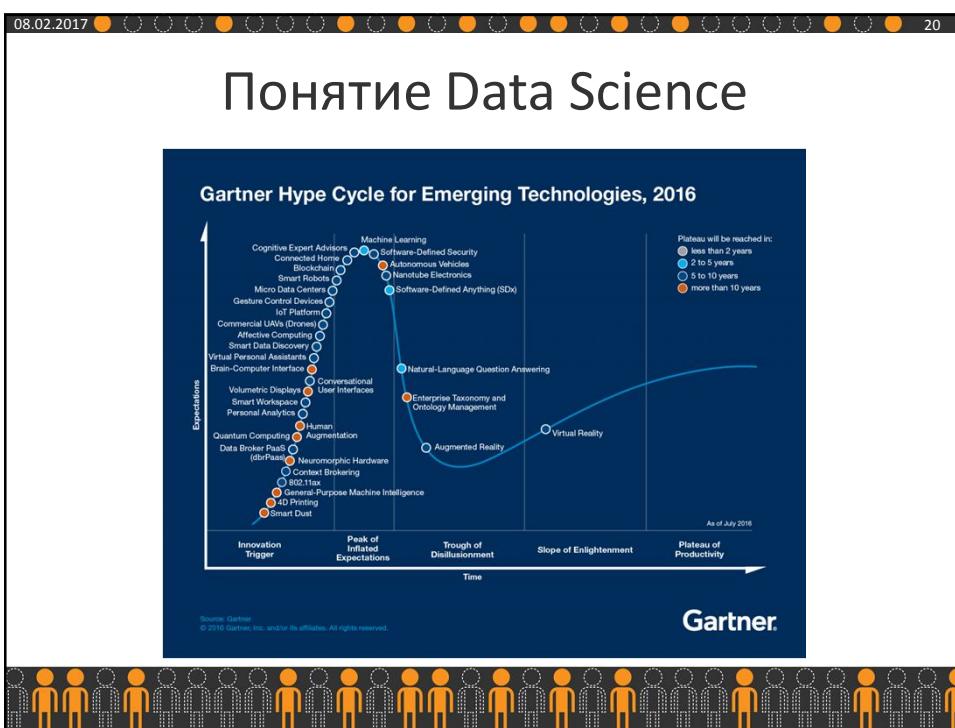
08.02.2017 ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● 19

Понятие Data Science

Data Science Is Multidisciplinary
By Brendan Tierney, 2012

The diagram shows the interconnectedness of various disciplines that contribute to Data Science. It includes Business Strategy, Domain Knowledge, Statistics, Pattern Recognition, Neurocomputing, Communications, AI, Machine Learning, Data Mining, KDD, Database & Data Processing, Problem Solving, Visualisations, Business Analyse, and Inquisitiveness.

08.02.2017 ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● 20



15.12.2016 ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● 21

Кто такой Data Scientist?

MODERN DATA SCIENTIST

Data Scientist, the sexiest job of the 21st century, requires a mixture of multidisciplinary skills ranging from an intersection of mathematics, statistics, computer science, communication and business. Finding a data scientist is hard. Finding people who understand what a data scientist is, is even harder. So here is a little cheat sheet on who the modern data scientist really is.

MATH & STATISTICS	PROGRAMMING & DATABASE	DOMAIN KNOWLEDGE & SOFT SKILLS	COMMUNICATION & VISUALIZATION
<ul style="list-style-type: none"> Machine learning Statistical modeling Experiment design Bayesian inference Supervised learning: decision trees, random forests, logistic regression Unsupervised learning: clustering, dimensionality reduction Optimization: gradient descent and variants 	<ul style="list-style-type: none"> Computer science fundamentals Scripting language e.g. Python Statistical computing packages, e.g. R Databases: SQL and NoSQL Relational algebra Parallel databases and parallel query processing MapReduce concepts Hadoop and HDFS Custom reducers Experience with tools like AWS 	<ul style="list-style-type: none"> Passionate about the business Curious about data Inference without authority Master mindset Problem solver Strategic, practical, creative, innovative and collaborative 	<ul style="list-style-type: none"> Able to engage with senior management Story telling skills Translate data-driven insights into decisions and actions Visual art design It packages like spark or lattice Knowledge of any visualization tools e.g. R, Matplotlib, D3.js, Tableau

Marketing Data Scientists are a group of practitioners in the area of a company marketing. (Our field of expertise includes marketing strategy and optimization, customer tracking and analysis, predictive analytics and econometrics, data mining and big data systems, marketing channel insights in Paid Search, SEO, Social, CRM and beyond.)

Marketing Data Scientists are a group of practitioners in the area of a company marketing. (Our field of expertise includes marketing strategy and optimization, customer tracking and analysis, predictive analytics and econometrics, data mining and big data systems, marketing channel insights in Paid Search, SEO, Social, CRM and beyond.)

Marketing Data Scientists are a group of practitioners in the area of a company marketing. (Our field of expertise includes marketing strategy and optimization, customer tracking and analysis, predictive analytics and econometrics, data mining and big data systems, marketing channel insights in Paid Search, SEO, Social, CRM and beyond.)

15.12.2016 ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● 22

Кто такой Data Scientist?

The diagram illustrates the Data Scientist role and its relationship to various disciplines and skill sets:

- Core Competencies:** Scientific method, simulation, programming, privacy & security, data & text mining, statistics, artificial intelligence, machine learning, natural language processing, graph analysis, information retrieval, data warehousing, math/stats, databases, business intelligence, big data, data product design, entrepreneurship, domain knowledge, ethics, data visualization, art & design, communication.
- Interdisciplinary Areas:** Computer Science (privacy & security, programming, cloud computing, distributed systems, technology & infrastructure), Analytics (scientific method, simulation, data & text mining, machine learning, natural language processing, graph analysis, information retrieval, data warehousing, math/stats, databases, business intelligence, big data).
- Skills Venn Diagram:**
 - Software Engineer:** C#, Software Development, HTML, XML, jQuery, CSS, .NET, Web Services, PHP, Agile Methodologies.
 - Data Engineer:** MySQL, Hadoop, Data Warehousing, ETL, Hive, Java, Unix, Oracle, Business Intelligence.
 - Data Scientist:** R, Statistics, MATLAB, Machine Learning, Data Mining, Programming, SAS, Research Analytics, Algorithms, Statistical Modeling.

Data Roles and Skill Sets

Marketing Data Scientists are a group of practitioners in the area of a company marketing. (Our field of expertise includes marketing strategy and optimization, customer tracking and analysis, predictive analytics and econometrics, data mining and big data systems, marketing channel insights in Paid Search, SEO, Social, CRM and beyond.)

15.12.2016 ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● 23

В чём отличие Data Science от смежных областей ИТ-индустрии?

A Venn diagram illustrating the interdisciplinary nature of Data Science. It consists of three overlapping circles: a blue circle labeled 'MATHEMATICS' (bottom left), a green circle labeled 'COMPUTER SCIENCE' (bottom right), and an orange circle labeled 'DOMAIN EXPERTISE' (top). The central area where all three circles overlap is labeled 'DATA SCIENCE'. Within this central area, the words 'MACHINE LEARNING' are written. Below the diagram, a small caption reads: 'Source: Palmer, Shelly. *Data Science for the C-Suite*. New York: Digital Living Press, 2015. Print.'

DATA SCIENCE

MATHEMATICS

COMPUTER SCIENCE

DOMAIN EXPERTISE

STATISTICAL RESEARCH

DATA PROCESSING

MACHINE LEARNING

Source: Palmer, Shelly. *Data Science for the C-Suite*.
New York: Digital Living Press, 2015. Print.

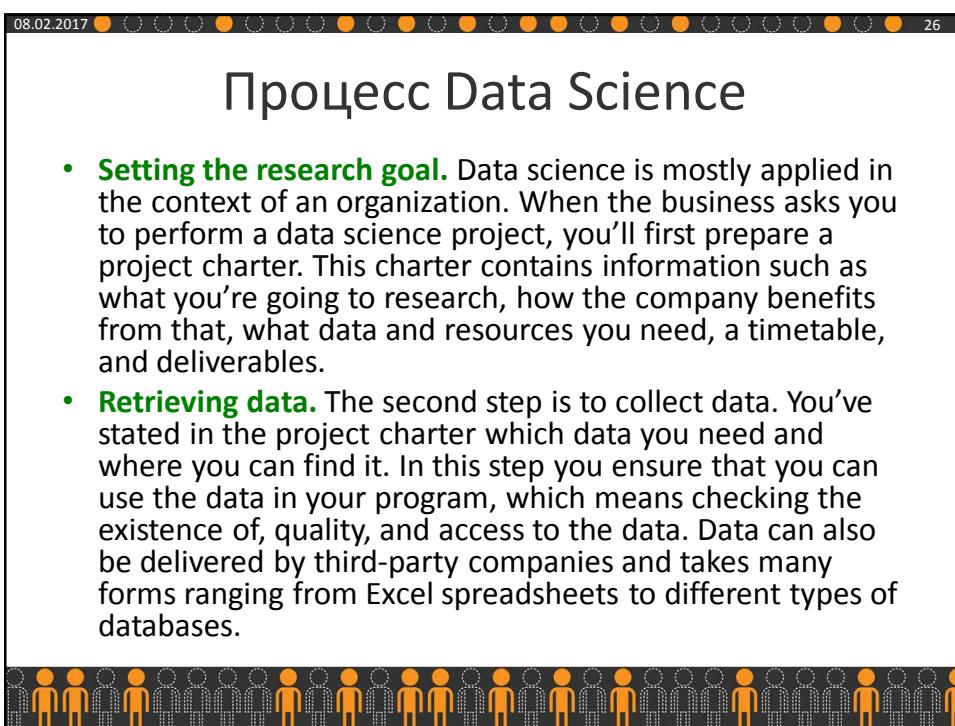
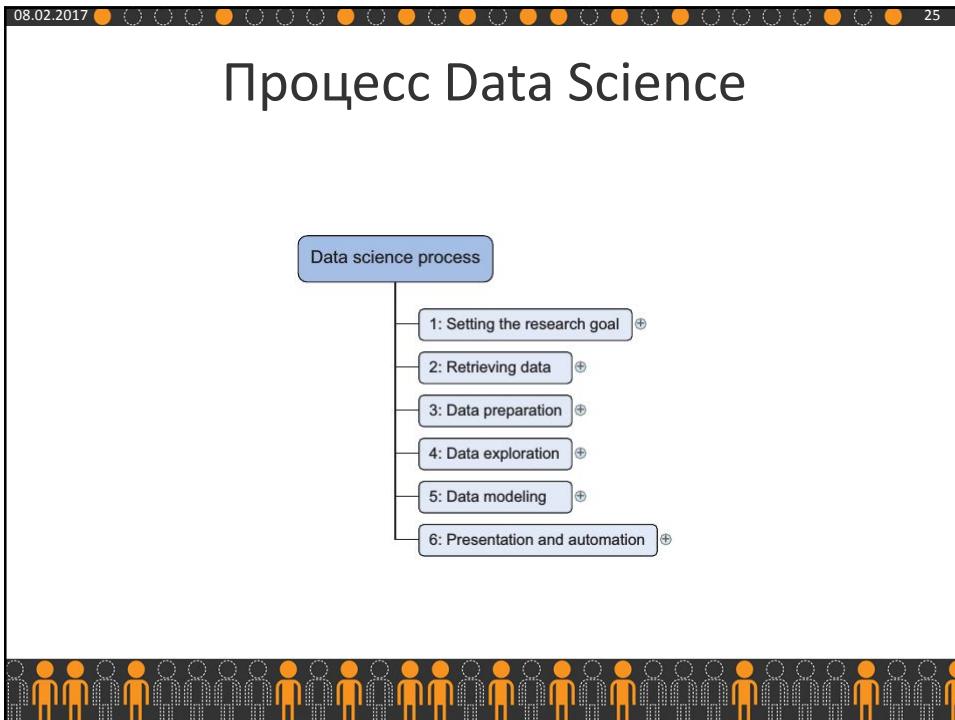
15.12.2016 ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● 24

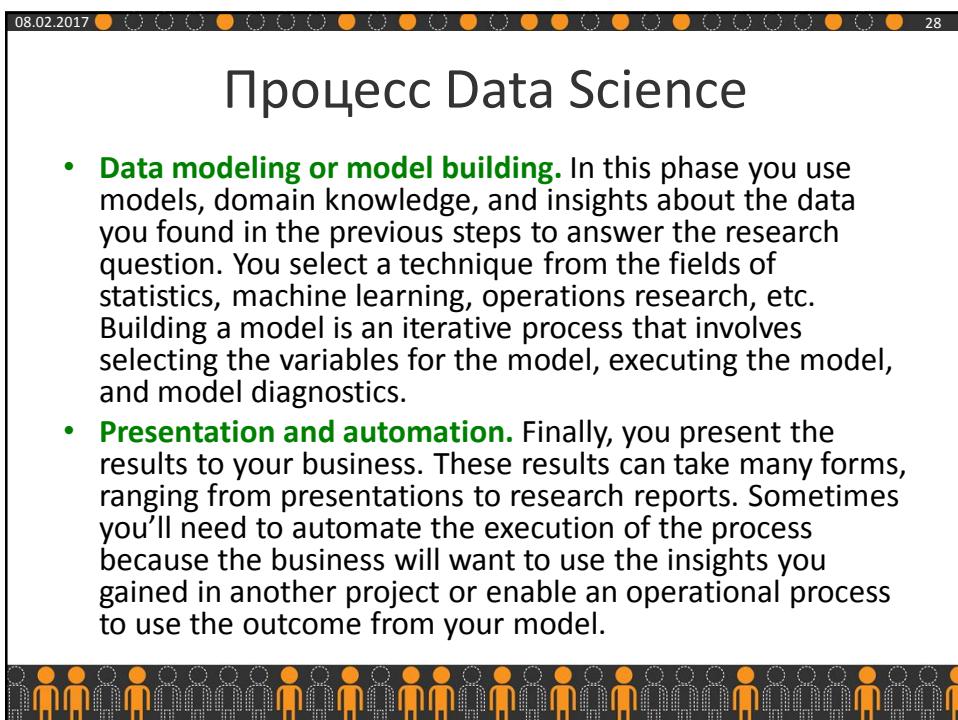
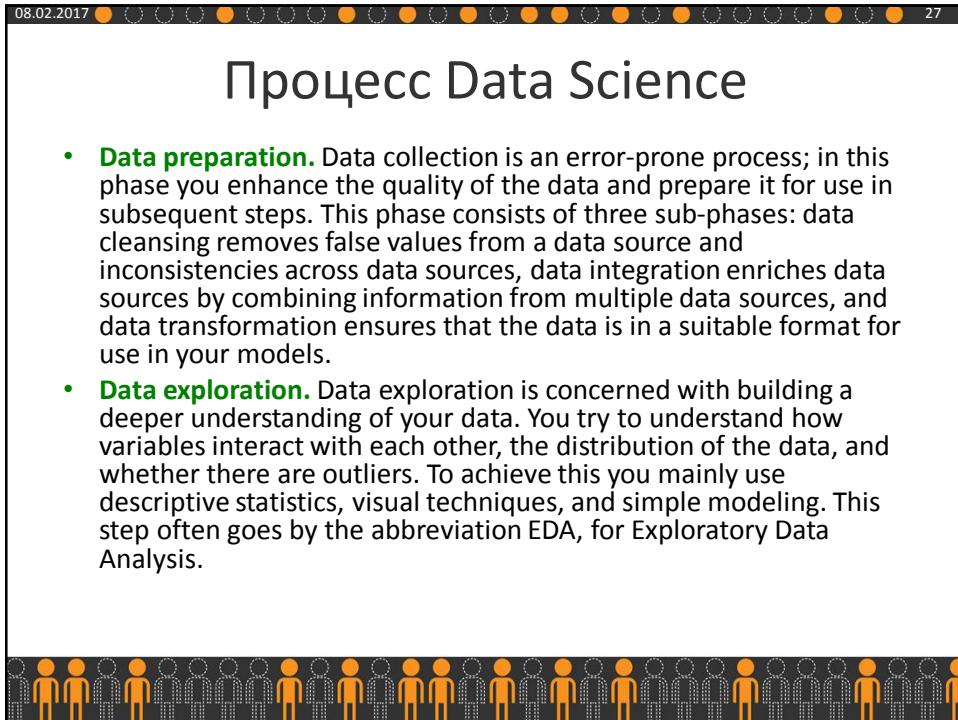
Отличие Data Science от Business Intelligence

Data Science vs. Business Intelligence

	Business Intelligence (BI)	Data Science
Data analysis	Yes	Yes
Statistics	Yes	Yes
Visualization	Yes	Yes
Data Sources	Usually SQL, often Data Warehouse	Less structured (logs, cloud data, SQL, noSQL, text)
Tools	Statistics, Visualization	Statistics, Machine Learning, Graph Analysis, NLP
Focus	Present and past	Future
Approach	Analytic	Scientific
Goal	Better strategic decisions	Advanced functionality

The two fields are closely related. In some ways Data Science is an evolution of BI.





08.02.2017 ● 29

Процесс Data Science

The previous description of the data science process gives you the impression that you walk through this process in a linear way, but in reality you often have to step back and rework certain findings.

For instance, you might find outliers in the data exploration phase that point to data import errors. As part of the data science process you gain incremental insights, which may lead to new questions.



08.02.2017 ● 30

Процесс Data Science

Data Science Process

```

graph TD
    Reality((Reality)) --> RawData[Raw Data Collected]
    RawData --> DataIsProcessed[Data Is Processed]
    DataIsProcessed --> CleanDataset[Clean Dataset]
    CleanDataset --> ModelsAlgorithms[Models & Algorithms]
    CleanDataset --> ExploratoryAnalysis[Exploratory Data Analysis]
    ExploratoryAnalysis --> ModelsAlgorithms
    ModelsAlgorithms --> MakeDecisions[Make Decisions]
    MakeDecisions --> CommunicateReport[Communicate Visualize Report]
    CommunicateReport --> DataProduct[Data Product]
    DataProduct --> Reality
    
```



08.02.2017 ● 31

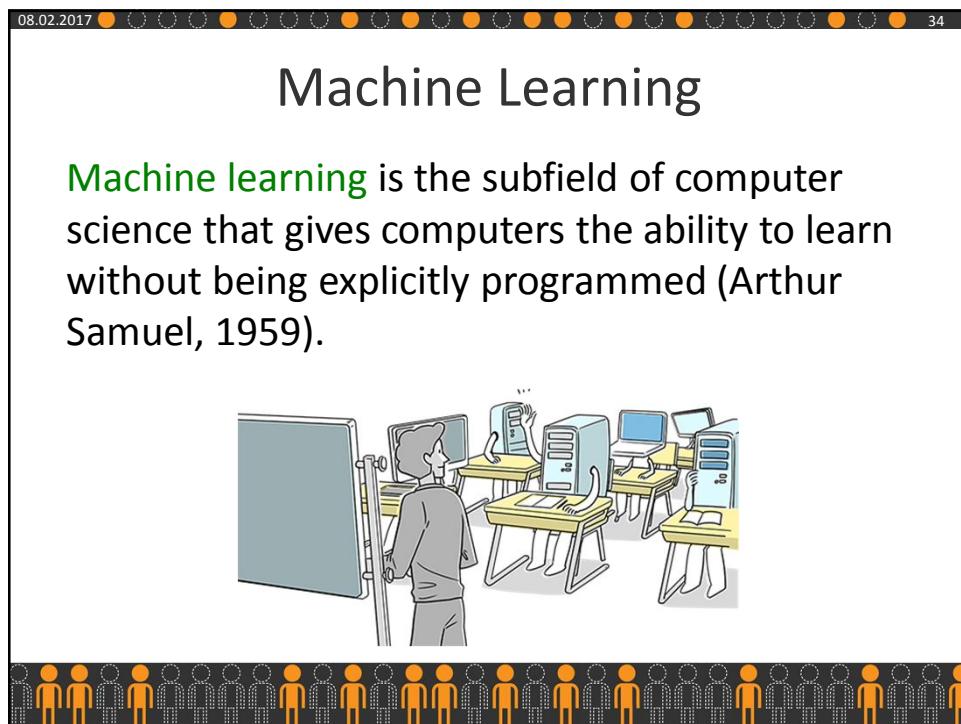
Технологии и инструментарий

<ul style="list-style-type: none"> • Linear Regression • Logistic Regression • Jackknife Regression • Density Estimation • Confidence Interval • Test of Hypotheses • Pattern Recognition • Clustering - (aka Unsupervised Learning) • Supervised Learning • Time Series • Decision Trees • Random Numbers • Monte-Carlo Simulation • Bayesian Statistics • Naive Bayes 	<ul style="list-style-type: none"> • Principal Component Analysis - (PCA) • Ensembles • Neural Networks • Support Vector Machine - (SVM) • Nearest Neighbors - (k-NN) • Feature Selection - (aka Variable Reduction) • Indexation / Cataloguing • (Geo-) Spatial Modeling • Recommendation Engine • Search Engine • Attribution Modeling • Collaborative Filtering • Rule System • Linkage Analysis 	<ul style="list-style-type: none"> • Association Rules • Scoring Engine • Segmentation • Predictive Modeling • Graphs • Deep Learning • Game Theory • Imputation • Survival Analysis • Arbitrage • Lift Modeling • Yield Optimization • Cross-Validation • Model Fitting • Relevancy Algorithm • Experimental Design
--	---	--

08.02.2017 ● 32

Технологии и инструментарий

Source Data	
Store Data	
Convert & ETL	
Transform Data	
Exploratory Analysis	
Model Build & Generate Insights	
Visualisation	
Model Execution in Production	



08.02.2017 ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● 35

Machine Learning

Basic tasks:

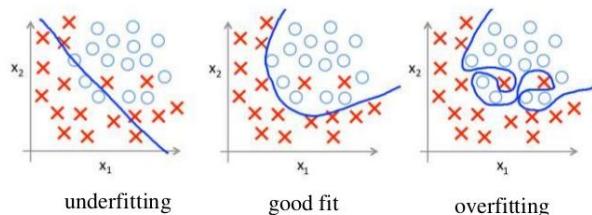
- In **classification**, inputs are divided into two or more classes, and the learner must produce a model that assigns unseen inputs to one or more (multi-label classification) of these classes. This is typically tackled in a supervised way. Spam filtering is an example of classification, where the inputs are email (or other) messages and the classes are "spam" and "not spam".
- In **regression**, also a supervised problem, the outputs are continuous rather than discrete.
- In **clustering**, a set of inputs is to be divided into groups. Unlike in classification, the groups are not known beforehand, making this typically an unsupervised task.
- **Collaborative filtering** is a technique used by recommender systems. It is a method of making automatic predictions (filtering) about the interests of a user by collecting preferences or taste information from many users (collaborating).



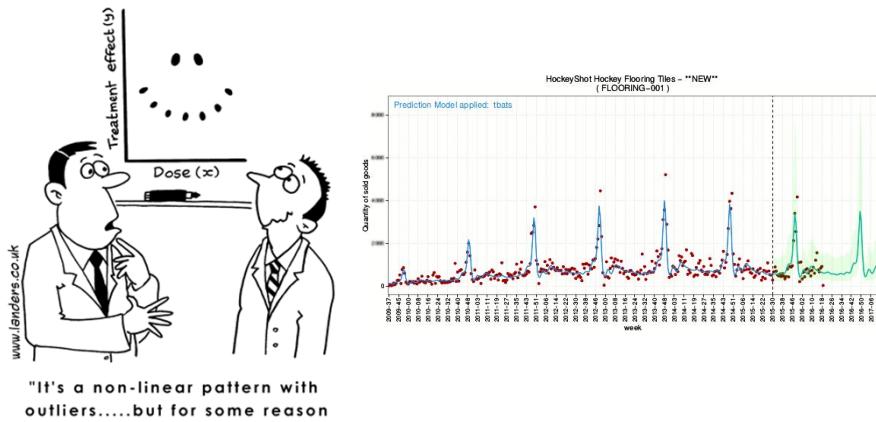
08.02.2017 ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● 36

Classification

Overfitting and underfitting



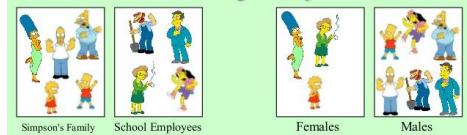
Regression



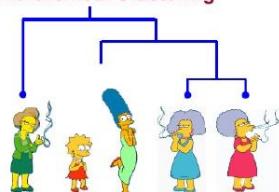
Clustering



Clustering is subjective



Hierarchical Clustering



Partition-based Clustering

