



Анализ данных -
знакомство с одной из самых
сексапильных IT-специальностей 21
века

© Kleiner Igor M.Sc. 2015
Школа обработки и анализа данных

Добро пожаловать



<http://geekbrains.ru/events/118>

ברוכים הבאים

Добро пожаловать



<http://geekbrains.ru/events/118>







ברוכים הבאים

Будем знакомы










Будем знакомы



Учеба	   	















Будем знакомы



Учеба	   	  







Будем знакомы



Учеба	   	  
Работа	    	 

Будем знакомы



Учеба	   	  
Работа	    	 
Специализация	Image processing Big data \ machine learning Stochastic optimization	Psychology of perception

Цель нашей встречи

Цель нашей встречи

- o Где я?
- o Кто здесь?
- o А что вы тут делаете?



Цель нашей встречи

- Где я?
- Кто здесь?
- А что вы тут делаете?



Цель нашей встречи

◦ Познакомиться с вселенной анализа данных



Цель нашей встречи

- Познакомиться с вселенной анализа данных в интересной и доступной форме



Цель нашей встречи

- Познакомиться с вселенной анализа данных в интересной и доступной форме
- Узнать как самостоятельно продолжить обучение в области работы с данными и их анализом



Цель нашей встречи



- o Познакомиться с вселенной анализа данных в интересной и доступной форме
- o Узнать как самостоятельно продолжить обучение в области работы с данными и их анализом
- o Получить удовольствие и хорошо провести время

ГОТОВЫ?



Анализ данных – сексапильная специальность?

o hbr.org/2012/10/data-scientist-the-sexiest-job-of-the-21st-century/

Data Scientist: The Sexiest Job of the 21st Century

by **Thomas H. Davenport** and **D.J. Patil**

FROM THE OCTOBER 2012 ISSUE

 SUMMARY  SAVE  SHARE  COMMENT  TEXT SIZE  PRINT  BUY COPIES **\$8.95**

When Jonathan Goldman arrived for work in June 2006 at LinkedIn, the business networking site, the place still felt like a start-up. The company had just under 8 million accounts, and the number was

ARTICLE | HARVARD BUSINESS REVIEW | OCTOBER 2012

Data Scientist: The Sexiest Job of the 21st Century

The Sexiest Job of the 21st Century: Data Analyst

Chris Morris, Special to CNBC.com



Data Scientist

B.P.G.S Global Ltd

Jerusalem, IL • Sep 10, 2015 • From www.drushim.co.il

[Similar](#)

[View](#)



Data Scientist

Intel

Sep 16, 2015 • From intel.taleo.net

[Similar](#)

[View](#)



Data Scientist

Check Point Software Technologies, Ltd.

Jerusalem, IL • Sep 25, 2015 • From www.drushim.co.il

[Similar](#)

[View](#)



Data Scientist

Check Point Software Technologies, Ltd.

Tel-Aviv, Israel • Sep 16, 2015

▶ [24 connections to the poster](#) • [Similar](#)

[View](#)



Data Scientist

Check Point Software Technologies, Ltd.

Tel-Aviv, Israel • Sep 7, 2015

▶ [24 connections to the poster](#) • [Similar](#)

[View](#)



Data Scientist (Job 1066)

Viber

Israel • Aug 30, 2015

▶ [13 connections to the poster](#) • [Similar](#)

[View](#)

дата дата дата

- o Более чем 7.9 зетабайт электронной информации существует в мире сегодня

$\times 1024$

Zettabyte

$\times 1024$

Exabyte

$\times 1024$

Petabyte

Terabyte



Данные Данные Данные

- o Более чем 7.9 зетабайт электронной информации существует в мире сегодня
- o 7 900 000 000 000 000 000 000 байт – это количество примерно эквивалентно информации, содержащейся в более чем 600 миллиардов фильмов в HD качестве

Данные Данные Данные

- Более чем 7.9 зетабайт электронной информации существует в мире сегодня
- 7 900 000 000 000 000 000 000 байт – это количество примерно эквивалентно информации, содержащейся в более чем 600 миллиардов фильмов в HD качестве
- Новые данные появляются с экспоненциальной скоростью

WHERE IS DATA COMING FROM?

Twitter users send out

277,000
tweets

Google processes more than

2 million
search queries

Facebook processes almost

350 GB of data

72 hours

of new video are uploaded to YouTube

EVERY MINUTE...

Individuals and organizations launch

571

new websites

Walmart processes almost

17,000

transactions

More than

100 million

new emails are generated

Sprint processes more than

250,000

phone calls



Данные и их анализ очень ПОЛЕЗНЫ

- o Медицина
- o Спорт
- o Финансы
- o Корпорации
- o Государственные учреждения
- o СМИ
- o ...

Что такое анализ данных?

- Наука (или искусство) об использовании данных, с целью строить модели, которые позволяют принимать лучшие решения и приносят пользу



Что такое анализ данных?

- Наука (или искусство) об использовании данных, для того чтобы строить модели, которые позволяют принимать лучшие решения и приносят пользу



“Science is what we understand well enough to explain to a computer. Art is everything else we do”




Примеры успешного применения анализа данных

- o eHarmony
- o The Framingham Heart Study - фремингемское исследование сердца
- o Выбор игроков в команду

eHarmony

- Сайт знакомств:
- *модус операнди*: создание пар для долгосрочных отношений
- *идея*: научный подход для поиска подходящих кандидатур
- нет поиска по анкетам

 Знакомства Мамба - крупнейший бесплатный сайт...
mamba.ru



eHarmony



- o Сайт знакомств:
 - o *модус операнди*: создание пар для долгосрочных отношений
 - o *идея*: научный подход, для поиска подходящих кандидатур
 - o нет поиска по анкетам
- o Общая прибыль сайта превысила 1 миллиард долларов
- o Около 4% браков в США это результат eHarmony

www.eharmony.com/press-release/31/



eHarmony

o Сайт знакомств:

пользователь при регистрации заполняет длинную анкету



eHarmony



o Сайт знакомств:

пользователь при регистрации заполняет длинную анкету



o проанализировав данные, сайт выдает подходящие анкеты пользователей, проживающих рядом

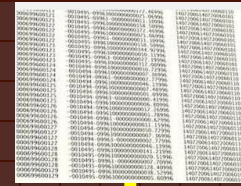
eHarmony

o Магия анализа данных



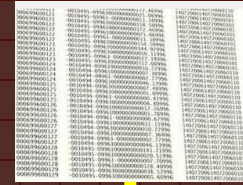
eHarmony

o Магия анализа данных



eHarmony

o Магия анализа данных



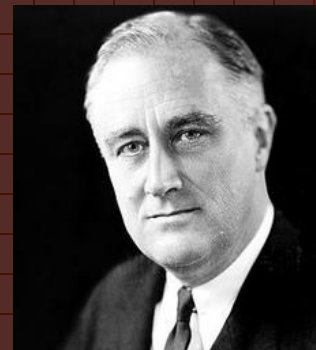
The Framingham Heart Study

- Исследование длится более 65 лет и является одним из самых продолжительных эпидемиологических исследований в истории медицины (фремингемское исследование сердца)



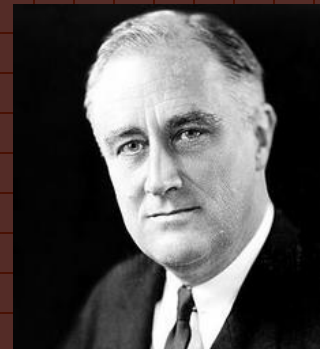
фремингемское исследование сердца

- Франклин Делано Рузвельт президент США 1933-1945
- Умер во время исполнения своих обязанностей в 1945



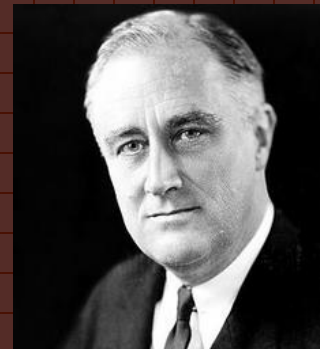
фремингемское исследование сердца

- Франклин Делано Рузвельт президент США 1933-1945
- Умер во время исполнения своих обязанностей в 1945
- Давление до 1933 года **140/100** - сегодня считается высоким давлением
- Давление за год до смерти **210/120** - сегодня считается опасным кризисом
- **260/150** давление за два месяца до смерти
- **300/190** в день смерти



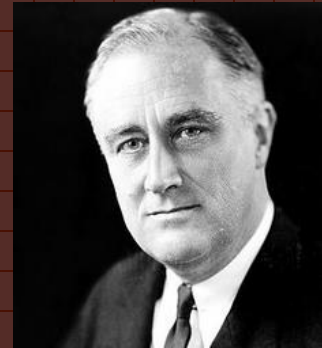
фремингемское исследование сердца

- Сегодня мы знаем об опасности высокого давления
- Откуда сегодня врачам известна эта информация?



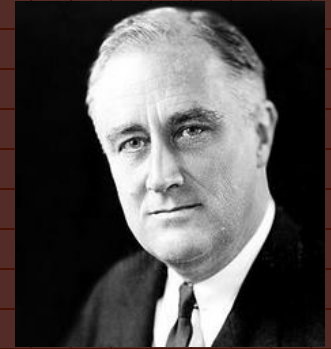
фремингемское исследование сердца

- Сегодня мы знаем об опасности высокого давления
- Откуда сегодня врачам известна эта информация?

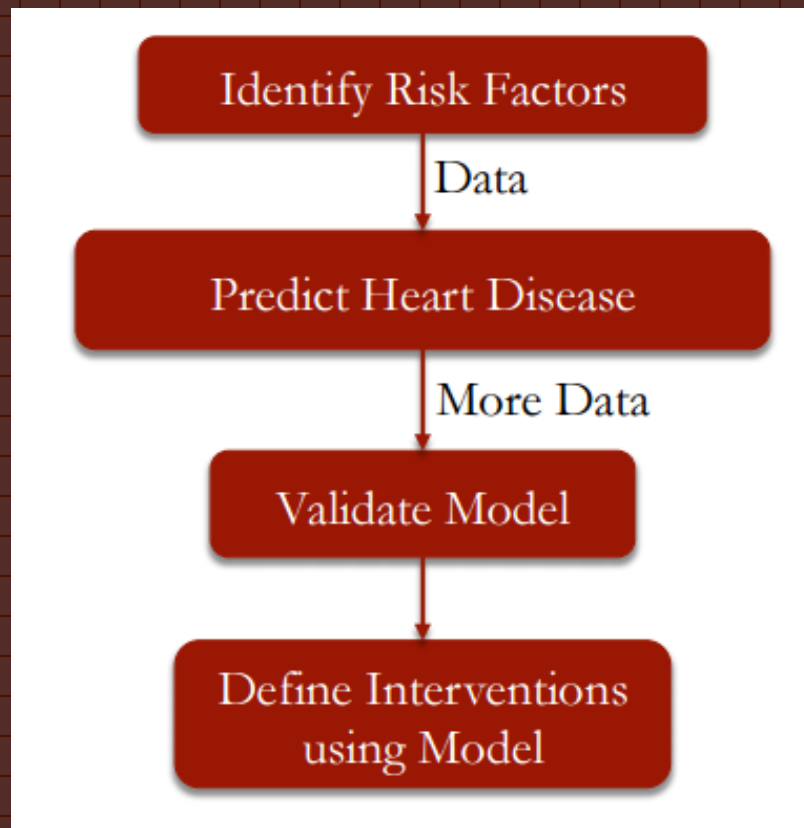


фремингемское исследование сердца

- 1948 год, город Фремингем
- 5209 участников
- участие в наблюдениях и тестах в течении длительного времени
- **Цель:** выявление факторов риска для болезней сердца



фремингемское исследование сердца



фремингемское исследование сердца

- Благодаря полученным данным и последующим исследованиям, учеными были обнаружены различные факторы риска:

курение

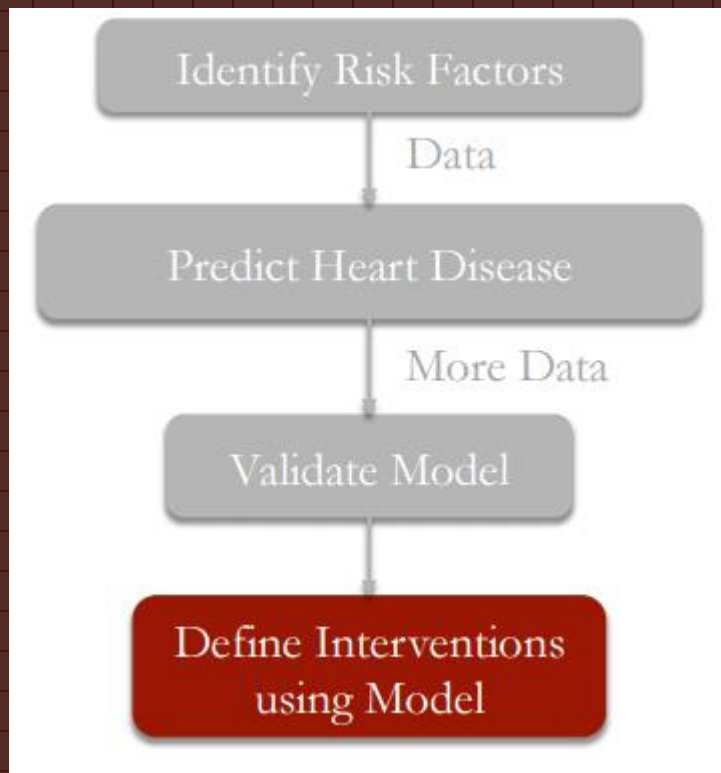
уровень холестерина

давление

уровень сахара в крови

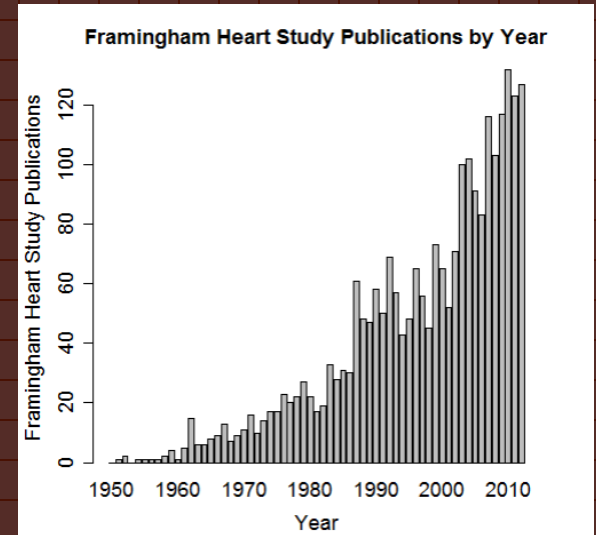
...

фремингемское исследование сердца



фремингемское исследование сердца

- Более 2400 исследований на основе полученных данных
- Выявление множества факторов риска



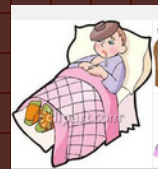
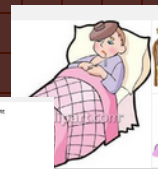
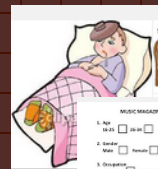
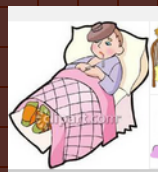
фремингемское исследование сердца

o Магия анализа данных



фремингемское исследование сердца

Магия анализа данных



MUSIC MAGAZINE QUESTIONNAIRE

1. Sex M F

2. Grade Fresh

3. Occupation Unemployed

4. Income (per month) \$100-\$199 \$200-\$299 \$300-\$399 \$400-\$499 \$500-\$599 \$600-\$699 \$700-\$799 \$800-\$899 \$900-\$999 \$1000 or more Other

5. What type of magazine do you read most often? News Sports Science Other

6. How often do you read a magazine? Daily Weekly Monthly Other

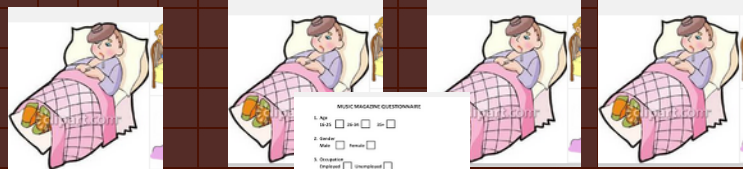
7. How much money do you spend on magazines each month? \$0 \$1-\$5 \$6-\$10 \$11-\$15 \$16-\$20 \$21-\$25 \$26-\$30 \$31-\$35 \$36-\$40 \$41-\$45 \$46-\$50 \$51-\$55 \$56-\$60 \$61-\$65 \$66-\$70 \$71-\$75 \$76-\$80 \$81-\$85 \$86-\$90 \$91-\$95 \$96-\$100 Other

8. If you have spent this much money on magazines each month, how many magazines do you read each month? None One Two Three Four Five Six Seven Eight Nine Ten Eleven Twelve Other



фремингемское исследование сердца

Магия анализа данных



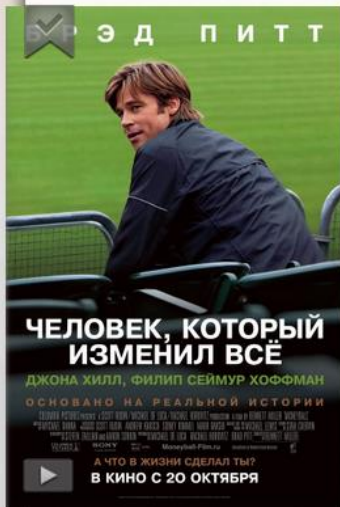
логистическая регрессия



Выбор лучшего игрока в команду moneyball

Человек, который изменил всё

Moneyball



год	2011
страна	США
слоган	«А что в жизни сделал ты?»
режиссер	Беннетт Миллер
сценарий	Стивен Зеллиан , Аарон Соркин , Стэн Червин , ...
продюсер	Майкл Де Лука , Рэйчел Хоровиц , Брэд Питт , ...
оператор	Уолли Пфистер
композитор	Майкл Дэнна
художник	Джесс Гончор , Брэд Рикер , Дэвид Скотт , ...
монтаж	Кристофер Теллефсен
жанр	драма , биография , спорт , ...

В главных ролях:

[Брэд Питт](#)
[Джона Хилл](#)
[Филип Сеймур Хоффман](#)
[Крис Пратт](#)
[Стефен Бишоп](#)
[Кэррис Дорси](#)
[Робин Райт](#)
[Рид Даймонд](#)
[Брент Дженингс](#)
[Кен Медлок](#)
...

Роли дублировали:

[Всеволод Кузнецов](#)
[Диомид Виноградов](#)
[Александр Новиков](#)
[Дмитрий Давыдов](#)

Ссылка

558 000 000

слова

Бейсбол для чайников

o <https://goo.gl/FS7uPg>

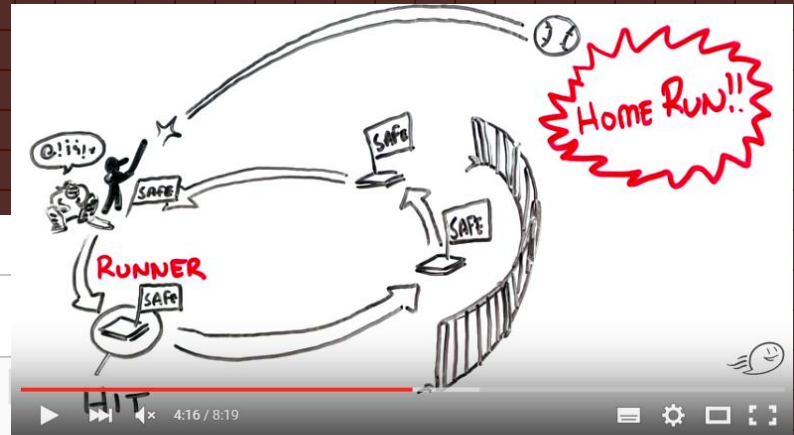
Бейсбол

Материал из Википедии — свободной энциклопедии

Бейсбóл (англ. *baseball*, от *base* — «база, основание» и *ball* — «мяч») — командная спортивная игра с бейсбольным мячом и битой. В состязаниях участвуют две команды по девять (иногда десять) игроков каждая.

Бейсбол наиболее популярен на Кубе, в США, в Венесуэле, в Японии, Китае и Южной Корее. В США, Японии, Чехии и других странах распространён также софтбол — упрощённый вариант бейсбола — игра, которую можно проводить в помещении и на небольших полях. На данный момент в бейсбол играют более чем в 120 странах мира.

К родственным бейсболу видам спорта относятся крикет, песаполо в Финляндии, ойна в Румынии и лапта в России.

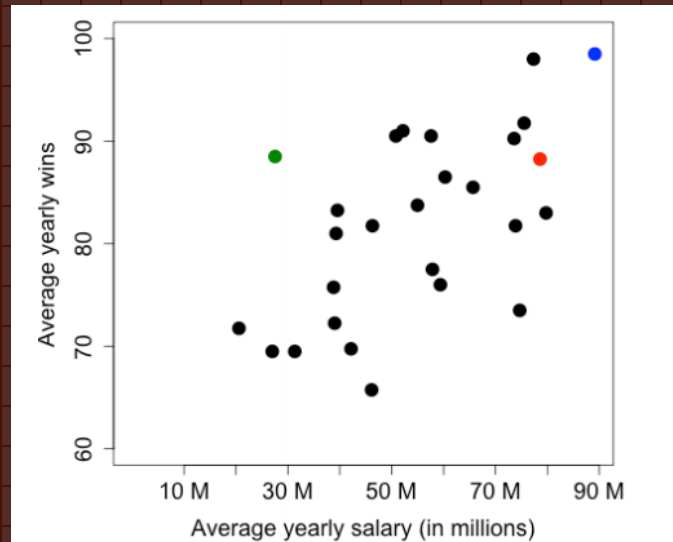


Baseball Rules Whiteboard Video Rules of Baseball



Выбор лучшего игрока в команду

- У богатых команд больше денег и они могут позволить купить лучших игроков



Выбор лучшего игрока в команду

○ У богатых команд больше денег и они могут позволить купить лучших игроков



○ Оклендская бедная команда после прихода нового менеджера стала показывать хорошие результаты

Выбор лучшего игрока в команду

○ У богатых команд больше денег и они могут позволить купить лучших игроков



○ Оклендская бедная команда после прихода нового менеджера стала показывать хорошие результаты

○ Что произошло?

Выбор лучшего игрока в команду

○ У богатых команд больше денег и они могут позволить купить лучших игроков



○ Оклендская бедная команда после прихода нового менеджера стала показывать хорошие результаты

○ Что произошло?



Выбор лучшего игрока в команду

○ У богатых команд больше денег и они могут позволить купить лучших игроков



○ Обработав множество параметров игроков, программа выявила тех игроков, которые были недооценены, т.е. качество игры которых было меньше заработка игроков их уровня



Выбор лучшего игрока в команду

- У богатых команд больше денег и они могут позволить купить лучших игроков
- Обработав множество параметров игроков, программа выявила тех игроков, которые были недооценены, т.е. качество игры которых было меньше заработка игроков их уровня
- Сегодня в любой команде высшей лиги есть свой статистик



Выбор лучшего игрока в команду

o Магия анализа данных



логистическая регрессия



Еще примеры

- предсказание решений высшего суда
- предсказание будущей цены вина
- предсказания цен на авиабилеты
- предсказания эпидемии гриппа на основе поисковых запросов
- ...
- ...

Еще примеры

- предсказание решений высшего суда
- предсказание будущей цены вина
- предсказания цен на авиабилеты
- предсказания эпидемии гриппа на основе поисковых запросов
- ...
- ...



предсказания эпидемии гриппа на основе поисковых запросов

Google предупредит о вспышке эпидемии гриппа

13.11.2008 12:31 Компания Google запустила новый интернет-сервис Google Flu Trends, который будет отслеживать вспышки эпидемии гриппа на территории Соединенных Штатов. Поисковик будет регистрировать все пользовательские запросы, связанные с гриппом, в частности вопро ...

По материалам: Радио Свобода: новости



предсказания эпидемии гриппа на основе поисковых запросов

Google предупредит о вспышке эпидемии гриппа

13.11.2008 12:31 Компания Google запустила новый интернет-сервис Google Flu Trends, который будет отслеживать вспышки эпидемии гриппа на территории Соединенных Штатов. Поисковик будет регистрировать все пользовательские запросы, связанные с гриппом, в частности, вопро ...

По матери

When Google got flu wrong

US outbreak foxes a leading web-based method for tracking seasonal flu.

Declan Butler

13 February 2013



PDF



Rights & Permissions



Дайте мне точку опоры



Дайте мне данные



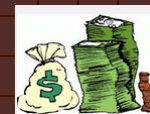
Дайте мне данные, компьютер



Дайте мне данные, компьютер и специалиста по анализу данных



Дайте мне данные,
компьютер и специалиста
по анализу данных и я
изменю жизнь людей



Магия анализа данных

данные



закономерности - предсказания

Магия анализа данных

данные



закономерности - предсказания

Магия анализа данных

данные



закономерности - предсказания

Детали дьявола анализа данных

- Поиск и сбор необходимых данных
- Приведение полученных данных в удобную для анализа форму
- Выбор подходящей модели для анализа данных
- Анализ данных
- Верификация полученных результатов
- Презентация полученных результатов и принятие решений

Детали дьявола анализа данных

- Define the question
- Define the ideal data set
- Determine what data you can access
- Obtain the data
- Clean the data
- Exploratory data analysis

- Statistical prediction/modeling
- Interpret results
- Challenge results
- Synthesize/write up results
- Create reproducible code
- Distribute results to other people

Краткие итоги

o Анализ данных это:

Краткие итоги

- Анализ данных это:
 - интересно
 - полезно
 - прибыльно

Вопросы

- Анализ данных
 - как изучить
 - что изучить
 - можно ли изучать самостоятельно
 - какие есть направления развития
 - что включает в себя анализ данных

Как изучать анализ данных

- o Университет
- o Онлайн курсы \ онлайн специализации
- o Учебная литература \ интернет

Как изучать анализ данных

- Университет:

- фундаментальные знания

- долго

- 1-4 релевантных курса за все время обучения

Как изучать анализ данных

- Университет
- Онлайн курсы \ онлайн специализации
 - множество бесплатных курсов
 - не все курсы одинаково хороши и полезны
 - есть очень хорошие курсы и специализации

Онлайн курсы

- o EDX – MIT – «Меч Аналитики»
- o edx.org/course/analytics-edge-mitx-15-071x-0
- o Достоинства курса:
 - o множество интересных примеров
 - o минимум теории,
 - o максимум практики
- o Язык программирования R



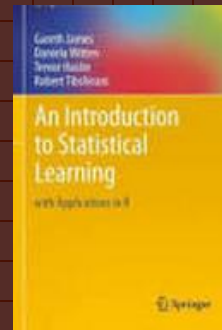
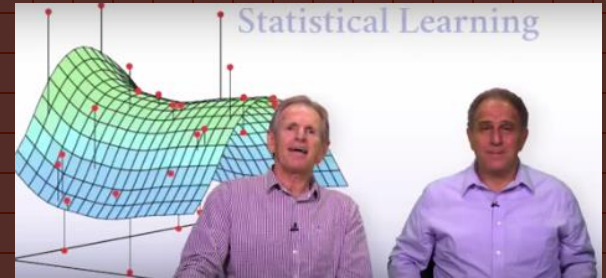
The Analytics Edge

Through inspiring examples and stories, discover the power of data and use analytics to provide an edge to your career and your life.



Онлайн курсы

- o Stanford – Statistical Learning
- o lagunita.stanford.edu/courses/HumanitiesandScience/StatLearning/Winter2015/about
- o Достоинства курса:
 - o немного теории
 - o много практики
 - o хороший учебник по курсу
- o Язык программирования: R



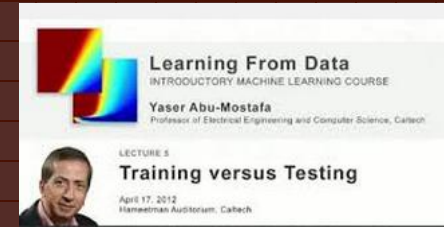
Онлайн курсы

- o Coursera, Stanford – Machine Learning
- o *coursera.org/learn/machine-learning*
- o Достоинства курса:
 - o удачное сочетание практики и теории
- o Язык программирования: Matlab, Octave



Онлайн курсы

- o Edx, Caltech, Learning from data
- o <https://work.caltech.edu/telecourse.html>
- o Достоинства курса:
 - o твердый теоретический фундамент
 - o основные теоретические моменты объяснены в интересной и доступной форме
 - o хороший учебник сопровождающий курс



«Анти онлайн курсы»

- Записи лекций
- Курс «Машинное обучение»
- Преподаватель — Константин Вячеславович Воронцов.

Онлайн специализации

o Coursera, Machine Learning Specialization

coursera.org/specializations/machine-learning

Язык программирования: *Питон*

1 Machine Learning Foundations: A Case Study Approach

2 Regression

3 Classification

4 Clustering & Retrieval

5 Recommender Systems & Dimensionality Reduction

6 Machine Learning Capstone: An Intelligent Application with Deep Learning

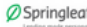






Онлайн специализации

- o Coursera, Big Data Specialization
coursera.org/specializations/big-data

- 1 Introduction to Big Data
- 2 Hadoop
- 3 Introduction to Big Data Analytics
- 4 Machine Learning With Big Data
- 5 Introduction to Graph Analytics
- 6 Big Data - Capstone Project

Соревнования

 www.kaggle.com

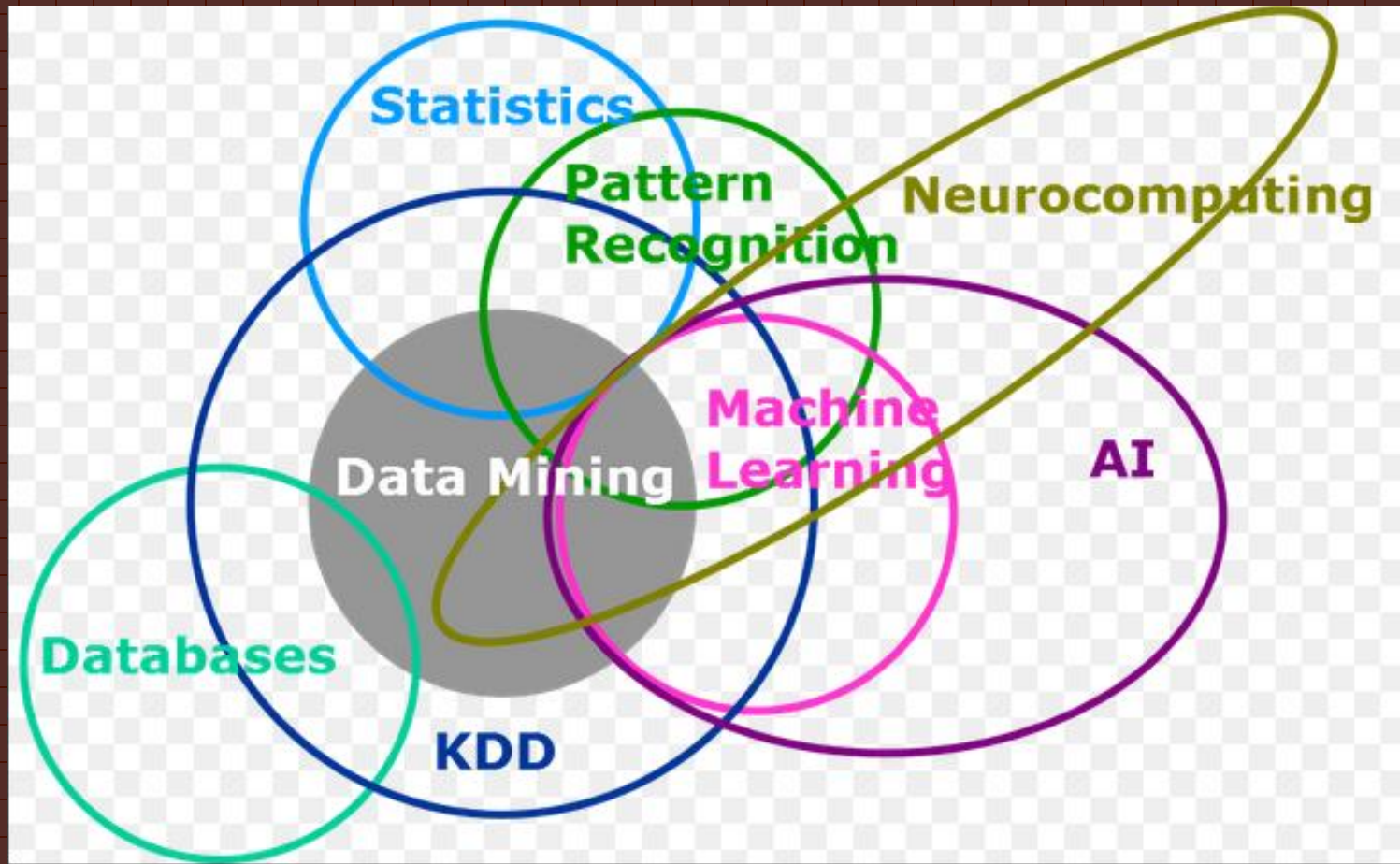
	Springleaf Marketing Response Determine whether to send a direct mail piece to a customer
	Western Australia Rental Prices  Predict rental prices for properties across Western Australia
	Coupon Purchase Prediction Predict which coupons a customer will buy
	Flavours of Physics: Finding $\tau \rightarrow \mu\mu$ Identify a rare decay phenomenon
	Truly Native? Predict which web pages served by StumbleUpon are sponsored
	Right Whale Recognition Identify endangered right whales in aerial photographs

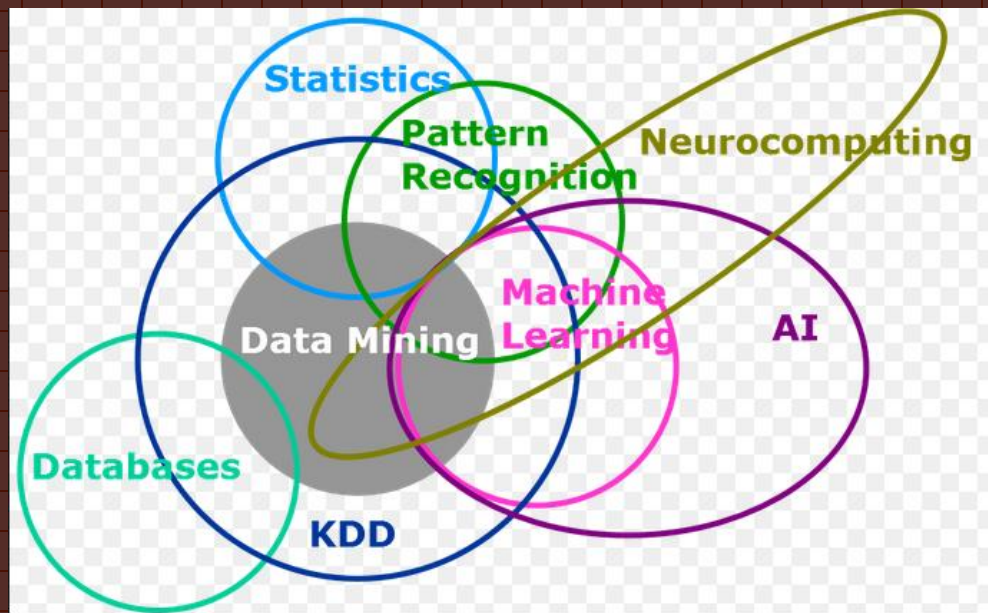
Вопросы

o Что в себя включает анализ данных?

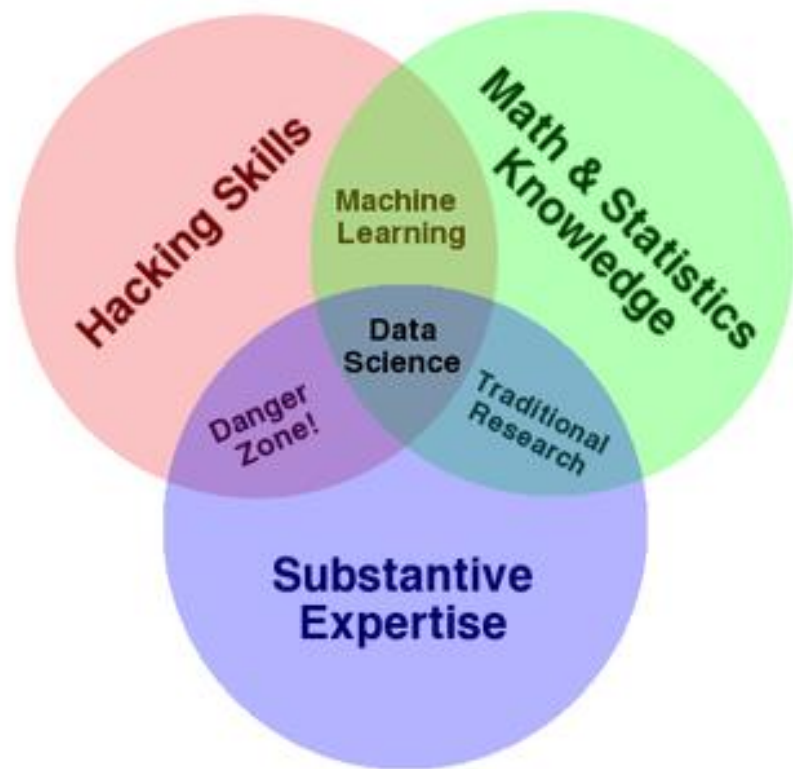
Вопросы

- Что в себя включает анализ данных?
- Анализ данных лежит на пересечениях множества областей наук





теория вероятностей
статистика
случайные процессы
структуры данных
ИИ
базы данных
параллельные вычисления
оптимизация
выпуклая оптимизация
линейное программирование
алгоритмы
структуры данных
финансовое моделирование
...
...



[Drew Conway](#)



«Каждая кухарка может управлять государством»

Ленин



«Каждая кухарка может управлять государством»

Ленин

Может ли каждая кухарка научиться
анализировать данные?

Вопрос

- На каких языках лучше заниматься анализом данных?

Вопрос

- На каких языках лучше заниматься анализом данных?
 - нет жестких правил
 - R
 - Python
 - Matlab \ Octave
 - C++
 - ...

Вопрос

- Необходимы ли навыки программирования для анализа данных?

Вопрос

- Необходимы ли навыки программирования для анализа данных?
- нет, но желательны

Вопрос

- Необходимы ли навыки программирования для анализа данных?
- нет, но желательны
- Существуют специальные программы позволяющие анализировать данные без знаний программирования

Вопрос

- Необходимы ли навыки программирования для анализа данных?
- Существуют специальные программы позволяющие анализировать данные без знаний программирования:
 - Excel, Rattle, RapidMiner и другие

Анализ данных на Excel

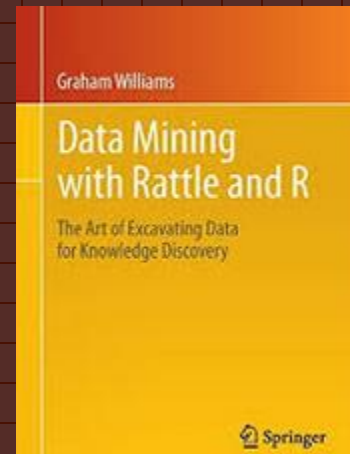


The image shows a central graphic with a green background. At the top, the Microsoft Excel logo (a white 'X' on a green square) is displayed next to the word 'Excel' in white. Below the logo, there are several devices: a desktop monitor, a laptop, and a tablet, all displaying colorful pie charts and data tables. A green banner with the word 'VERIFIED' in white capital letters is positioned below the devices. To the right of the banner is a green ribbon icon with a white checkmark.

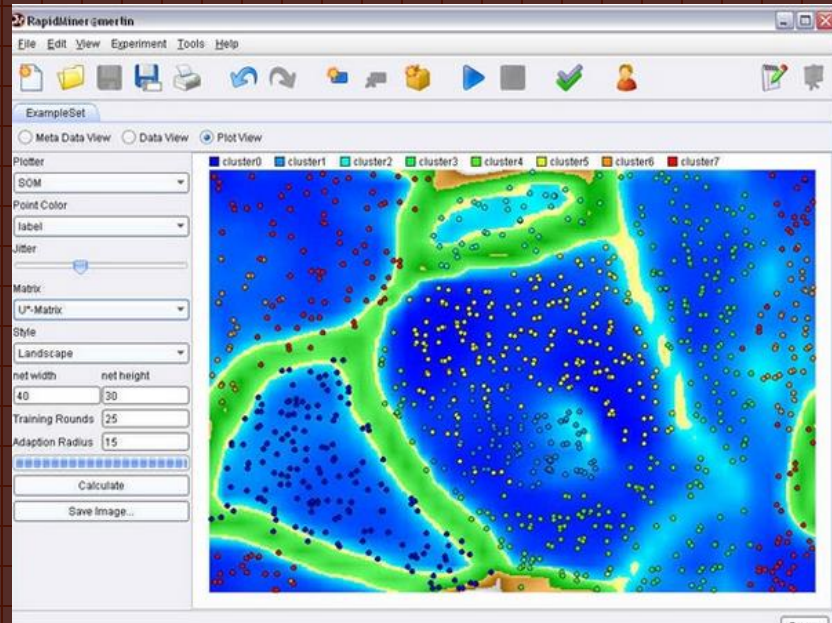
Microsoft
DAT206x

Excel for Data Analysis and
Visualization

Анализ данных на Rattle



Анализ данных RapidMiner



Анализ данных под микроскопом за 5 минут

Machine Learning



Анализ данных под микроскопом за 5 минут

1. Данные:

id	problem_id	subject_id	start	stop	time_left	answ
2	1	498	17 1307119989	1307120016	2369	A
3	2	150	15 1307119991	1307120009	2376	D
4	3	313	18 1307119994	1307120009	2378	E
5	4	12	13 1307119995	1307120019	2366	B
6	5	273	14 1307119996	1307120028	2357	A
7	6	101	19 1307119996	1307120021	2364	B
8	7	105	18 1307119998	1307120048	2337	B
9	8	162	12 1307120004	1307120042	2343	C
10	9	70	15 1307120011	1307120038	2347	C
11	10	300	16 1307120012	1307120092	2293	B
12	11	494	17 1307120017	1307120075	2310	D
13	12	357	13 1307120021	1307120118	2267	A
14	13	522	19 1307120025	1307120152	2233	D
15	14	232	14 1307120030	1307120158	2227	C
16	15	344	15 1307120041	1307120117	2268	B
17	16	160	17 1307120079	1307120249	2138	D
18	17	516	16 1307120094	1307120159	2226	B
19	18	472	12 1307120119	1307120170	2215	A
20	19	43	15 1307120122	1307120140	2245	C
21	20	353	13 1307120144	1307120199	2186	C
22	21	218	15 1307120152	1307120272	2113	E
23	22	69	16 1307120163	1307120188	2197	D
24	23	562	16 1307120190	1307120301	2084	D
25	24	121	19 1307120253	1307120294	2091	E
26	25	297	15 1307120277	1307120342	2043	B
27	26	495	13 1307120281	1307120353	2032	E
28	27	94	14 1307120288	1307120343	2042	E
29	28	22	18 1307120310	1307120365	2020	C
30	29	64	19 1307120310	1307120385	2000	B
31	30	502	16 1307120323	1307120336	2049	B
32	31	44	16 1307120339	1307120352	2033	A
33	32	315	14 1307120348	1307120362	2023	B
34	33	285	15 1307120352	1307120553	1832	E
35	34	550	13 1307120356	1307120444	1941	B
36	35	92	14 1307120368	1307120397	1988	B
37	36	395	16 1307120377	1307120426	1959	D
38	37	267	17 1307120382	1307120515	1870	E
39	38	257	14 1307120401	1307120427	1958	C
40	39	312	19 1307120407	1307120548	1837	D
41	40	321	18 1307120431	1307120449	1936	A
42	41	220	16 1307120437	1307120510	1874	A

Анализ данных под микроскопом за 5 минут

1. Где взять данные?



Анализ данных под микроскопом за 5 минут

1. Где взять данные?

- United Nations <http://data.un.org/>
- U.S. <http://www.data.gov/>
 - [List of cities/states with open data](#)
- United Kingdom <http://data.gov.uk/>
- France <http://www.data.gouv.fr/>
- Ghana <http://data.gov.gh/>
- Australia <http://data.gov.au/>
- Germany <https://www.govdata.de/>
- Hong Kong <http://www.gov.hk/en/theme/psi/datasets/>
- Japan <http://www.data.go.jp/>

Анализ данных под микроскопом за 5 минут

1. Где взять данные?

- [Stanford Large Newtork Data](#)
- [UCI Machine Learning](#)
- [KDD Nugets Datasets](#)
- [CMU Statlib](#)
- [Gene expression omnibus](#)
- [ArXiv Data](#)
- [Public Data Sets on Amazon Web Services](#)

Анализ данных под микроскопом за 5 минут

◦ Перед анализом данные надо подготовить:

Анализ данных под микроскопом за 5 минут

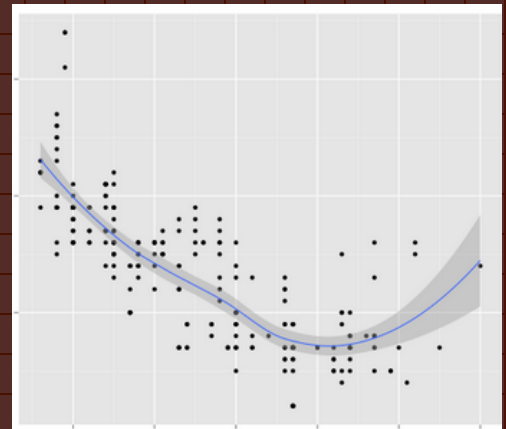
- Перед анализом данные надо подготовить:
 - препроцессинг: нормализация, ...
 - обработать отсутствующие значения
 - привести данные в удобный для анализа вид

○ TIDY DATA

messy				tidy		
	treatmenta	treatmentb		name	trt	result
John Smith	—	2		John Smith	a	—
Jane Doe	16	11		Jane Doe	a	16
Mary Johnson	3	1		Mary Johnson	a	3
	John Smith	Jane Doe	Mary Johnson	John Smith	b	2
treatmenta	—	16	3	Jane Doe	b	11
treatmentb	2	11	1	Mary Johnson	b	1

Анализ данных под микроскопом за 5 минут

- Предварительное знакомство с данными с помощью графической информации может помочь в дальнейшем анализе
- Exploratory Data Analysis



Анализ данных под микроскопом за 5 минут

- Анализ данных, выбор лучшей модели,
тестирование качества выбранной модели,
использование модели на новых данных

Анализ данных под микроскопом за 5 минут

o Анализ данных, выбор лучшей модели

<u>Unsupervised</u>	<u>Supervised</u>
<ul style="list-style-type: none">• Clustering & Dimensionality Reduction<ul style="list-style-type: none">o SVDo PCAo K-means	<ul style="list-style-type: none">• Regression<ul style="list-style-type: none">o Linearo Polynomial• Decision Trees• Random Forests
<ul style="list-style-type: none">• Association Analysis<ul style="list-style-type: none">o Apriorio FP-Growth• Hidden Markov Model	<ul style="list-style-type: none">• Classification<ul style="list-style-type: none">o KNNo Treeso Logistic Regressiono Naive-Bayeso SVM

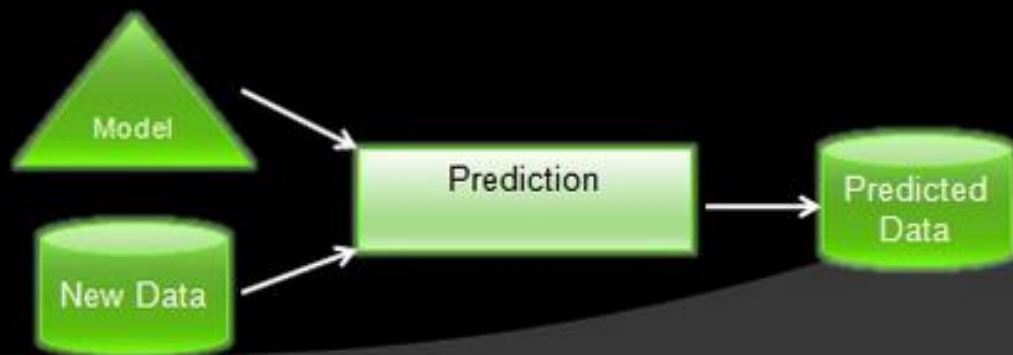
Анализ данных под микроскопом за 5 минут

- o Анализ данных, выбор лучшей модели
 - o регрессия линейная
 - o дискриминантный анализ
 - o логистическая регрессия
 - o сплайны
 - o случайные деревья
 - o случайные леса
 - o РСР
 - o метод опорных векторов
 - o бустинг
 - o метод ближайших соседей
 - o ...
 - o ...

Phase 1) Learning



Phase 2) Prediction



Вопрос

- Как вы думаете какой этап занимает больше всего времени? (как правило)
 - скачать данные
 - подготовить данные к анализу
 - выбор лучшей модели
 - представление результатов анализа

Вопрос

- Как вы думаете какой этап занимает больше всего времени? (как правило)
 - скачать данные
 - **подготовить данные к анализу**
 - выбор лучшей модели
 - представление результатов анализа

Пример модели для анализа данных

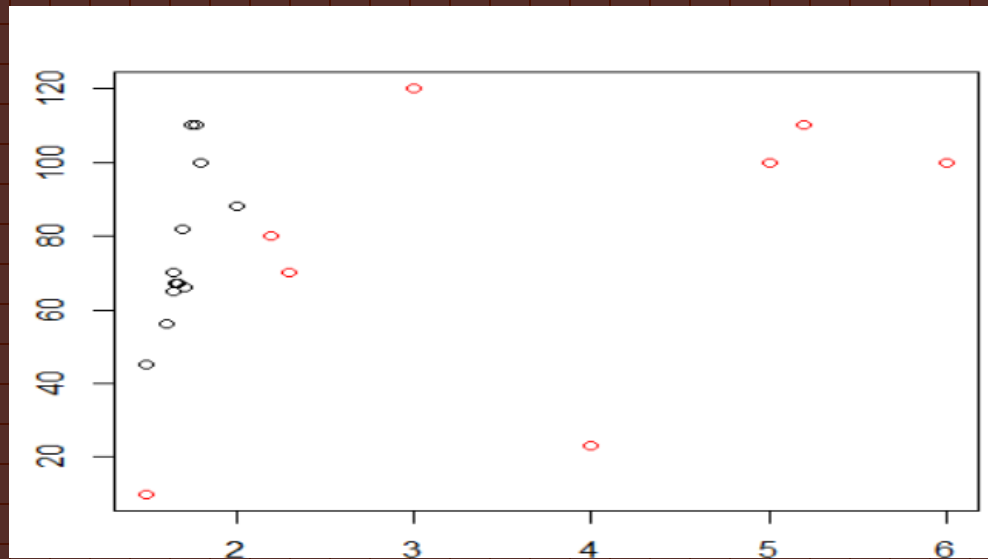
Задача классификации

Цель: научить систему различать людей и пришельцев с марса



Человек или пришелец

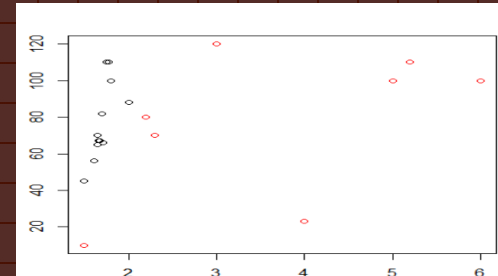
Данные для обучения: таблица содержащая рост и вес 20 кандидатов, 12 из которых люди и 8 из которых пришельцы с Марса

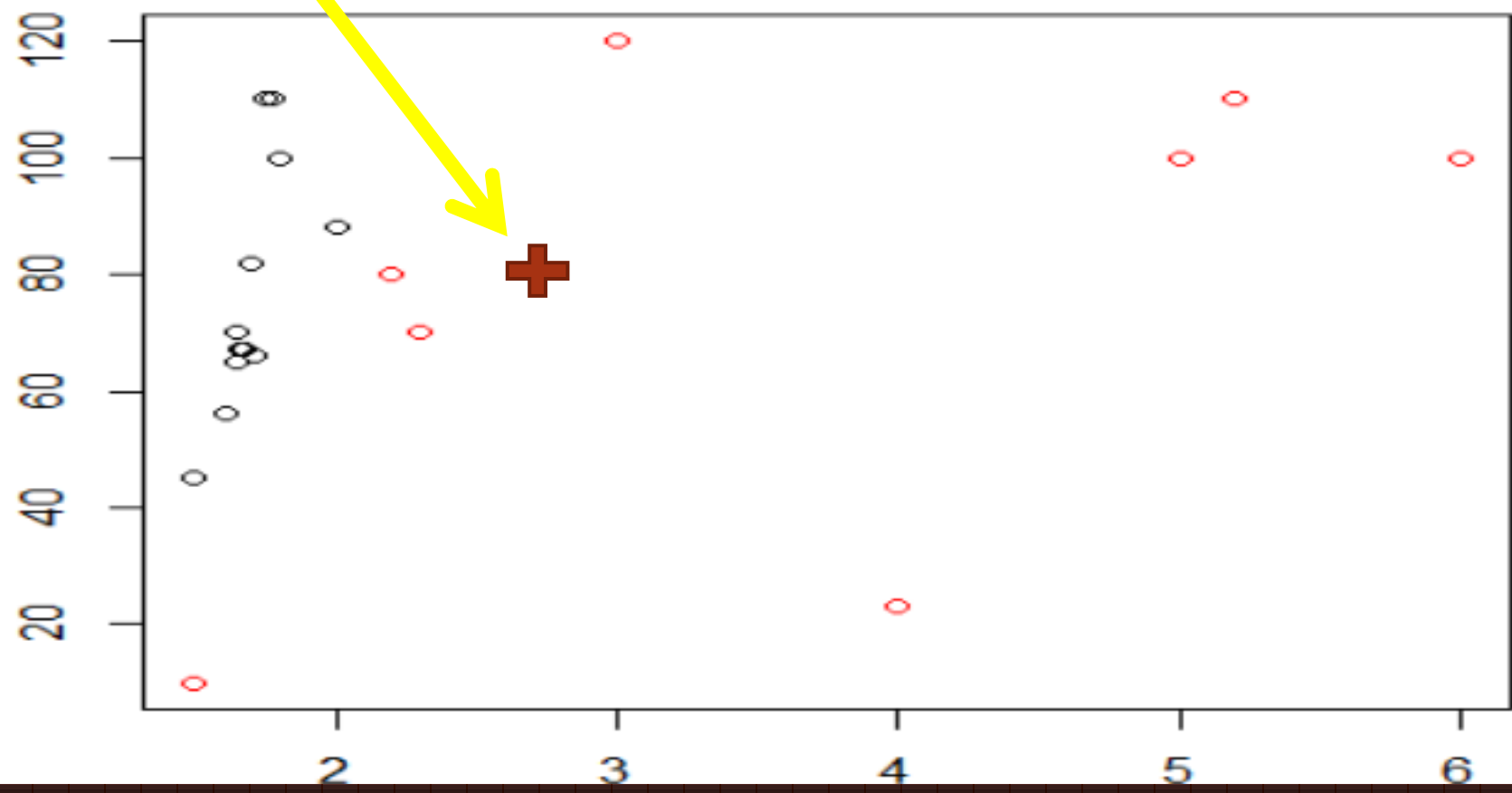


Человек или пришелец

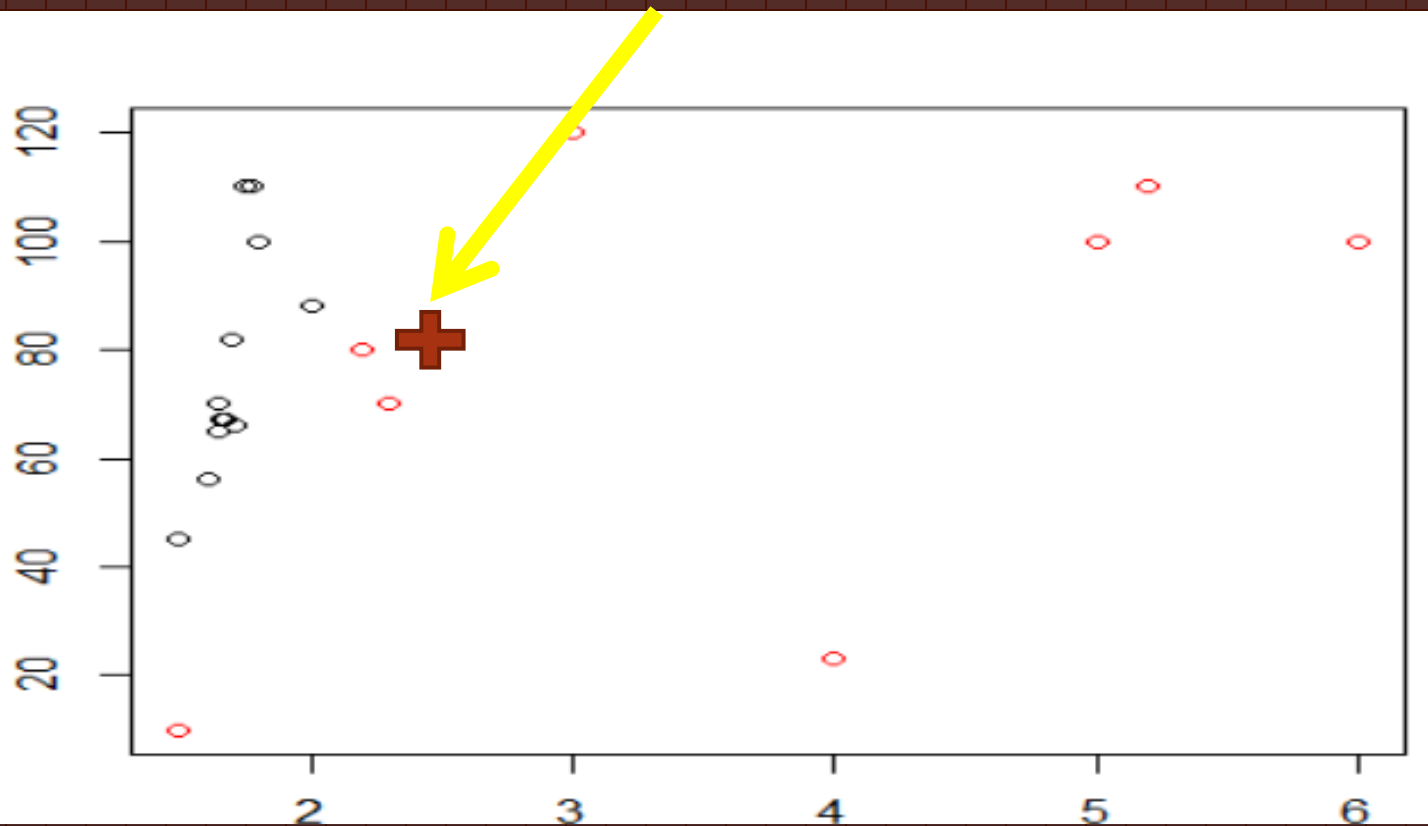
Данные для обучения: таблица содержащая рост и вес 20 кандидатов, 12 из которых люди и 8 из которых пришельцы с Марса

Перед нами новый персонаж, как понять человек это или пришелец?



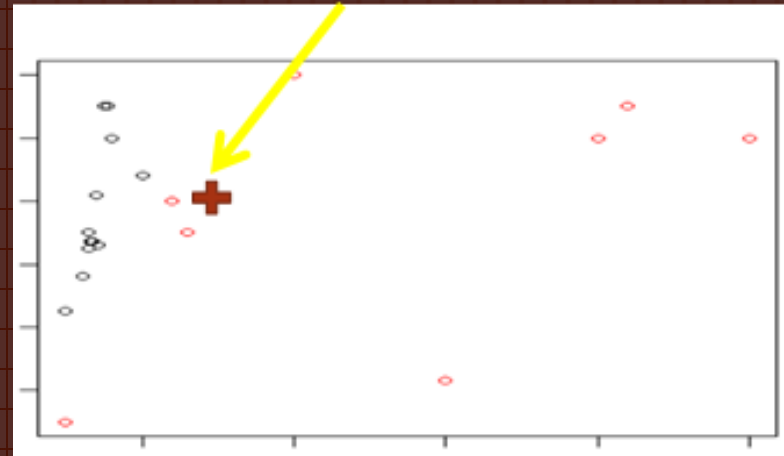


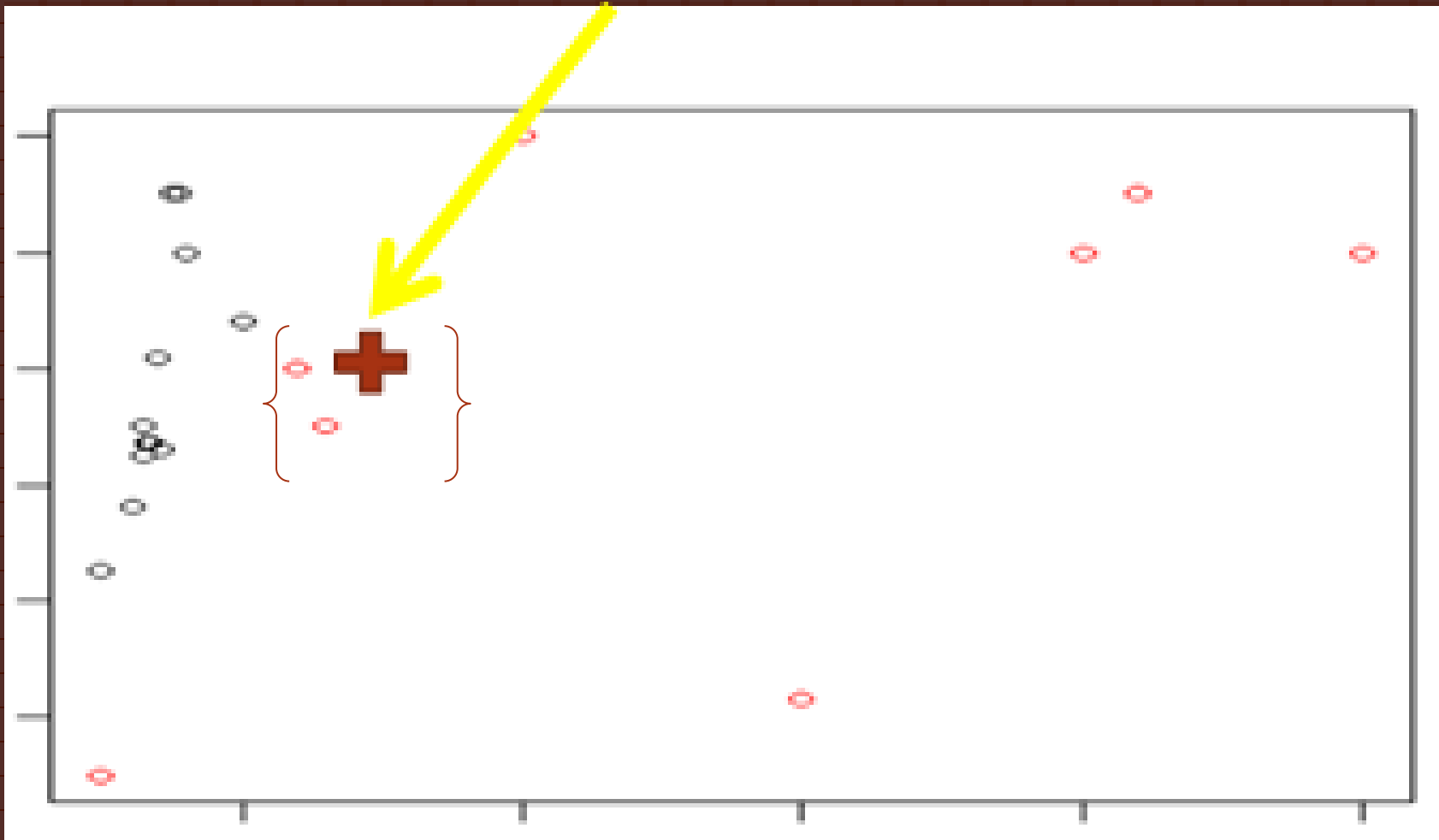
Кто это?

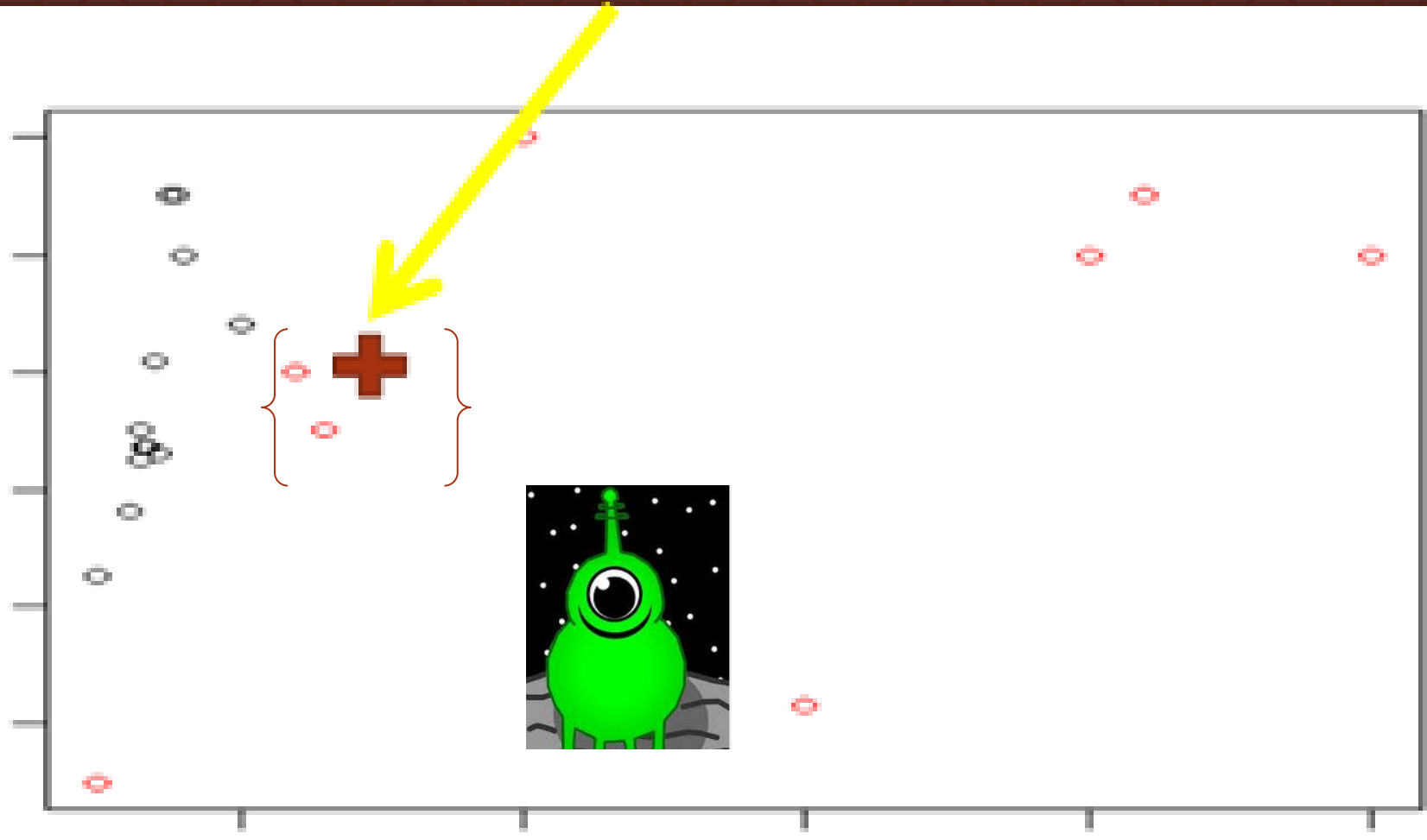


Метод ближайшего соседа

- «скажи мне кто твой друг и я скажу кто ты»
- «дурак дурака видит издалека»
- «Каковы соседи, таков и ты»





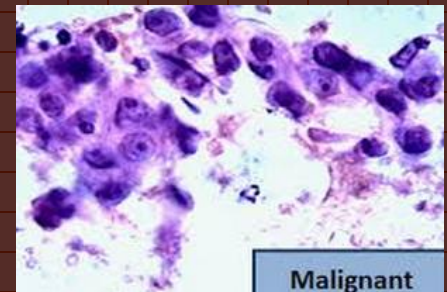


Метод ближайших соседей

- Несмотря на свою простоту, метод достаточно эффективен при определенных условиях
- Метод хуже работает в пространствах высокой размерности (почему?)

Пример

- Анализ данных при диагностике рака груди
- Wisconsin breast cancer data
- Данные 569 примеров биопсий, каждая из которых характеризуется 30 параметрами
- 31 параметр - тип опухоли



Пример

o Загрузка данных

```
> wbcd <- read.csv("wisc_bc_data.csv",  
stringsAsFactors = FALSE)
```

Пример

◦ Нормализация данных

```
> wbcd_n <- as.data.frame(lapply(wbcd  
[2:31], normalize))
```


Пример

○ Анализ результатов

wbc_d_test_labels	wbc_d_test_pred		Row Total
	Benign	Malignant	
Benign	61 1.000 0.968 0.610	0 0.000 0.000 0.000	61 0.610
Malignant	2 0.051 0.032 0.020	37 0.949 1.000 0.370	39 0.390
Column Total	63 0.630	37 0.370	100

Краткие итоги

○ Анализ данных это:

○ круто

○ интересно

○ доступно

○ выгодно

BIG DATA

$\times 1024$

Zettabyte

$\times 1024$

Exabyte

$\times 1024$

Petabyte

Terabyte



WHERE IS DATA COMING FROM?

Twitter users send out

277,000
tweets

Google processes more than

2 million
search queries

Facebook processes almost

350 GB of data

72 hours

of new video are uploaded to YouTube

EVERY MINUTE...

Individuals and organizations launch

571

new websites

Walmart processes almost

17,000

transactions

More than

100 million

new emails are generated

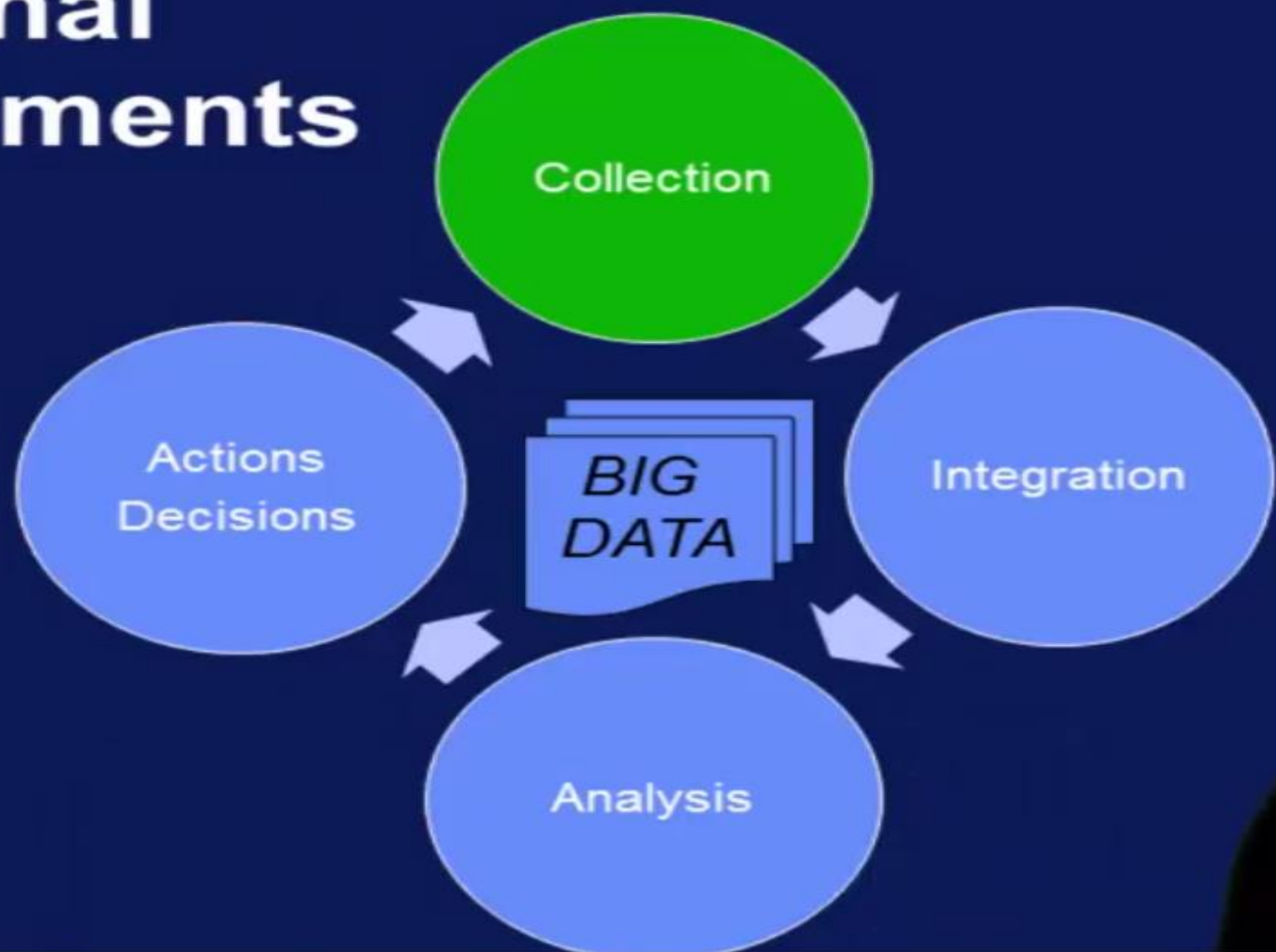
Sprint processes more than

250,000

phone calls



Functional Requirements





Big Data: The Moving Parts

Increasing Age & Maturity



Fast Data

- Hadoop
- Vertica
- MapReduce
- Esper
- kdb
- Greenplum
- ETL
- Netezza
- ECL
- Teradata

Big Analytics

- Hive
- SciPy
- Mahout
- MATLAB
- Revolution R
- SPSS
- AMPL
- SAS

Deep Insight

- unsupervised learning
- social media analytics
- sentiment analysis
- predictive modeling
- BPO
- BI
- network analysis
- visualization
- simulation

- Business Objectives**
- mass customization of services
 - quicker response to market trends
 - identifying real-time cost optimizations
 - faster, more accurate decision making
 - better and more holistic R&D
 - autonomic supply chain management

From <http://blogs.zdnet.com/Hinchcliffe>

the growth of data will be exponential for the foreseeable future



the amount of data stored by the average company today

BIG DATA LANDSCAPE, VERSION 3.0

Exit: Acquisition or IPO

Infrastructure

NoSQL Databases

MapReduce

Cloud

NEWSQL Databases

Column Services

Management / Monitoring

MPP Databases

Graph Databases

Data Transformation

Cloud Sourcing

App Dev.

Analytics

Analysis Platforms

For Business Analysis

Data Sourcing Platforms / Tools

BI Platforms

Unstructured Data

Data Visualization

Machine Learning

Location / People / Events

Big Data Search

Cloud Sourcing

Statistical Computing

Log Analytics

Cloud Sourcing

SMB

Applications

Optimization

Marketing

Finance

Government / Regulation

Education / Learning

Health

Publisher Tools

AI Optimization

Human Capital

Security

Industries

Cross Infrastructure / Analytics

Open Source

Data Sources

HADOOP 1.0



Pig
(data flow)

Hive
(sql)

Other
(cascading)

MapReduce
(cluster resource management
& data processing)

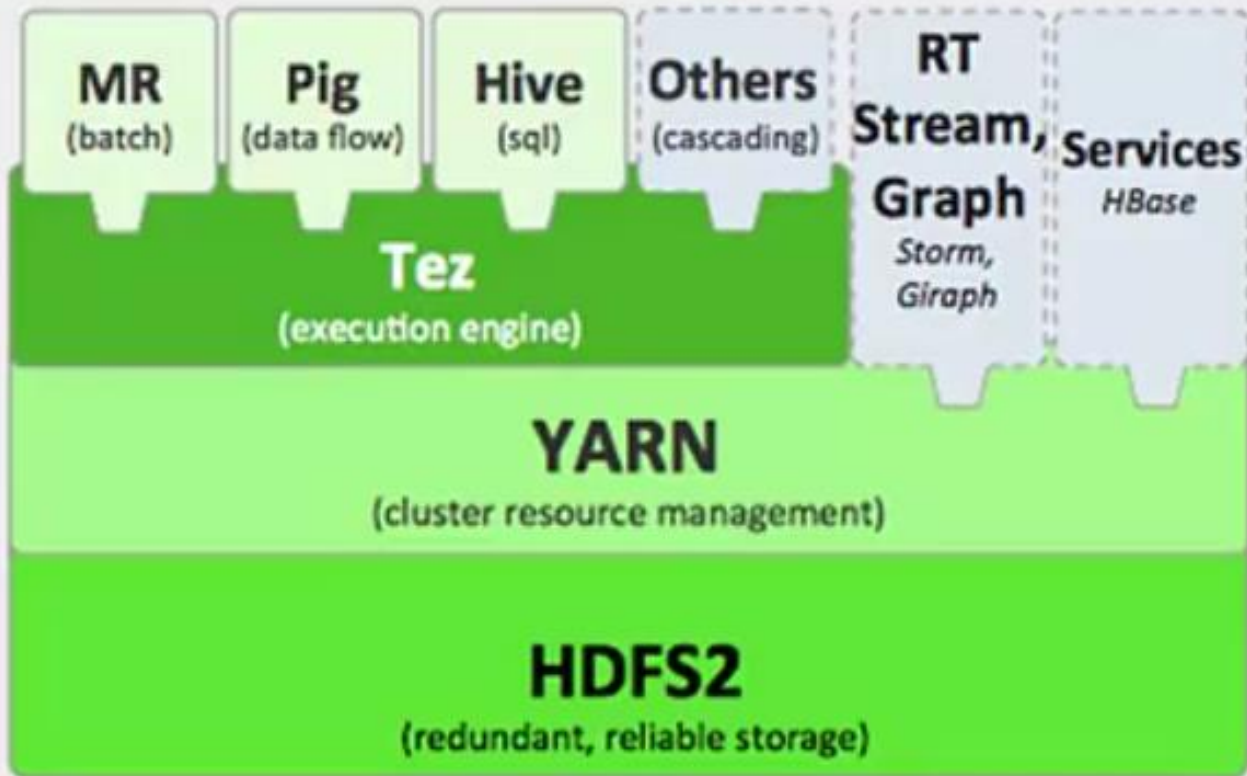
HDFS
(redundant, reliable storage)

Other applications

Processing

Data

HADOOP 2.0



Flume

Log Collector



Sqoop

Data Exchange



Zookeeper

Coordination



HDFS

Hadoop Distributed File System



Oozie

Workflow



Pig

Scripting



Mahout

Machine Learning



R Connectors

Statistics

Hive

SQL Query



Hbase

Columnar Store



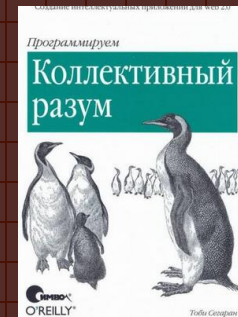
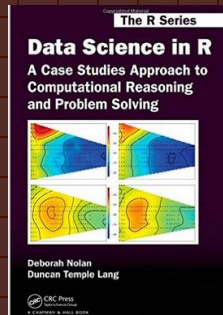
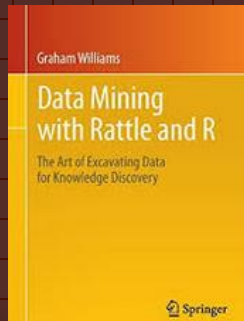
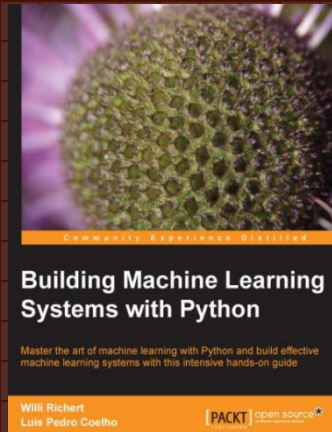
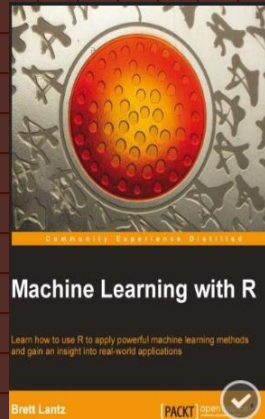
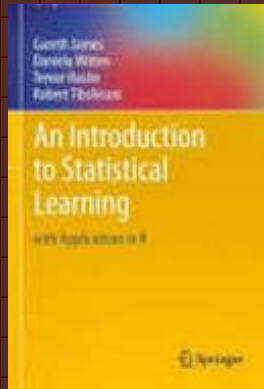
Apache Hadoop Ecosystem

Provisioning, Managing and Monitoring Hadoop Clusters

Ambari



Уголок Библиофила



ОБРАТНАЯ СВЯЗЬ

igkleiner@gmail.com

Ваши вопросы и обратная связь суть лучший
источник мотивации

Благодарности



Благодарности

Образовательный IT-портал
GeekBrains

Благодарности

- Клейнер Надежда
- Бородин Захар
- Гольцман Александр
- Дубинский Игаль
- Гликин Григорий

Ответы на вопросы слушателей