



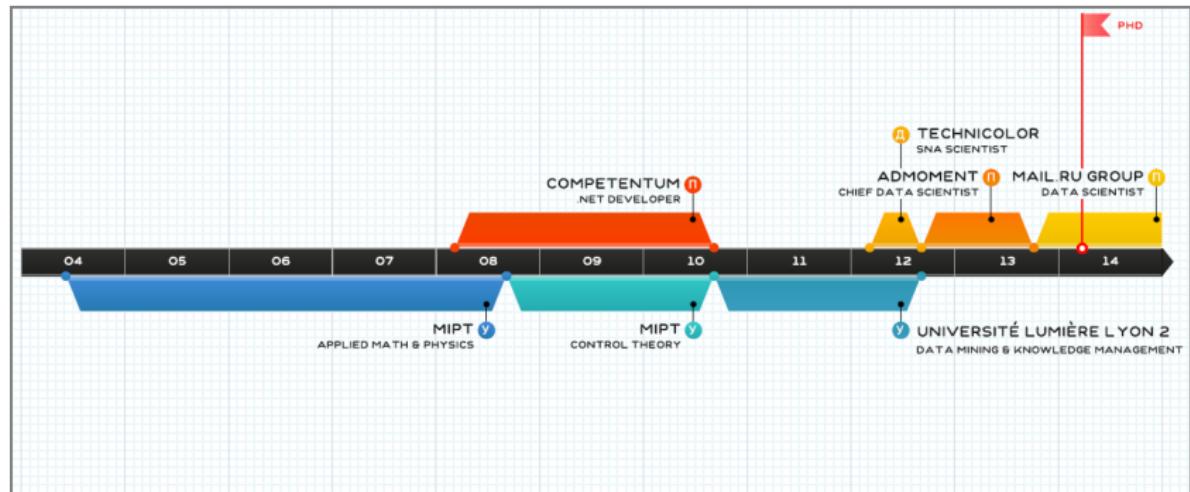
ТЕХНОСФЕРА

Лекция 1 Задачи Data Mining

Николай Анохин

24 сентября 2014 г.

Николай Анохин



e-mail: n.anokhin@corp.mail.ru

тел.: +7 (903) 111-44-60

План лекции

Структура курса

Что такое Data Mining

Процесс Data Mining

Exploratory data analysis

Структура курса

Модуль 1

1. Задачи Data Mining (Николай Анохин)
2. Задача кластеризации и ЕМ-алгоритм (Николай Анохин)
3. Различные алгоритмы кластеризации (Николай Анохин)^H
4. Задача классификации (Николай Анохин)
5. Naive Bayes (Николай Анохин)
6. Линейные модели (Николай Анохин)
7. Метод опорных векторов (Николай Анохин)^{HP}

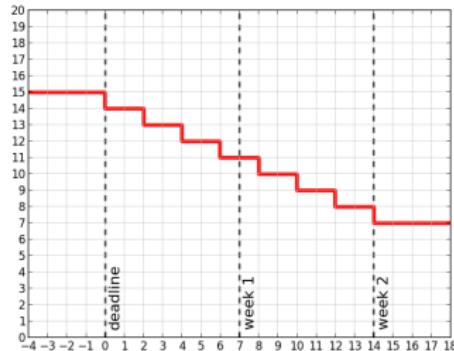
Модуль 2

1. Снижение размерности пространства (Владимир Гулин)
2. Алгоритмические композиции 1 (Владимир Гулин)
3. Алгоритмические композиции 2 (Владимир Гулин)^H
4. Нейросети, обучение с учителем (Павел Нестеров)^H
5. Нейросети, обучение без учителя (Павел Нестеров)
6. Нейросети, глубокие сети (Павел Нестеров)

Контроль знаний

дз

4 домашних задания на 6-8 часов самостоятельной работы каждое
(4x15 баллов)



Экзамен

Презентация и защита семестрового проекта
(40 баллов)

Правила

- + Можно задавать вопросы по ходу лекции
- + Можно входить и выходить, не мешая коллегам
- Нельзя нарушать порядок в аудитории
- Нельзя разговаривать по телефону
- ▶ Общение с преподавателем на “Вы”

Ваши правила?

DM как KDD

Data Mining

Процесс извлечения знаний из различных источников данных, таких как базы данных, текст, картинки, видео и т.д. Полученные знания должны быть достоверными, полезными и интерпретируемыми.

DM как моделирование

Data Mining

Процесс построения модели, хорошо описывающей закономерности, которые порождают данные.

Подходы к построению моделей

- ✓ статистический
- ✓ на основании машинного обучения
- ✗ вычислительный

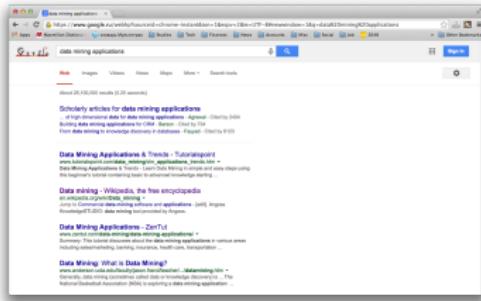
Data Scientist

Person who is better at statistics than any software engineer and better at software engineering than any statistician
(J. Wills, Data Scientist at Cloudera Inc.)

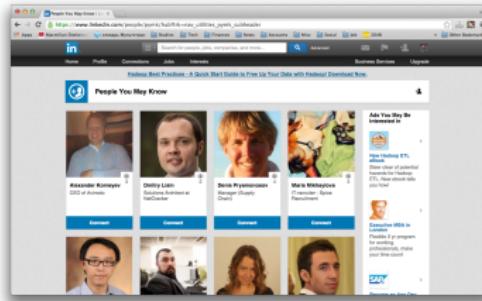
Data-

- ✖ -architecture
- ✖ -acquisition
- ✓ -analysis
- ✖ -archiving

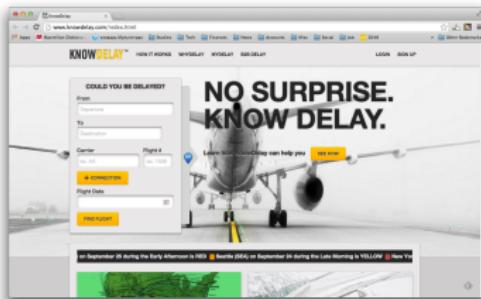
Success stories



(a) Google



(b) LinkedIn



(c) KnowDelay



(d) Handwritten digits

Fail (Программа TIA)

- ▶ Наблюдаем 10^9 человек
- ▶ Человек в среднем посещает отель раз в 100 дней
- ▶ Есть 10^5 отелей на 100 человек каждый
- ▶ Проверим посещения за 1000 дней

Вероятность для конкретной пары встретиться в отеле в конкретный день:

$$p_1 = \left(\frac{1}{100} \right)^2 \cdot 10^{-5} = 10^{-9}$$

Всего пар людей

$$n_{pp} = C_2^{10^9} \approx \frac{(10^9)^2}{2} = 5 \cdot 10^{17}$$

а пар дней

$$n_{pd} = C_2^{10^3} \approx \frac{(10^3)^2}{2} = 5 \cdot 10^5$$

Ожидаемое количество “подозрительных” встреч в отелях

$$N = p_1^2 n_{pp} n_{pd} = 250000 >> 10$$

Принцип Бонферрони

Вычислить количество рассматриваемых событий при предположении их полной случайности. Если это количество намного превосходит количество событий, о котором идет речь в задаче, полученные результаты нельзя будет считать достоверными.

Cross Industry Standard Process for Data Mining

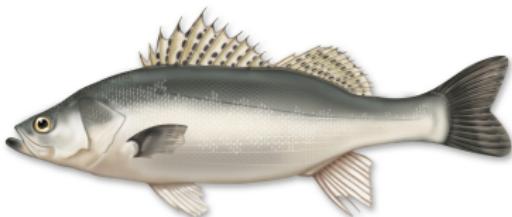
CRISP-DM

- ▶ SPSS
- ▶ Teradata
- ▶ Daimler AG
- ▶ NCR Corporation
- ▶ OHRA
- ▶ IBM

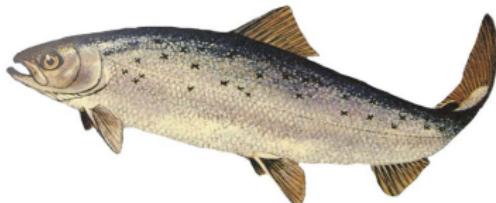
Другие процессы: KDD, SEMMA



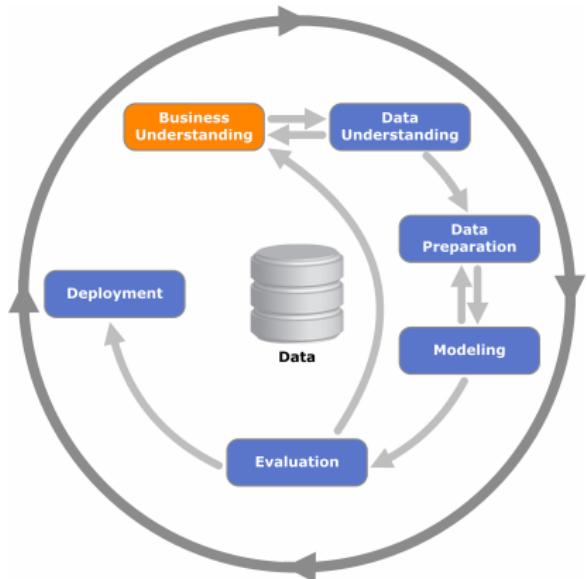
Задача: на рыболовном
предприятии автоматизировать
сортировку улова



(a) Сибас



(b) Лосось



Признаки

\mathcal{D} – множество объектов (data set)

$d \in \mathcal{D}$ – обучающий объект

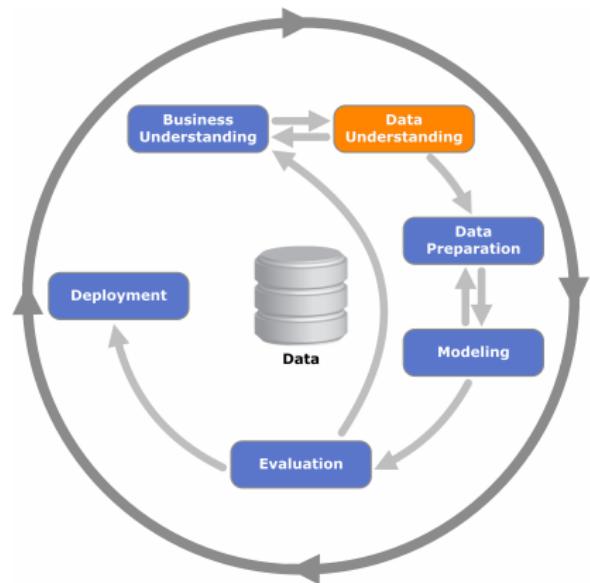
$\phi_j : D \rightarrow F_j$ – признак

Виды признаков

- ▶ Бинарные/Binary
 $F_j = \{true, false\}$
- ▶ Номинальные/Categorical
 F_j – конечно
- ▶ Порядковый/Ordinal
 F_j – конечно, упорядочено
- ▶ Количественный/Numerical
 $F_j = \mathbb{R}$

Представление d_i :

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

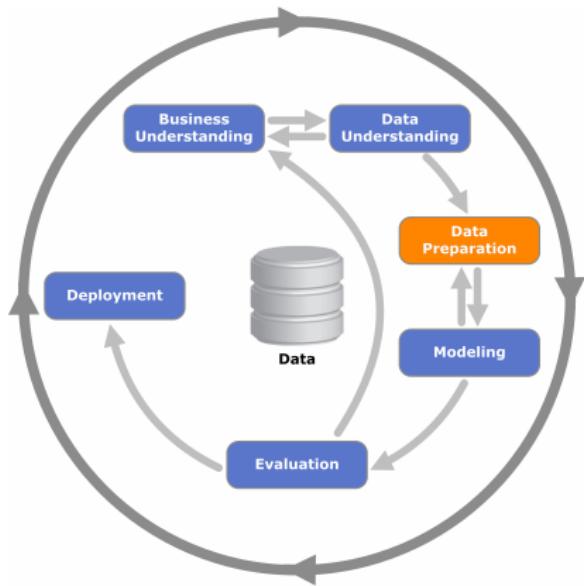


- ▶ Удаление шума
- ▶ Заполнение отсутствующих значений
- ▶ Трансформация факторов
- ▶ Выбор факторов
- ▶ Использование априорных знаний

В итоге:

$$X = \begin{pmatrix} x_1 \\ \dots \\ x_N \end{pmatrix}, \quad x_i \in \mathcal{X}$$

$$Y = \begin{pmatrix} y_1 \\ \dots \\ y_N \end{pmatrix}, \quad y_i \in \mathcal{Y}$$



Модель

семейство параметрических функций вида

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

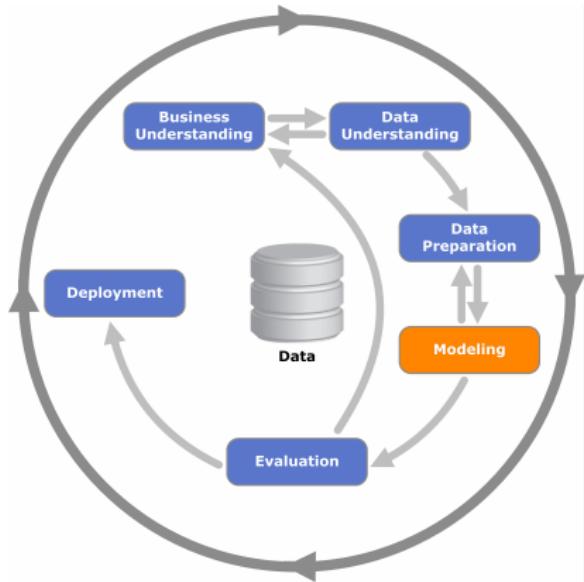
Алгоритм обучения

выбор наилучших параметров θ^*

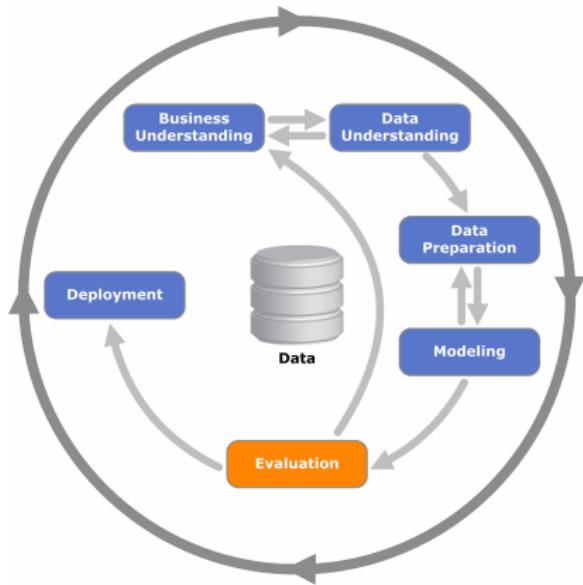
$$A(X, Y) : (\mathcal{X} \times \mathcal{Y})^N \rightarrow \Theta$$

В итоге:

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



- ▶ точность или аппроксимация?
- ▶ bias или variance?
- ▶ интерпретируемость или качество?





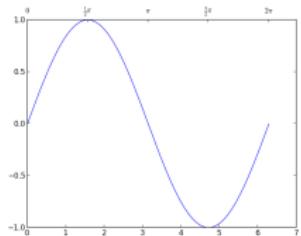
Exploratory data analysis

EDA направлен на предварительное изучение данных

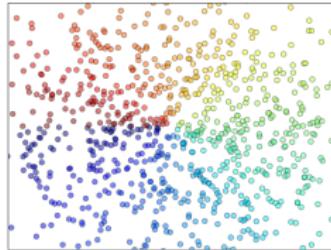
- ▶ формирование гипотез относительно структуры данных
- ▶ выбор необходимых инструментов анализа

Особенность метода состоит в визуализации и поиске важных характеристик и тенденций

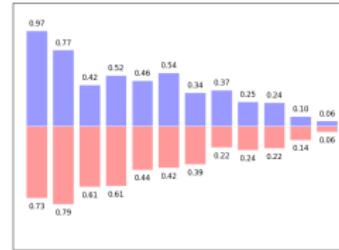
Примеры



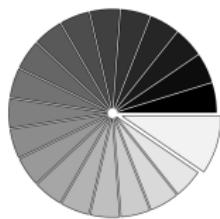
(a) Plot



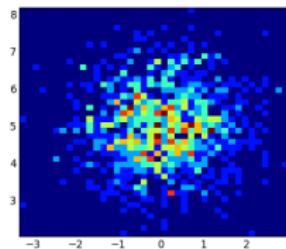
(b) Scatter



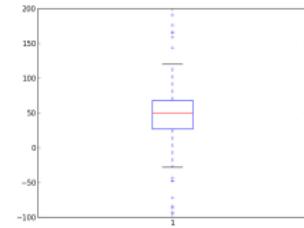
(c) Barplot



(d) Piechart



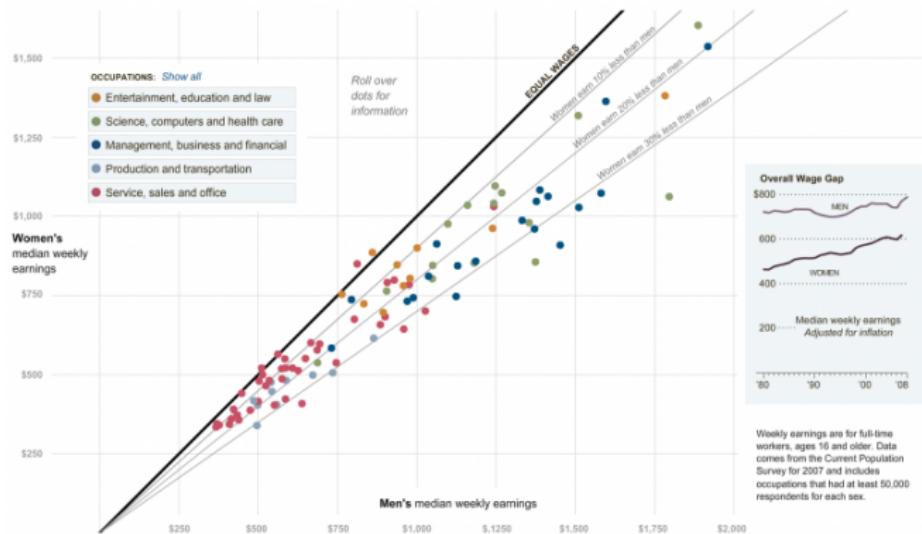
(e) Heatmap



(f) Boxplot

Полезные советы

- ▶ Все познается в сравнении
 - ▶ Причинно-следственные связи
 - ▶ Размерность имеет значение
(больше-лучше)
- ▶ Не избегать пояснений
 - ▶ Content is king



Задача

Дано: информация о посещении пользователями интернет-сайтов

Требуется: исследовать распределение количества уникальных пользователей в зависимости от популярности сайта

Пошаговая инструкция

- Скачать файл с данными (400M) <http://bit.ly/1xd0rHD>
- Скачать и запустить шаблон кода на python
<http://bit.ly/1B6fEc0>

```
$ python eda.py -h
$ python eda.py -n 20 eda.dat
```

- Заполнить функцию `top_domain_user_counts`
- Заполнить функцию `plot_log_top_domains` и посмотреть на распределение в log-log шкале
- Заполнить функцию `plot_linear_model`, попутно осознав параметры модели
- Выписать полученное распределение $n_{users}(rank)$

Вопросы

