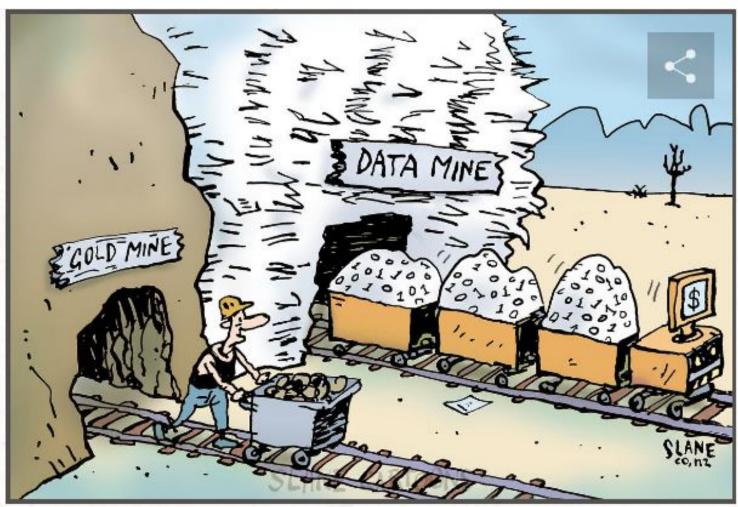




Введение в большие данные

Корчагин Сергей Алексеевич, кандидат физико-математических наук, Заместитель руководителя Департамента анализа данных и машинного обучения SAKorchagin@fa.ru Москва, 2021



https://www.slanecartoon.com/-/galleries/privacy-1/-/medias/f882b7b9-5893-4572-99f1-1c1 7c7ebbe81-data-mine/share





Большие данные (Від Data, биг дата) — это структурированные и неструктурированные данные огромных объемов и разнообразия, а также методы их обработки, которые позволяют распределено анализировать информацию.







Клиффорд Линч, 2008





17

Business Intelligence

Business Intelligence = BI = Бизнес – аналитика (rus) — это набор IT-технологий для сбора, хранения и анализа данных, позволяющих предоставлять пользователям достоверную аналитику в удобном формате, на основе которой можно принимать эффективные решения для управления бизнес-

Входные данные: разнородные отчеты

Отчет первого магазина

	Α	В	C
1	Товар	Количество	Примечания
2	Бум Тул	100 шт.	Договор 100500
3	Кока-кола литровая	84 шт.	Закупить ещё
4	Кока-кола полу-литровая	12 шт.	-
5	Шетки зуб	145 шт.	-
6	Прочее	134 шт.	-

Отчет второго магазина

Продукция	Кол-во	Комментарии
Бумага туалетная		105 Обсудить скидки
³ Coca-cola 1л		14 Нет
4 Coca-cola 0.5л		12 Нет
з Зубные щётки		54 <u>Нет</u>
 Другое 		42 Разобрать

Отчет третьего магазина

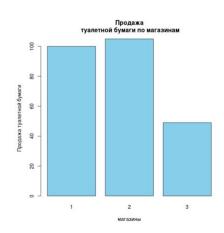
1	Изделие	Продано	Примечание
2	Туалетная бумага	49 рулонов	Закупить двуслойные
3	Кола литровая	14 бутылок	Обсудить опт
4	Кола поллитра	12 бутылок	
5	Щётка зубная	54 упаковок	
6	Остальное	42 штук	

Выходные данные: общая таблица и визуализация

Сводная таблица

	A	В	c	D
1	Товар	Магазин 1	Магазин 2	Магазин 3
2	Туалетная бумага	100	105	49
3	Кола 1л	84	14	14
4	Кола 0.5 л	12	12	12
5	Зубные щётки	145	54	54
6	Прочее	134	42	42

Аналитика по сводным данным





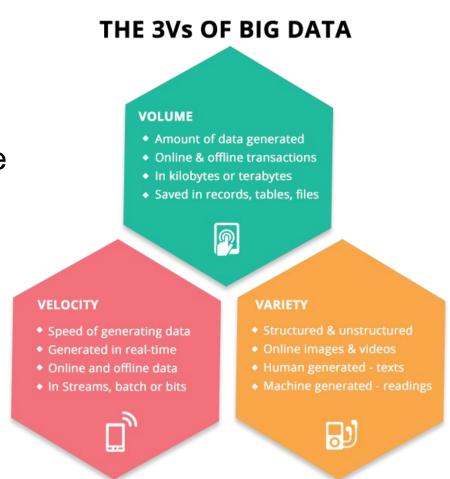
Традиционная аналитика	Big data аналитика
Постепенный анализ небольших пакетов данных	Обработка сразу всего массива доступных данных
Редакция и сортировка данных перед обработкой	Данные обрабатываются в их исходном виде
Старт с гипотезы и ее тестирования относительно данных	Поиск корреляций по всем данным до получения искомой информации
Данные собираются, обрабатываются, хранятся и лишь затем анализируются	Анализ и обработка больших данных в реальном времени, по мере поступления





Правило VVV — три признака или свойства, которыми большие данные должны обладать:

- Volume объем
- 2. Velocity скорость
- Variety многообразие



данных

💸 Функции и задачи больших д

Функция	Задача
Big Data — собственно массивы необработанных данных	Хранение и управление большими объемами постоянно обновляющейся информации
Data mining — процесс обработки и структуризации данных, этап аналитики для выявления закономерностей	Структурирование разнообразных сведений, поиск скрытых и неочевидных связей для приведения к единому знаменателю
Machine learning — процесс машинного обучения на основе обнаруженных связей в процессе анализа	Аналитика и прогнозирование на основе обработанной и структурированной информации



- **Интернет** соцсети, блоги, СМИ, форумы, сайты, интернет вещей (<u>loT</u>).
- Корпоративные данные транзакционная деловая информация, архивы, базы данных.
- Показания устройств — датчиков, приборов, а также метеорологические данные, данные СОТОВОЙ СВЯЗИ И Т.Д. -





- Горизонтальная масштабируемость
- Отказоустойчивость
- Локальность данных



🧩 Что можно сделать с помощью больших данных?

Предсказать победителя

Оскара!





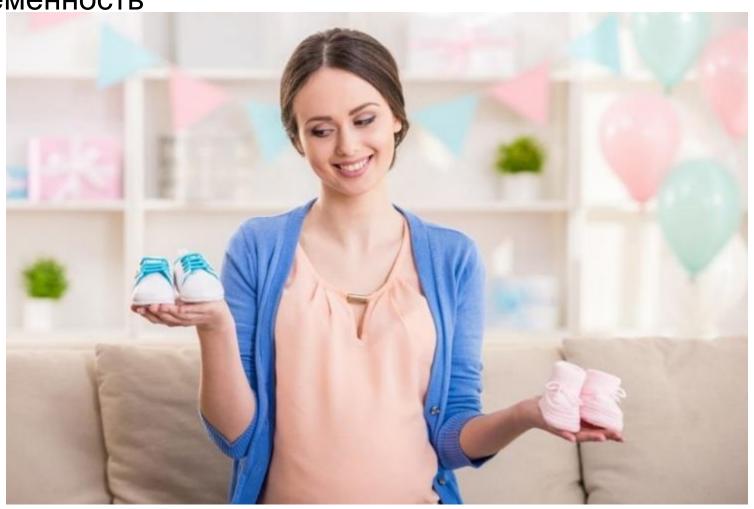


Найти военную базу **HATO**





Диагностировать беременность



🏞 Области применения технологий больших данных

Анализ данных опросов:

- Мониторинг общественного мнения и анализ социальноэкономической ситуации
- Определение проблем, формирующих кризисную ситуацию
- Анализ реакции населения на внедрение различных федеральных и региональных программ
- Анализ экономического положения и уровня жизни населения



🧩 Области применения технологий больших данных



Предвыборные исследования

- Анализ эффективности политической рекламы
- Анализ средств массовой информации
- Выявление наиболее эффективных средств влияния на мнения различных групп избирателей
- Диагностика предвыборной ситуации
- Анализ основных проблем избирателей



🧩 Области применения технологий больших данных

Общественная безопасность

- Анализ преступности
- Отслеживание уровня рецидивизма



🌣 Области применения технологий больших данных

Образование

- Планирование школьных округов
- Отслеживание успеваемости учащихся, выявление факторов способствующих повышению успеваемости
- Администрирование контроль за уровнем выполнения обязательных простосим и тостор



🧩 Области применения технологий больших данных



- Анализ рынка труда понимание состава и структуры рабочей силы
- Анализ заявлений о приеме на работу разработка профилей претендентов.



🧩 Области применения технологий больших данных

Анализ прибыли

- Оценка соответствия размеров уплаченных налогов и имущества
- Анализ мошенничеств





Области применения технологий больших данных

Здравоохранение

- Отслеживание болезней и создание отчетов о случаях заболеваний
- Эпидемиология выявление причин заболеваний и территории их

распространения, а также контроль заболеваемости

- Медицинская помощь - определение профилей тех, кому часто требуется медицинская помощь

 Профилактика медицинского вк



🧩 Области применения технологий больших данных

Окружающая среда

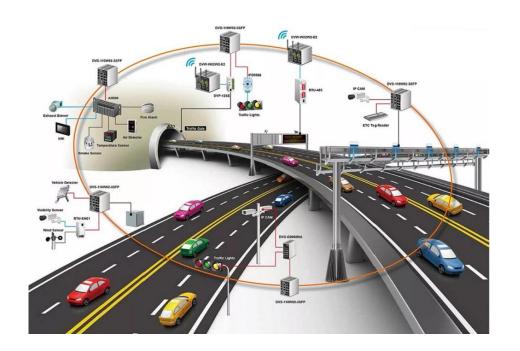
- Анализ экосистем выяснение факторов, влияющих на здоровье экосистемы
- Оценка качества воды/воздуха осуществление контроля за





Транспорт

- Планирование наиболее эффективных маршрутов для лучшей организации транспортных и пассажирских потоков
- Создание отчетов о дорожных происшествиях для выяснения факторов, влияющих на происшествия
- Моделирование программ поддержания надлежащего состояния дорожного покрытия, прогнозирование возможного ремонта дорог.



🧦 Области применения технологий больших данных



Стратегическое планирование

- Анализ удовлетворенности клиентов и изучения изменений потребностей общественности
- Оценка программ понимание факторов успешной реализации программы
- Профилирование населения более эффективное направление действия программы на определенные слои населения
- Анализ затрат выявления наиболее эффективных программ
- Анализ результатов выполнения программ



Технологии больших данных в маркетинге



- Создание точных портретов целевых потребителей.
- Предсказание реакции потребителей на маркетинговые сообщения.
- Максимальная персонализация рекламных сообщений.
- Увеличение кросс-продаж, повторных продаж, ремаркетинга.
- Поиск и определение причин популярности востребованных товаров и продуктов.
- Совершенствование продуктов и услуг, повышение лояльности клиентов.
- Повышение качества обслуживания.
- . Предупреждение мо
- . Снижение издержек



циками и клиентами.

Большие данные в бизнесе

Поставщики инфраструктуры — решают задачи хранения и предобработки данных.







Датамайнеры — разработчики алгоритмов, которые помогают заказчикам

RNHS TOWARD LIQUIN TO ADOPT HUR.



Yandex Data Factory





Системные интеграторы — компании, которые внедояют системы анализа

больших данны





Потребители — компании, которые покупают программно-аппаратные комплексы и заказывают алгоритмы — компании, которые покупают программно-аппаратные комплексы и заказывают алгоритмы — компании, которые покупают программно-аппаратные комплексы и заказывают алгоритмы — компании, которые покупают программно-аппаратные комплексы и заказывают алгоритмы — компании, которые покупают программно-аппаратные комплексы и заказывают алгоритмы — компании, которые покупают программно-аппаратные комплексы и заказывают алгоритмы — комплексы и заказывают алгоритмы и заказывают и заказывают алгоритмы и заказ



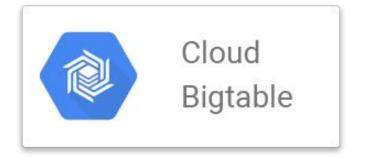






Yandex Data Factory







Рамблер/

💸 Выгоды использования Big Data в бизнесе

- Упрощается планирование.
- Увеличивается скорость запуска новых проектов.
- Повышаются шансы проекта на востребованность.
- Можно оценить степень удовлетворенности пользователей.
- Проще найти и привлечь целевую аудиторию.
- Ускоряется взаимодействие с клиентами и контрагентами.
- Оптимизируются интеграции в цепи поставок.
- Повышается качество клиентского сервиса, скорость взаимодействия.
- Повышается лояльность текущих клиентов.

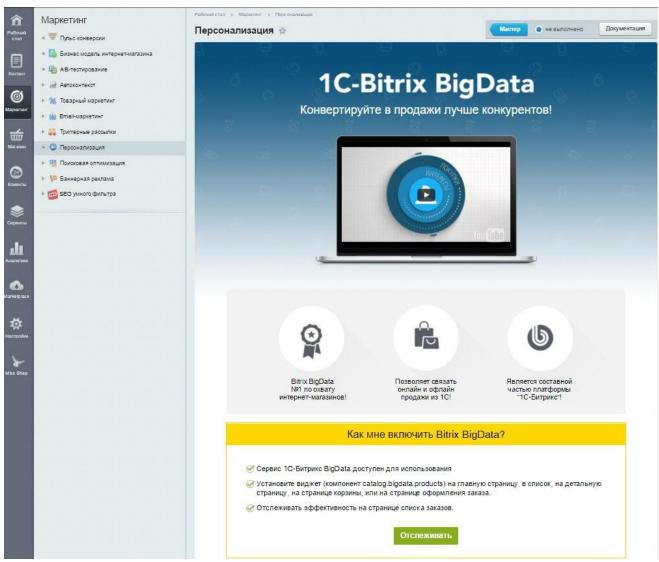


🔀 Драйверы и ограничители Big Data в России

Драйверы	Ограничители
Драиверы	Ограничители
Высокий спрос на Big Data для повышения конкурентоспособности с помощью возможностей технологий	Необходимость обеспечивать безопасность и конфиденциальность данных
Развитие методов обработки медиафайлов на мировом уровне	Нехватка квалифицированных кадров
Реализация отраслевого плана по импортозамещению программного обеспечения	В большинстве российских компаний объем накопленных информационных ресурсов не достигает уровня Big Data
Тренд на использование услуг российских провайдеров и системных интеграторов	Новые технологии сложно внедрять в устоявшиеся информационные системы компаний
Создание технопарков, которые способствуют развитию информационных технологий	Высокая стоимость технологий
Государственная программа по внедрению грид-систем — виртуальных суперкомпьютеров, которые распространяются по кластерам и связываются сетью	Заморозка инвестиционных проектов в России и отток зарубежного капитала
Перенос на территорию России серверов,	Рост цен на импортную продукцию

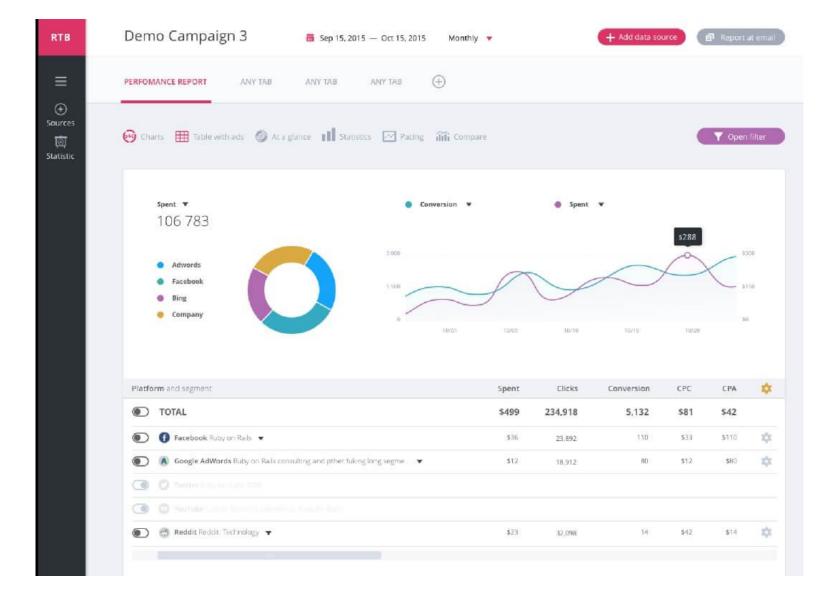
🎎 Сервисы Big Data

«1С-Битрикс BigData»





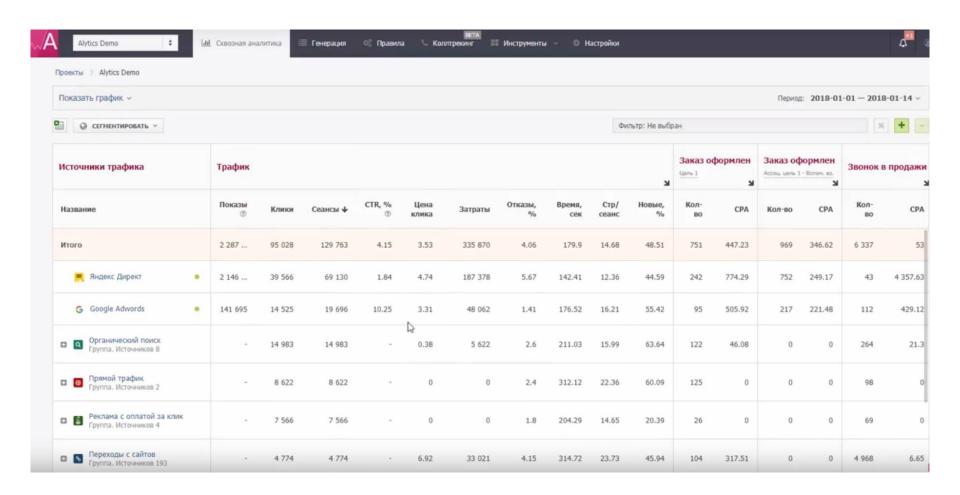
RTB Media







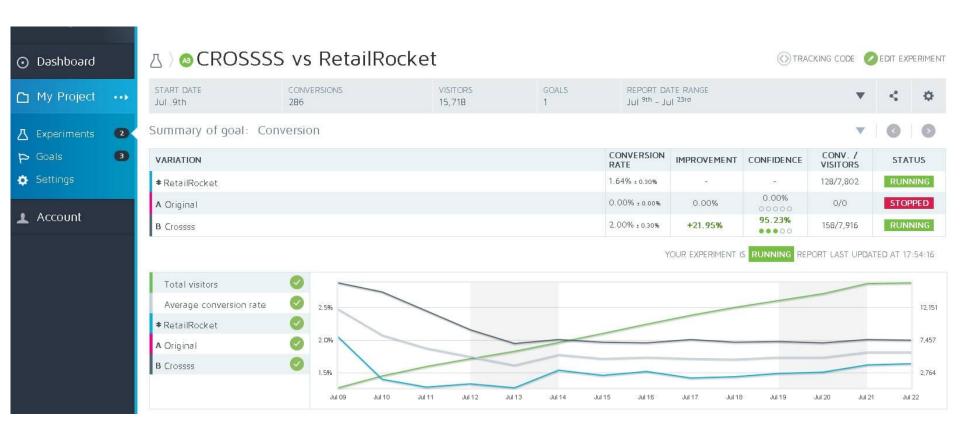
Alytics







Crossss







- NoSQL;
- MapReduce;
- Hadoop;
- R;
- Python;
- Аппаратные решения.













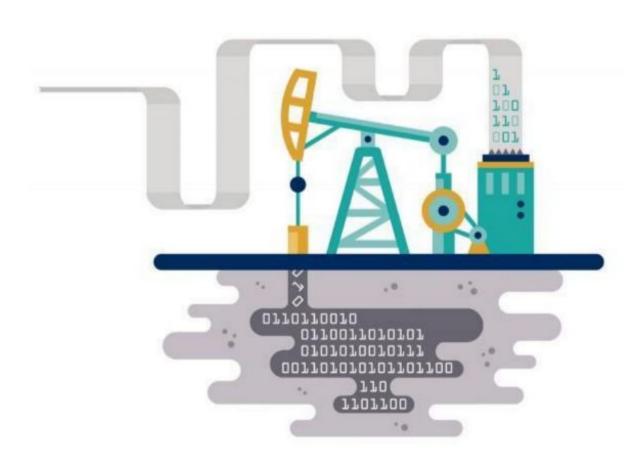
Основными методами и техниками анализа, применимыми к Большим данным, являются следующие:

- •Методы класса Data Mining
- •Краудсорсинг
- •Смешение и интеграция данных
- •Машинное обучение
- •Визуализация аналитических данных



Data Mining







Data Mining – это сочетание широкого математического инструментария (от классического статистического анализа до новых кибернетических методов) и последних достижений в сфере информационных технологий

Data Mining (добыча данных, интеллектуальный анализ данных, глубинный анализ данных) — собирательное название, используемое для обозначения совокупности методов обнаружения в данных ранее неизвестных, нетривиальных, практически полезных и доступных интерпретации знаний, необходимых для принятия решений в различных сферах человеческой деятельности.



Термин введён Григорием Пятецким-Шапиро в 1989 году.



Data Mining - мультидисциплинарная область, возникшая и развивающаяся на базе таких наук как прикладная статистика, распознавание образов, искусственный интеллект, теория баз данных



🧱 Модели представления знаний Data Mining



1	• ассоциативные правила	
2	• деревья решений	
3	• кластеры	
3	• математические функции	



- искусственные нейронные сети
- деревья решений, символьные правила
- методы ближайшего соседа и k-ближайшего соседа
- метод опорных векторов
- байесовские сети
- линейная регрессия
- корреляционно-регрессионный анализ
- иерархические методы кластерного анализа
- неиерархические методы кластерного анализа, в том числе алгоритмы k-средних и k-медианы
- эволюционное программирование и генетические алгоритмы
- метод ограниченного перебора
- эволюционное программирование и генетические алгоритмы
- разнообразные методы визуализации данных и множество других методов.



- точность
- масштабируемость
- интерпретируемость
- проверяемость
- трудоемкость
- гибкость
- быстрота
- популярность





АЛГОРИТМ	точность	МАСШТАБИРУЕМОСТЬ	ИНТЕРПРЕТИРУЕМОСТЬ	пригодность кисп.
линейная регрессия	нейтральная	EMCOKSE	высокая / нейтральная	EPICORSA
нейронные сети	высокая	HURRAN	HMSKRR	низкая
методы визуапизации	высокая	очень низкая	высокая	высокая
деревья решений	HHIKAR	BAICOKRE	BPICOKBE	высокая / нейтральная
нейронные сети	высокая	нейтральная	HERRAR	высокая / нейтральная
k-ближайшего соседа	RESERVE	очень низкая	высокая / нейтральная	нейтральная

АЛГОРИТМ	трудоемкость	РАЗНОСТОРОННОСТЬ	БЫСТРОТА	популярность
линейная регрессия	нейтральная	нейтральная	высокая	низкая
нейронные сети	нейтральная	низкая	очень низкая	HHSKRR
методы визуализации	очень высокая	никая	чрезвычайно низкая	высокая / нейтральная
деревья решений	высокая	BPICOKBE	высокая / нейтральная	высокая / нейтральная
нейронные сети	низкая / нейтральная	нейтральная	низкая / нейтральная	нейтральная
k-б.1иж айшего соседа	нейтральная низкая	никая	EPICONSE	низкая

Таблица 1 Сравнительная характеристика методов Data Mining

💸 Классификация методов Data Mining





- кластерный анализ
- метод ближайшего соседа
- метод к-ближайшего соседа
- рассуждение по аналогии

Методы группы:

- логические методы
 - нечеткие запросы и анализы
 - символьные правила
 - деревья решений
 - генетические алгоритмы.
- методы визуализации
- методы кросс-табуляции
- методы, основанные на уравнениях

🧩 Подход к обучению математических моделей Data Mining



Статистические методы Data Mining



- дескриптивный анализ и описание исходных данных
- анализ связей (корреляционный и регрессионный анализ, факторный анализ, дисперсионный анализ)
- многомерный статистический анализ (компонентный анализ, дискриминантный анализ, многомерный регрессионный анализ, канонические корреляции и др.)
- анализ временных рядов (динамические модели и прогнозирование)

Кибернетические методы Data Mining



- искусственные нейронные сети (распознавание, кластеризация, прогноз);
- эволюционное программирование (в т.ч. алгоритмы метода группового учета аргументов);
- генетические алгоритмы (оптимизация);
- ассоциативная память (поиск аналогов, прототипов);
- нечеткая логика:
- деревья решений;
- системы обработки экспертных знаний.





Описательные методы



- алгоритм к-средних
- к-медианы
- иерархические методы кластерного анализа
- самоорганизующиеся карты Кохонена
- методы кросс-табличной визуализации

Прогнозирующие методы



- нейронные сети
- деревья решений,
- линейная регрессия
- метод ближайшего соседа
- метод опорных векторов



Aberdeen Group: " Data Mining - технология добычи полезной информации из баз данных. Однако в связи с существенными различиями между инструментами, опытом и финансовым состоянием поставщиков продуктов, предприятиям необходимо тщательно оценивать предполагаемых разработчиков Data Mining и партнеров.

ABERDEENGROUP



Herb Edelstein: «Недавнее исследование компании Two Crows показало, что Data Mining находится все еще на ранней стадии развития. Многие организации интересуются этой технологией, но лишь некоторые активно внедряют такие проекты. Удалось выяснить еще один важный момент: процесс реализации Data Mining на практике оказывается более сложным, чем ожидается».



- 1. Data Mining не может заменить аналитика
- 2. Сложность разработки и эксплуатации приложения Data Mining
- 3. Квалификация пользователя
- 4. Извлечение полезных сведений невозможно без хорошего понимания сути данных
- 5. Сложность подготовки данных
- 6. Большой процент ложных, недостоверных или бессмысленных результатов
- 7. Высокая стоимость
- 8. Наличие достаточного количества репрезентативных данных



Перспективы технологии Data Mining

- выделение типов предметных областей с соответствующими им эвристиками, формализация которых облегчит решение соответствующих задач Data Mining, относящихся к этим областям;
- создание формальных языков и логических средств, с помощью которых будут формализованы рассуждения и автоматизация которых станет инструментом решения задач Data Mining в конкретных предметных областях;
- создание методов Data Mining, способных не только извлекать из данных закономерности, но и формировать некие теории, опирающиеся на эмпирические данные ;
- преодоление существенного отставания возможностей инструментальных средств Data Mining от теоретических достижений в этой области





Области, где применения технологии Data Mining, скорее всего, будут успешными, имеют такие особенности:

- требуют решений, основанных на знаниях;
- имеют изменяющуюся окружающую среду;
- имеют доступные, достаточные и значимые данные ;
- обеспечивают высокие дивиденды от правильных решений.





Международная конференция по Knowledge Discovery Data Mining (International Conferences on Knowledge Discovery and Data Mining).

Среди наиболее известных WWW-источников - сайт www.kdnuggets.com, который ведет один из основателей Data Mining Григорий Пиатецкий-Шапиро.

Периодические издания по Data Mining: Data Mining and Knowledge Discovery, KDD Explorations, ACM-TODS, IEEE-TKDE, JIIS, J. ACM, Machine Learning, Artificial Intelligence.

Материалы конференций: ACM-SIGKDD, IEEE-ICDM, SIAM-DM, PKDD, PAKDD, Machine learning (ICML), AAAI, IJCAI, COLT (Learning Theory).

Краудсорсинг







Краудсорсинг — привлечение к решению какой-либо проблемы большой группы людей



В 2003 году **Луис фон Ах (Luis von Ahn)** вместе со своими коллегами впервые предложил понятие "человеческих вычислений

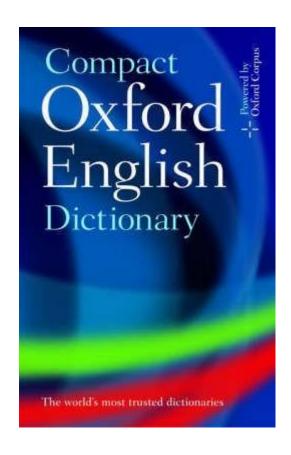


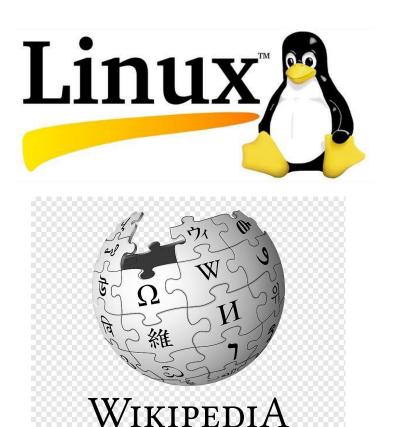
в 2006 году термин "Краудсорсинг" (crowdsourcing) был сформулирован редактором журнала "Wired" **Джеффом** Xay (Jeff Howe)





Краудсорсинг - это мобилизация ресурсов людей посредством информационных технологий с целью решения задач, стоящих перед бизнесом, государством и обществом в целом.





The Free Encyclopedia



- По сфере жизни (бизнес, социальный, политический)
- П. По типу решаемых задач (создание продукта (контента), голосование, поиск решения, поиск людей, сбор информации, сбор мнений, тестирование, служба поддержки, сбор средств Краудфандинг).

^{*}http://crowdsourcing.ru/article/what is the crowdsourcing







По сфере жизни:

- Бизнес
- 2) Социальный или общественный
- 3) Политический или государственный





Классификация Краудсорсинга

По типу решаемых задач:

- 1) Создание продукта (контента) <u>99designs</u>, <u>TopCoder</u>, <u>Witmart</u>, <u>Tongal</u>, <u>Audiodraft</u>
- 2) Голосование
- 3) Поиск решения Kaggle, CrowdFlower, InnoCentive, Academy of Ideas, Wazoku
- 4) Поиск людей, например Liza Alert
- 5) Сбор информации, например Zooniverse
- 6) Сбор мнений, например <u>Chaordix</u>, <u>Innopinion</u> и <u>AnswerTap</u>
- 7) Тестирование, например <u>uTest</u>,
- 8) Сбор средств Краудфандинг



- 1) Большой охват
- 2) Вовлечение пользователей
- 3) Разнообразие выбора
- 4) Единственно возможный вариант
- 5) Фиксированные сроки
- 6) Экономия финансовых ресурсов







CITYCELEBRITY CROWDSOURCING PLATFORM







СПАСИБО ЗА ВНИМАНИЕ!

Сергей Алексеевич Корчагин SAKorchagin@fa.ru 2021