

# Краткое введение в data mining

Николай Анохин

# Data Mining как KDD

*Knowledge Discovery in Databases (KDD) – это процесс получения точных, неизвестных, потенциально полезных и интерпретируемых закономерностей из данных.<sup>1</sup>*

---

<sup>1</sup>U. Fayyad, G. Piatetsky-Shapiro, P. Smyth. From data mining to knowledge discovery: an overview. 1996

# Data Mining как моделирование

*Data Mining – процесс построения модели, хорошо описывающей закономерности, которые порождают данные.*

Подходы к построению моделей

- ▶ статистический
- ▶ машинное обучение
- ▶ вычислительный

## Пример 1. Красная икра на новогодний стол

настоящая	446	521	550	315	613	292	469	658	255	310
искусственная	372	351	361	398	348	457	370	473	475	435

### Статистический подход

$$\begin{cases} p(\text{цена}|\text{настоящая}) \sim \mathcal{N}(\text{цена}|\mu_r, \sigma_r) \\ p(\text{цена}|\text{искусственная}) \sim \mathcal{N}(\text{цена}|\mu_a, \sigma_a) \end{cases} \xrightarrow{MLE} \begin{cases} \mu_r = 443, \sigma_r = 136 \\ \mu_a = 404, \sigma_a = 49 \end{cases}$$

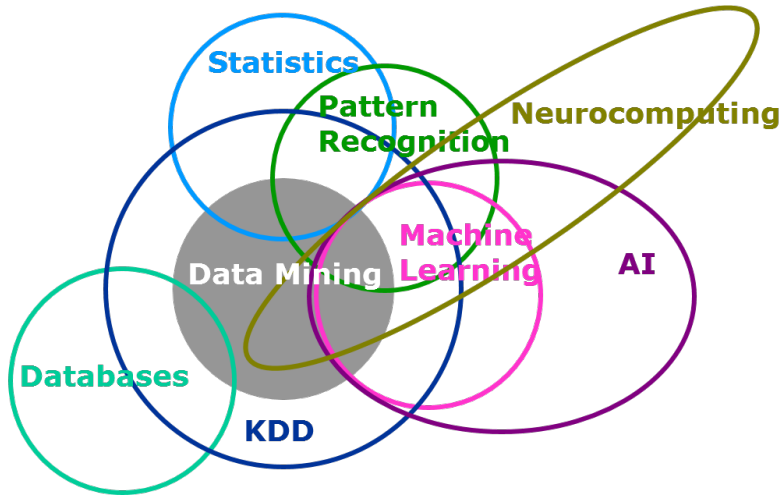
### Машинное обучение

Обучаем линейный SVM:  $\text{цена} > 482 \Rightarrow \text{настоящая}$

### Вычислительный подход

Подсчитываем параметры данных:  $\langle \text{цена}_r \rangle = 443, \langle \text{цена}_a \rangle = 404$

## Data Mining – область на пересечении дисциплин<sup>2</sup>



---

<sup>2</sup>Looking backwards, looking forwards: SAS, data mining, and machine learning

# Data Mining – область тысячи имен

1960-е Data Fishing, Data Dredging

1980-е Knowledge Discovery in Databases

1990-е Data Mining, Database mining<sup>TM</sup>

2000-е Data Analytics, Data Science<sup>34</sup>

---

<sup>3</sup>Data Scientist is a Data Analyst who lives in California

<sup>4</sup>A data scientist is someone who is better at statistics than any software engineer and better at software engineering than any statistician.

## Некоторые важные события в истории Data Mining

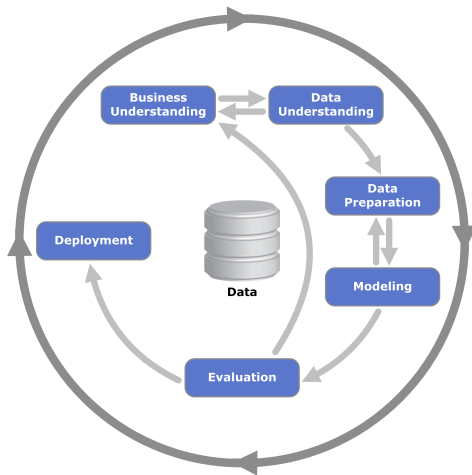
- 1989 IJCAI-89 Workshop on Knowledge Discovery in Databases
- 1995 ACM SIGKDD Conference on Knowledge Discovery and Data Mining
- 2003 Программа Total Information Awareness
- 2001 Leo Breiman's "Statistical Modeling: The Two Cultures"
- 2005 Doug Cutting и Mike Cafarella разработали пакет обработки данных Hadoop
- 2007 Первый релиз библиотеки scikit-learn
- 2010 Заработал сайт Kaggle – платформа для проведения соревнований по Data Science
- 2012 Harvard Business Review публикует статью Data Scientist: The Sexiest Job of the 21st Century
- 2013 Первая встреча сообщества Moscow Data Science<sup>5</sup> в московском офисе Mail.Ru Group

---

<sup>5</sup><http://www.meetup.com/Moscow-Data-Science/>

# CRISP-DM

(Cross Industry Standard Process for Data Mining)





## Пример 2. Игра в гольф

### **Business understanding**

- ▶ понимание задачи с точки зрения бизнеса
- ▶ сбор требований и ограничений
- ▶ постановка задачи в терминах Data Mining

### **Data understanding**

- ▶ первичный сбор данных
- ▶ ознакомление с данными и понимание их специфики

### **Data preparation**

- ▶ формирование финального набора данных

## Игра в гольф: данные

Outlook	Temperature	Humidity	Wind	Play
Sunny	85	85	false	no
Sunny	80	90	true	no
Overcast	83	86	false	yes
Rainy	70	96	false	yes
Rainy	68	80	false	yes
Rainy	65	70	true	no
Overcast	64	65	true	yes
Sunny	72	95	false	no
Sunny	69	70	false	yes
Rainy	75	80	false	yes
Sunny	75	70	true	yes
Overcast	72	90	true	yes
Overcast	81	75	false	yes
Rainy	71	91	true	no

## Признаки

$\mathcal{D}$  – множество, содержащее все рассматриваемые в задаче объекты

$d \in \mathcal{D}$  – объект,  $\phi_j : \mathcal{D} \rightarrow F_j$  – признак

Виды признаков

- ▶ Бинарные/Binary

$$F_j = \{true, false\}$$

- ▶ Номинальные/Categorical

$F_j$  – конечно

- ▶ Порядковые/Ordinal

$F_j$  – конечно, упорядочено

- ▶ Количественные/Numerical

$$F_j = \mathbb{R}$$

Признаковое представление объекта  $d_i$

$$\mathbf{x}_i = (\phi_1(d_i), \dots, \phi_n(d_i)) \in \mathcal{X}$$

# Моделирование

- ▶ Перебор различных вариантов моделирования
- ▶ Настройка параметров моделей

## Модель

семейство параметрических функций вида

$$H = \{h(\mathbf{x}, \theta) : \mathcal{X} \times \Theta \rightarrow \mathcal{Y}\}$$

## Алгоритм обучения

выбор наилучших параметров  $\theta^*$  с использованием обучающей выборки

$$A(X, Y) : (\mathcal{X} \times \mathcal{Y})^N \rightarrow \Theta$$

В итоге:

$$h^*(\mathbf{x}) = h(\mathbf{x}, \theta^*)$$