

Исторические исследования в контексте науки о данных:  
информационные ресурсы, аналитические методы и  
цифровые технологии

# О методологии машинного обучения

*Воронцов Константин Вячеславович*

д.ф.-м.н., профессор РАН,  
руководитель лаборатории Машинного интеллекта МФТИ

[k.v.vorontsov@phystech.edu](mailto:k.v.vorontsov@phystech.edu)

# Докладчик: *Воронцов Константин Вячеславович*

http://www.MachineLearning.ru/wiki?title=User:Vokov

Участник:Vokov

## Участник:Vokov

**Воронцов Константин Вячеславович**

профессор РАН, д.ф.-м.н.,  
руководитель [лаборатории машинного интеллекта МФТИ](#),  
проф. каф. «Интеллектуальные системы» [ФУПМ МФТИ](#),  
с.н.с. отдела «Интеллектуальные системы» [Вычислительного центра ФИЦ ИУ РАН](#),  
доц. каф. «Математические методы прогнозирования» [ВМК МГУ](#),  
преподаватель [Школы анализа данных Яндекса](#),  
зам. директора по науке [ЗАО «Форексис»](#), [www.forecsys.ru](#),  
один из идеологов и [Администраторов](#) ресурса [MachineLearning.RU](#),  
прочие подробности — на подстранице [Curriculum vitae](#).

- [Профиль ORCID = 0000-0002-4244-4270](#)
- [Профиль SCOPUS ID = 6507982932](#)
- [Профиль WoS ResearcherID = G-7857-2014](#)
- [Профиль Google Scholar](#)
- [Профиль DBLP](#)
- [Профиль РИНЦ ID = 15081](#)
- [Профиль в системе ИСТИНА](#)
- [Профиль MathNet.ru](#)

[Мне можно написать письмо.](#)

### Содержание [убрать]

- 1 Учебные материалы
  - 1.1 Курсы лекций
  - 1.2 Рекомендации для студентов и аспирантов
- 2 Интервью
  - 2.1 Российский радиоуниверситет, Радио России
  - 2.2 Газеты, журналы, электронные СМИ
  - 2.3 Видеоинтервью
- 3 Доклады на конференциях и семинарах
- 4 Научные интересы
  - 4.1 Анализ текстов и информационный поиск
  - 4.2 Диагностика заболеваний по ЭКГ
  - 4.3 Теория обобщающей способности
  - 4.4 Комбинаторная (перестановочная) статистика
  - 4.5 Прогнозирование объёмов продаж
  - 4.6 Другие проекты и семинары
- 5 Публикации
- 6 Софт
- 7 Аспиранты и студенты
  - 7.1 Бакалаврские диссертации
  - 7.2 Магистерские диссертации

# О методологии машинного обучения

## 1. Задачи машинного обучения

- Бум искусственного интеллекта и нейронных сетей
- Постановки задач и терминология машинного обучения
- Примеры задач машинного обучения

## 2. Методология машинного обучения

- Нейронные сети и глубокое обучение
- Обучение как задача оптимизации
- Типология и методология машинного обучения

## 3. Проблемы и перспективы применения

- Особенности практического применения DS/AI/ML
- Необходимые условия применения DS/AI/ML
- Мифы об искусственном интеллекте

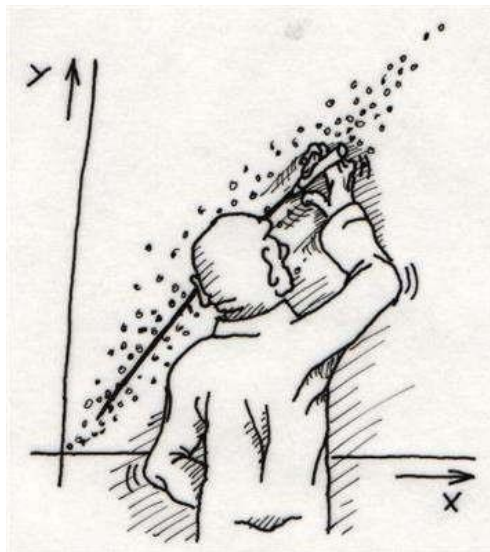
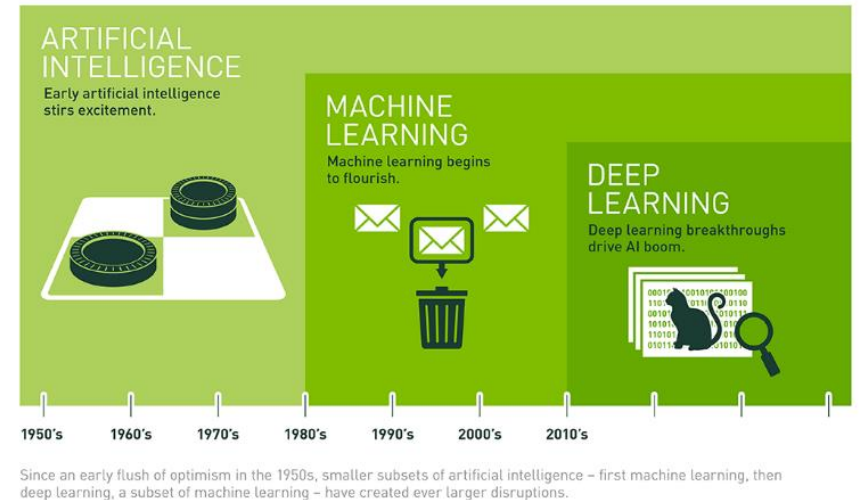
«Четвёртая технологическая революция строится на вездесущем и мобильном Интернете, *искусственном интеллекте* и *машинном обучении*» (2016)

Клаус Мартин Шваб,  
президент Всемирного  
экономического форума



# Машинное обучение (Machine Learning, ML)

- одна из ключевых информационных технологий будущего
- наиболее успешное направление ИИ, вытеснившее экспертные системы и инженерию знаний



- проведение функции через заданные точки в сложно устроенных пространствах
- математическое моделирование в условиях, когда знаний мало, данных много
- тысячи различных методов и алгоритмов
- около 100 000 научных публикаций в год



# Задача машинного обучения с учителем

## Этап №1 – обучение (train)

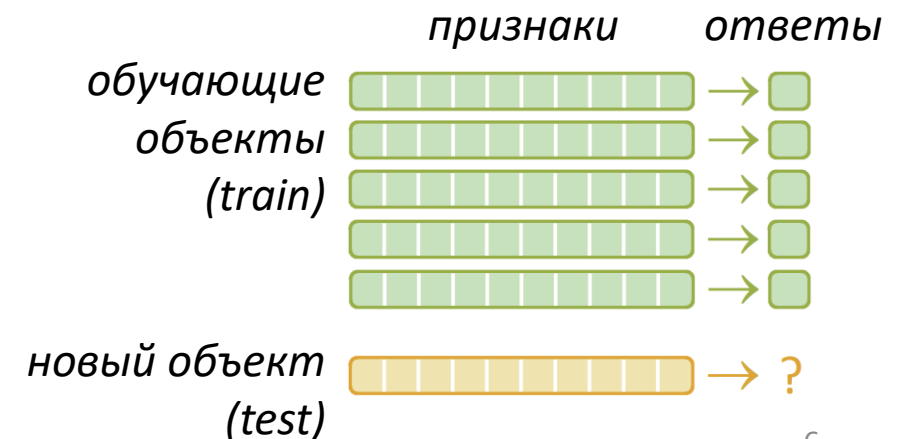
- **На входе:**  
*данные* – выборка пар «*объект* → *ответ*»,  
каждый объект описывается набором *признаков*
- **На выходе:**  
*модель*, предсказывающая ответ по объекту

Задача поставлена,  
если у неё есть «**ДНК**»:

- **Дано**
- **Найти**
- **Критерий**

## Этап №2 – применение (test)

- **На входе:**  
*данные* – **новый объект**
- **На выходе:**  
*предсказание* ответа на **новом объекте**



# Примеры задач машинного обучения

- **Медицинская диагностика:**

**объект** – данные о пациенте на текущий момент

**ответ** – диагноз / лечение / риск исхода



- **Поиск месторождений полезных ископаемых:**

**объект** – данные о геологии района

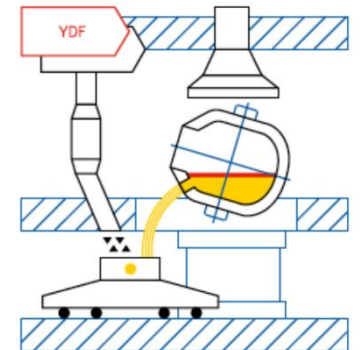
**ответ** – есть/нет месторождение



- **Управление технологическими процессами:**

**объект** – данные о сырье и управляющих параметрах

**ответ** – количество/качество полезного продукта



# Примеры задач ML в бизнесе

- **Кредитный скоринг:**

**объект** – данные о заёмщике

**ответ** – решение по кредиту & вероятность дефолта



- **Предсказание оттока клиентов:**

**объект** – данные о клиенте на момент времени  $t$

**ответ** – уйдёт ли клиент к моменту времени  $t + \Delta$



- **Прогнозирование объёмов продаж:**

**объект** – данные о продажах на момент времени  $t$

**ответ** – объём спроса в интервале от  $t$  до  $t + \Delta$





# Примеры задач ML в интернет-сервисах

- **Информационный поиск в Интернете:**

**объект** – данные о паре «запрос и документ»

**ответ** – оценка релевантности документа запросу



- **Продажа рекламы в Интернете:**

**объект** – данные о тройке «пользователь, страница, баннер»

**ответ** – оценка вероятности клика (CTR, Click-Through Rate)

- **Рекомендательные системы в Интернете / TV:**

**объект** – данные о паре «пользователь, товар / фильм»

**ответ** – оценка вероятности покупки / просмотра



# Примеры задач ML в LegalTech

- **Поиск схожей судебной практики:**

**объект** – текст иска, акта или обращения заявителя

**ответ** – ранжированный список схожих дел



- **Рекомендательный сервис:**

**объект** – пара «описание дела, профиль юриста/фирмы»

**ответ** – ранжированный список консультантов



- **Предсказание судебного решения:**

**объект** – описание дела, документы по делу

**ответ** – вероятность выиграть дело



# Примеры задач с не векторными данными

- **Статистический машинный перевод:**

**объект** – предложение на естественном языке

**ответ** – его перевод на другой язык

*Прогресс в этих  
областях связан с  
«большими данными»  
(англ. «Big Data»)*

- **Перевод речи в текст:**

**объект** – аудиозапись речи человека

**ответ** – текстовая запись речи

*...очень важное уточнение:*

***с аккуратными***

*большими данными*

- **Компьютерное зрение:**

**объект** – динамика сцены в видеопоследовательности

**ответ** – решение (объехать, остановиться, игнорировать)

# О методологии машинного обучения

## 1. Задачи машинного обучения

- Бум искусственного интеллекта и нейронных сетей
- Постановки задач и терминология машинного обучения
- Примеры задач машинного обучения

## 2. Методология машинного обучения

- Нейронные сети и глубокое обучение
- Обучение как задача оптимизации
- Типология и методология машинного обучения

## 3. Проблемы и перспективы применения

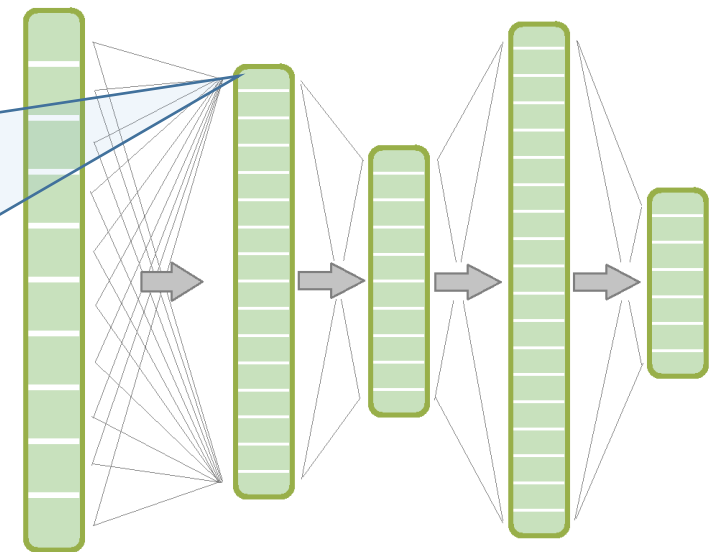
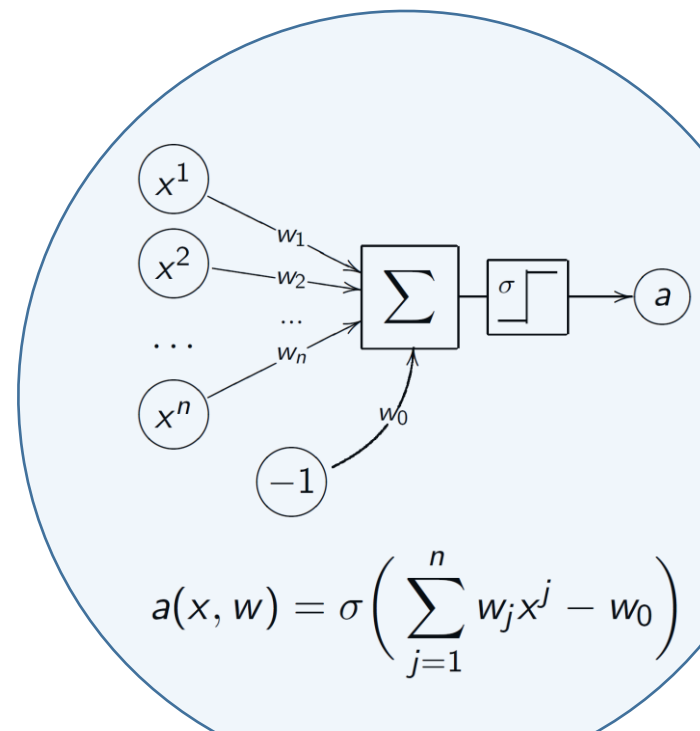
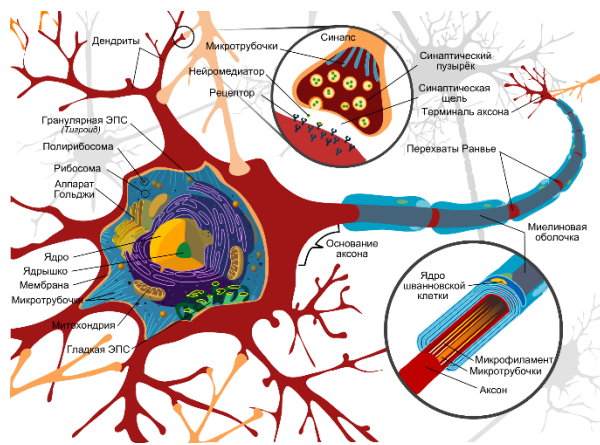
- Особенности практического применения DS/AI/ML
- Необходимые условия применения DS/AI/ML
- Мифы об искусственном интеллекте

# Искусственные нейронные сети

На каждом слое сети вектор объекта преобразуется в новый вектор

Эти преобразования обучаемые, их параметры входят в  $w$

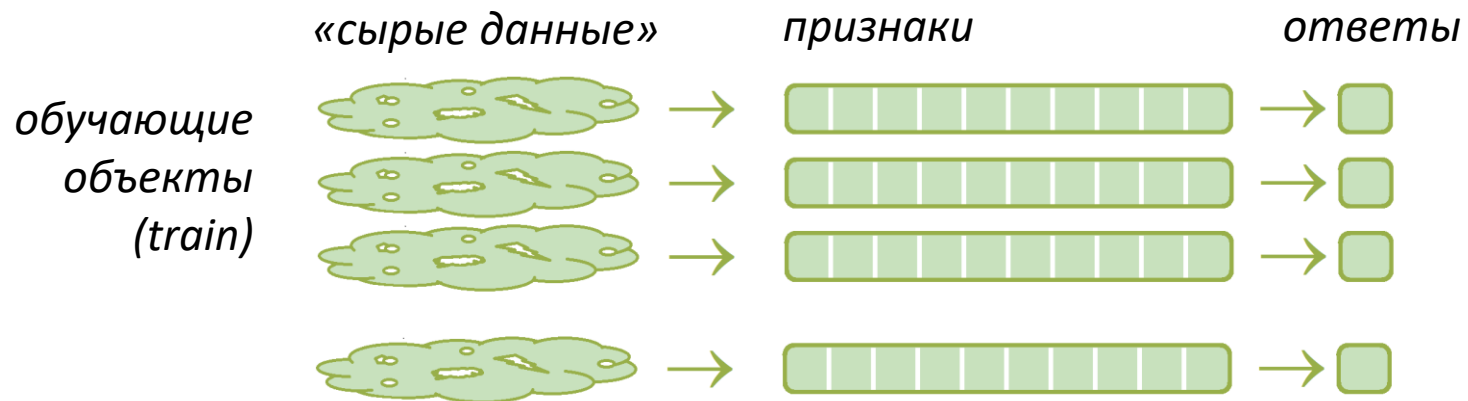
Каждое преобразование (нейрон) – взвешенная сумма признаков



# Глубокие нейронные сети

**Вход:** сложно структурированные «сырые» данные объектов

**Выход:** векторные представления объектов и ответы

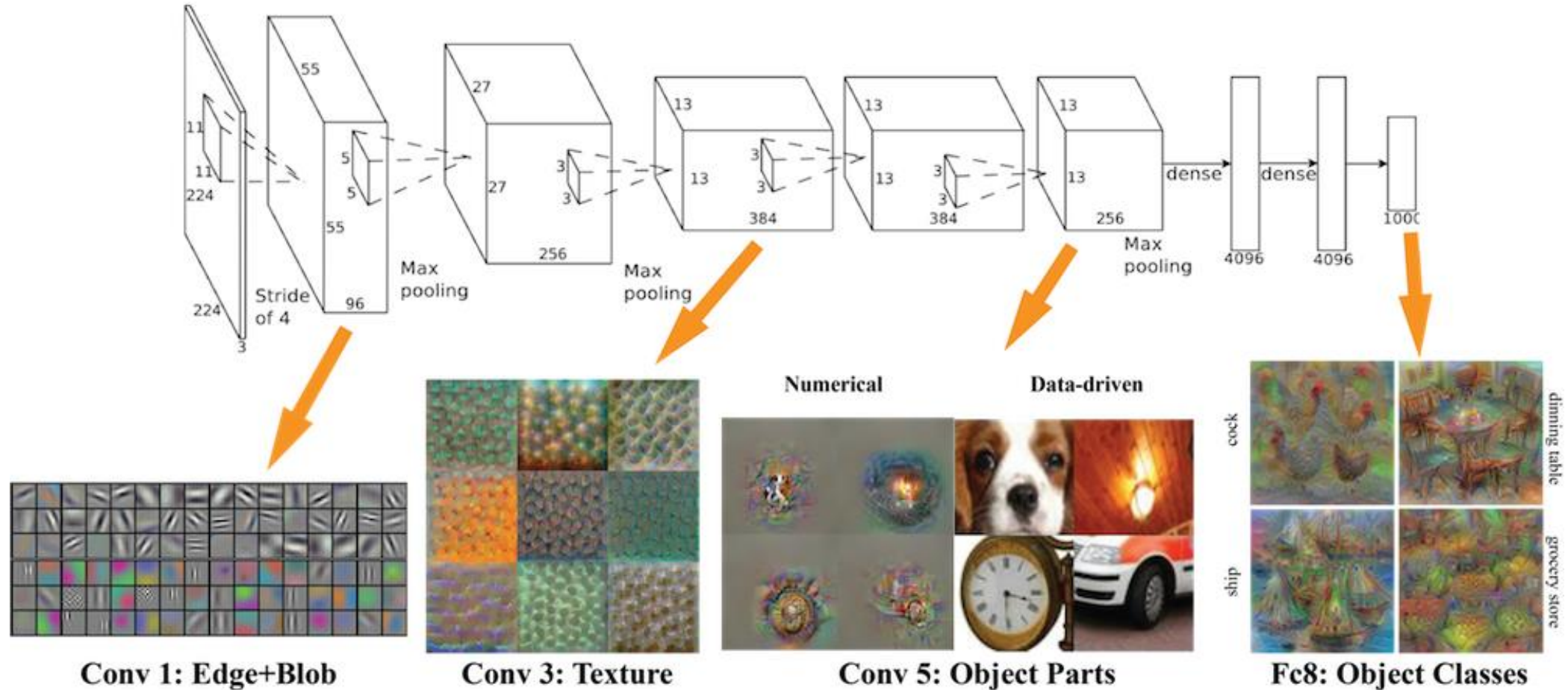


*Deep Learning – это всего лишь обучаемая векторизация сложных объектов*

**Примеры** сложно структурированных объектов:  
тексты, изображения, видео, временные ряды, транзакции, графы, ...



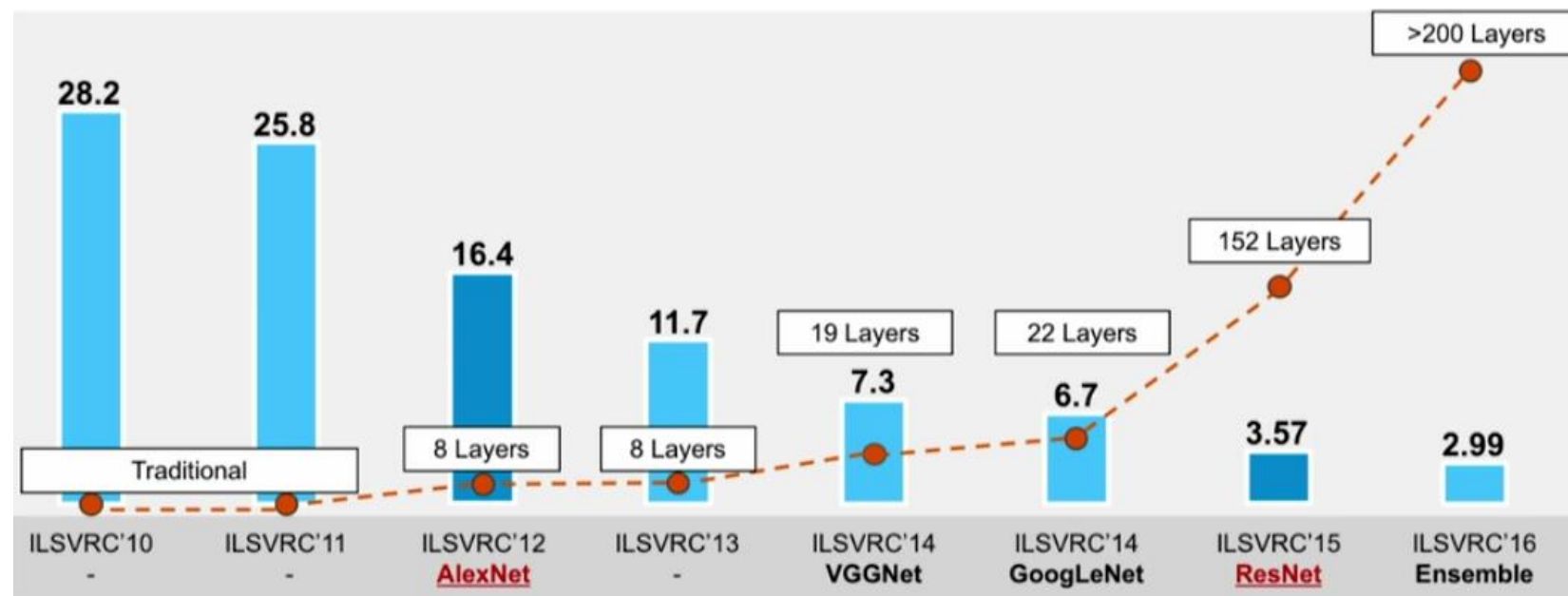
# Глубокие свёрточные нейронные сети для классификации изображений



# Роль больших данных

**ImageNet:** открытая выборка 14М изображений, 20К категорий

IMAGENET

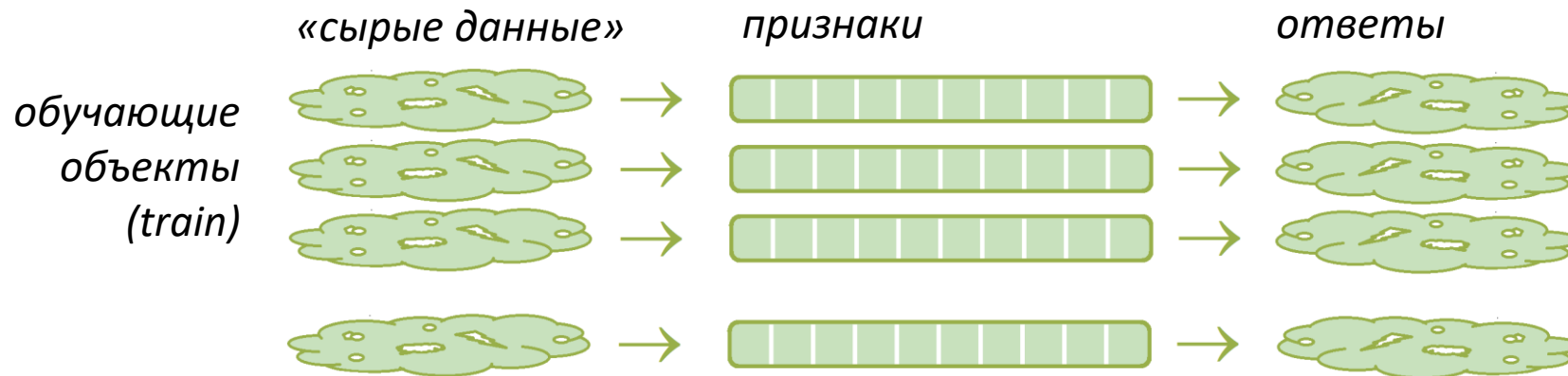


Старт в 2009 г. Человеческий уровень ошибок 5% пройден в 2015 г.

# Нейронные сети для синтеза объектов

**Вход:** сложно структурированные объекты

**Выход:** сложно структурированные ответы



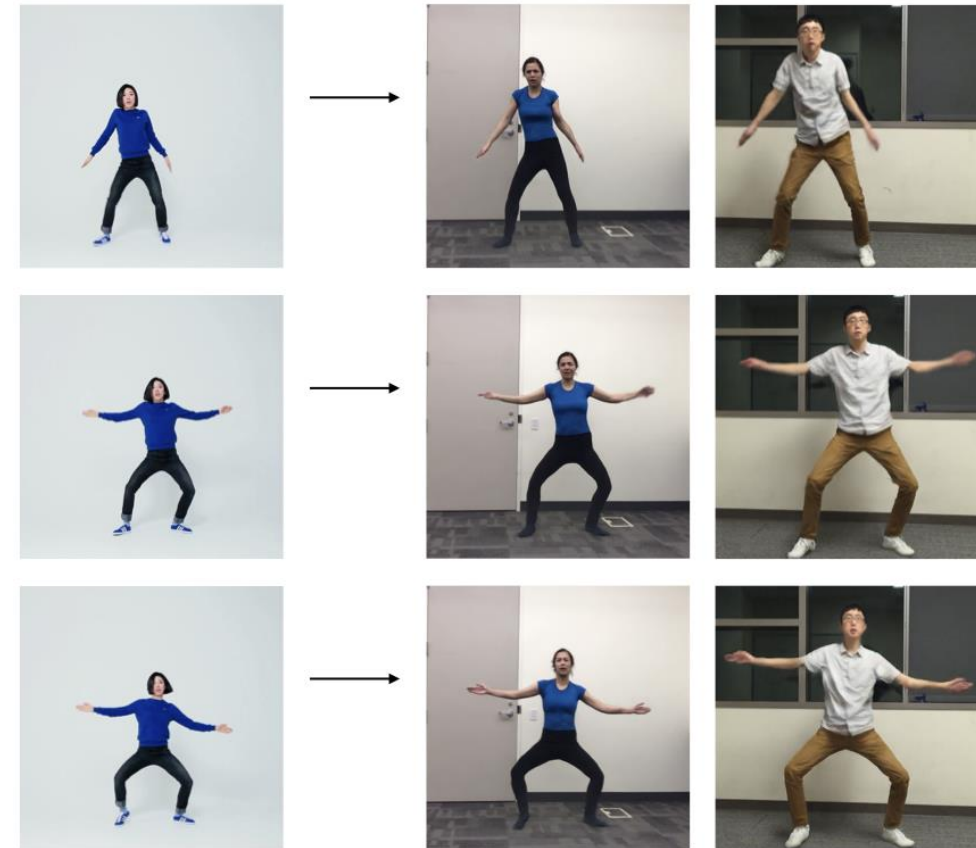
**Примеры:** синтез изображений, перенос стиля, машинный перевод, суммаризация текстов

**Модели:** seq2seq, CNN, RNN, LSTM, GAN, BERT, GPT-3 и др.

# Синтез изображений и видео



(d) input image (e) output 3d face (f) textured 3d face



Source Subject

Target Subject 1

Target Subject 2

Caroline Chan, Shiry Ginosar, Tinghui Zhou, Alexei A. Efros. Everybody Dance Now. ICCV-2019.



# Машинное обучение – это оптимизация

$x$  – вектор объекта обучающей выборки

$w$  – параметры модели

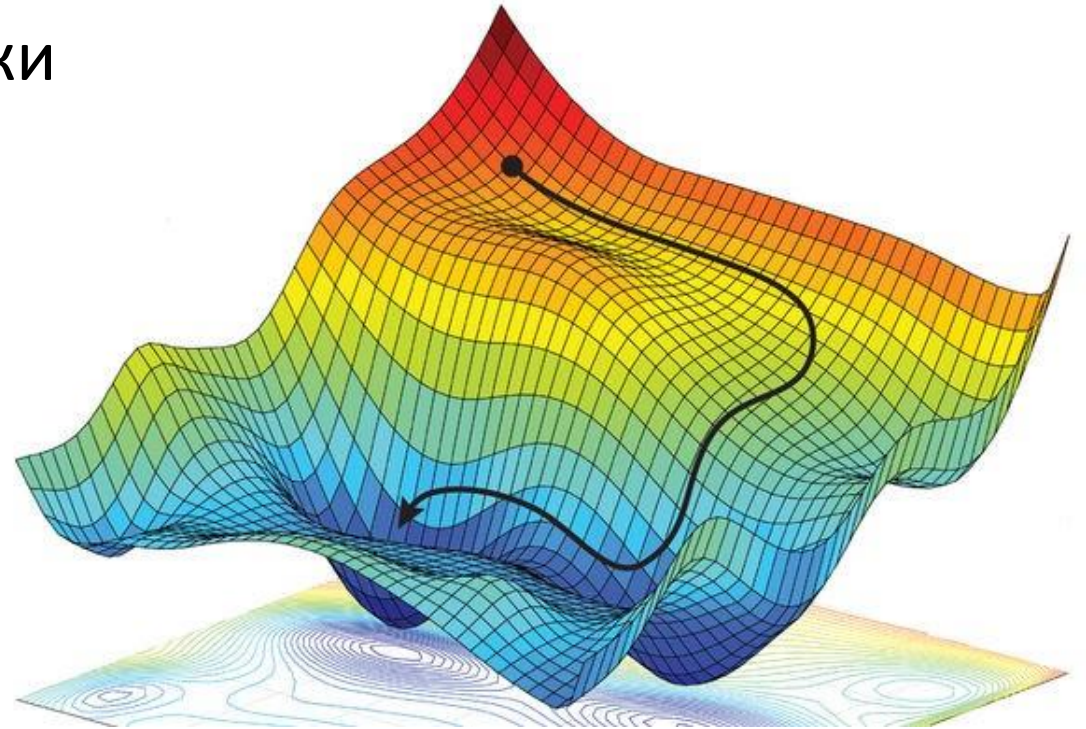
$\text{Loss}(x, w)$  – функция потерь

$Q(w)$  – критерий качества модели

Задача на этапе обучения модели:

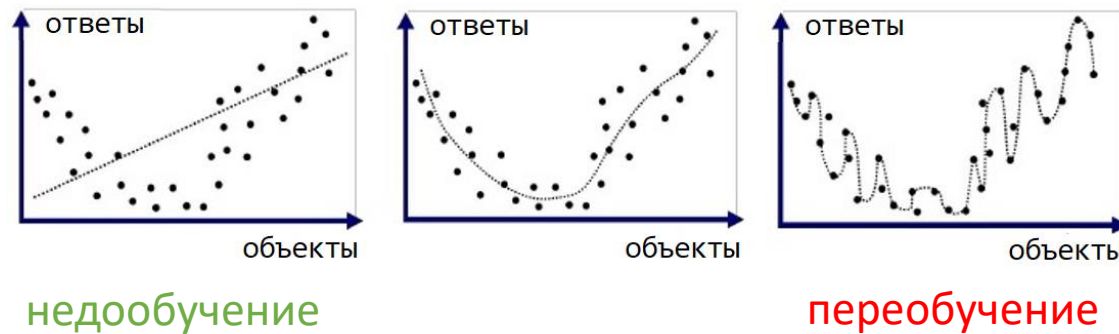
$$Q(w) = \sum_x \text{Loss}(x, w) \rightarrow \min$$

Способ решения – численные методы оптимизации

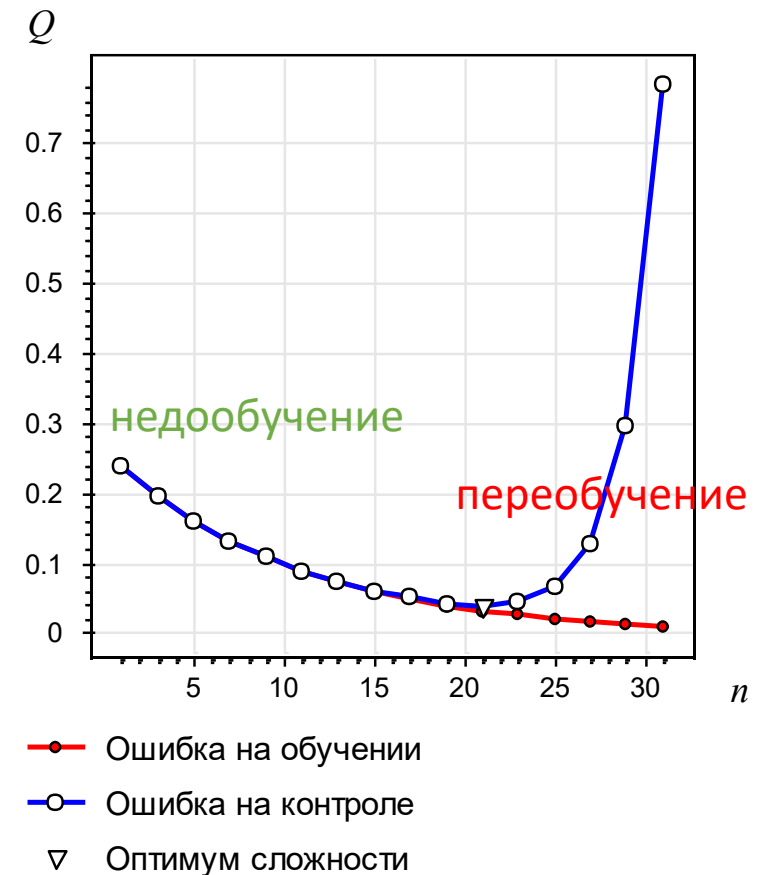


# Переобучение – основная трудность ML

Причина переобучения – избыточная сложность модели



- **Внутренние критерии:**  
для оптимизации параметров модели
- **Внешние критерии:**  
для оценивания обобщающей способности модели и контроля *переобучения*





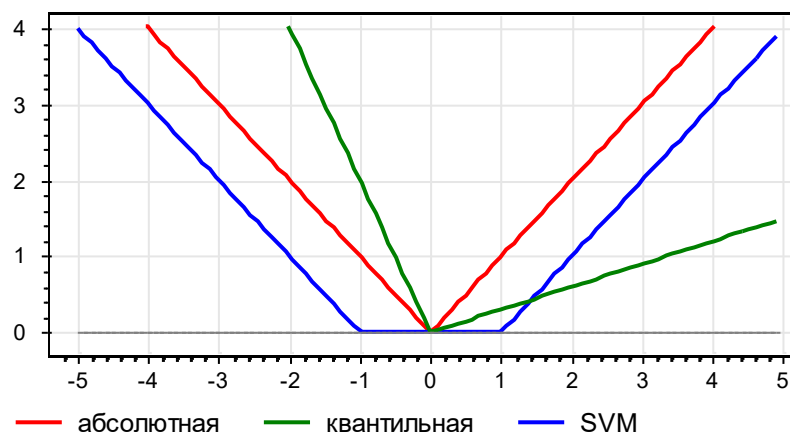
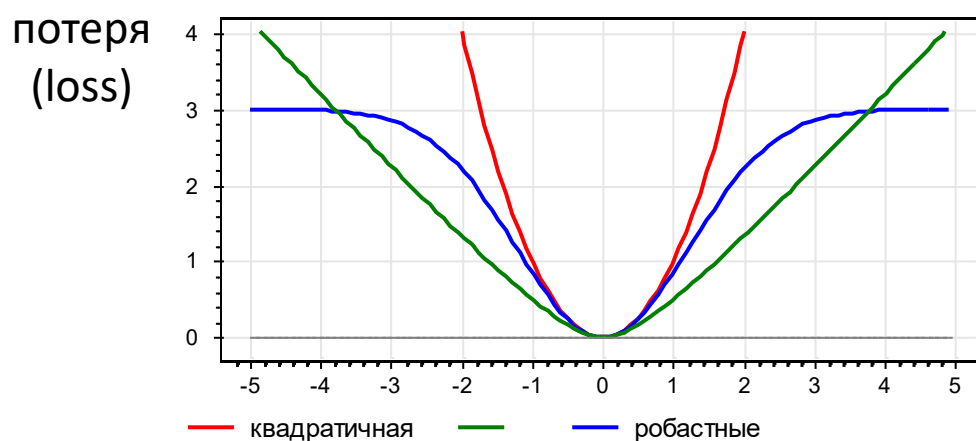
# Восстановление регрессии (regression)

$x$  – вектор объекта обучающей выборки,  $y$  – числовой ответ

$a(x, w)$  – модель регрессии с параметрами  $w$

Например,  $a(x, w) = \sum_j w_j x_j$  – линейная модель регрессии

$\text{Loss}(x, w) = (a(x, w) - y)^2$  – квадратичная функция потерь



НЕВЯЗКА  
(error)

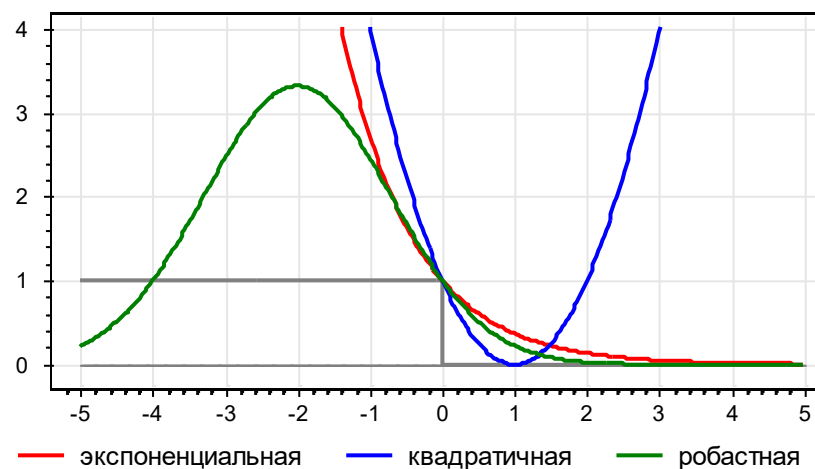
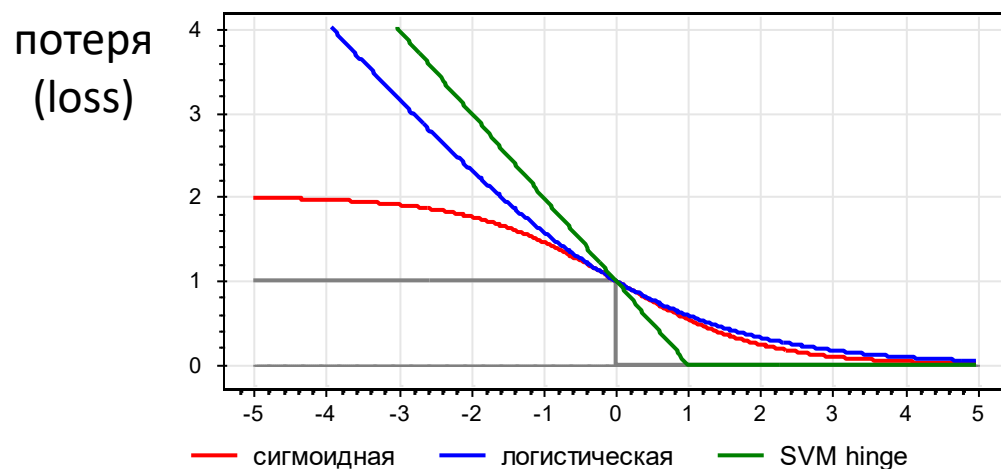
# Классификация (classification)

$x$  – вектор объекта обучающей выборки,  $y$  – ответ (+1 или -1)

$a(x, w)$  – модель классификации с параметрами  $w$

Например,  $a(x, w) = \text{sign}(\sum_j w_j x_j)$  – линейная модель

$\text{Loss}(x, w) = \max(0, 1 - y \sum_j w_j x_j)$  – функция потерь hinge



отступ  
(margin)

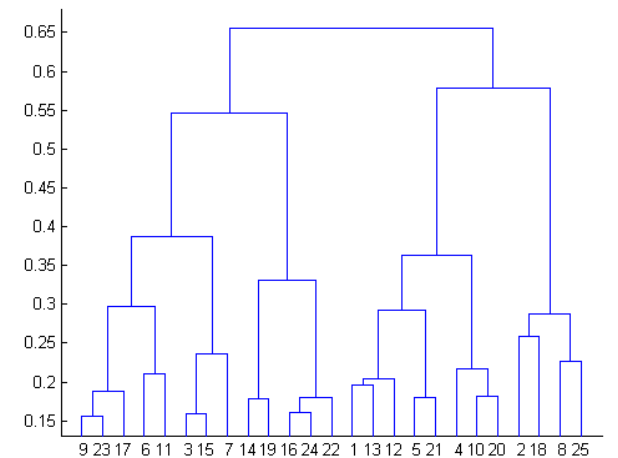
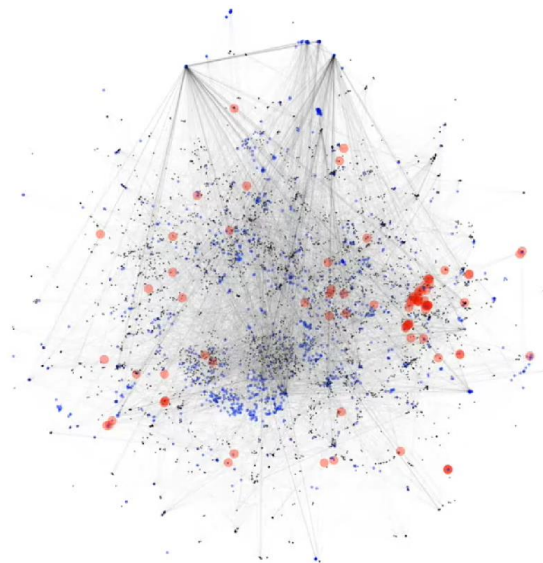
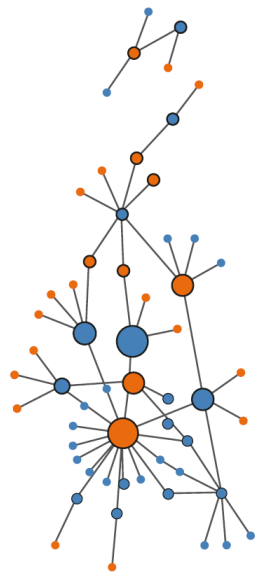
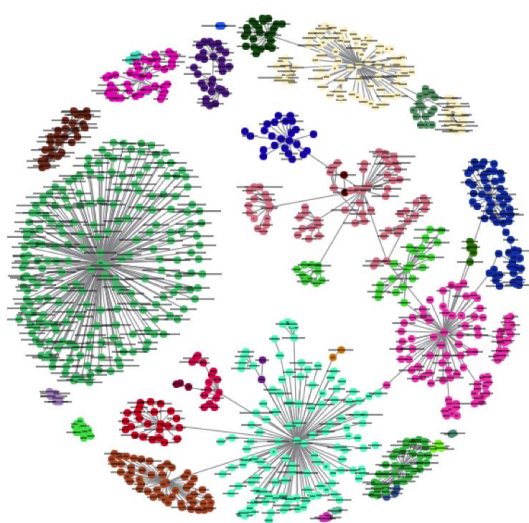
# Кластеризация (clustering)

$x$  – вектор объекта обучающей выборки, ответов не дано

$a(x, w)$  – ближайший к  $x$  центр кластера

$w = \{c_1, \dots, c_K\}$  – векторы центров всех кластеров

$\text{Loss}(x, w) = \min_k \|x - c_k\|$  – расстояние до ближайшего кластера



# Ранжирование (learning to rank)

$x$  – вектор пары «запрос-документ»,  $y$  – оценка релевантности

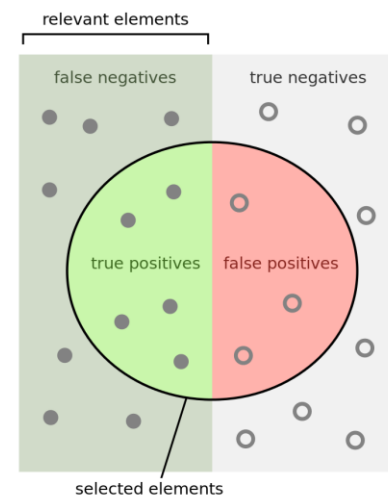
$a(x, w)$  – модель ранжирования документов по запросу, параметр  $w$

Например,  $a(x, w) = \sum_j w_j x_j$  – линейная модель

$\text{Loss}(x, x', w) = \max(0, 1 - [y > y'](a(x, w) - a(x', w)))$  – ф.потерь

историческая информатика

- Информатика историческая** litres. Без подписок  
litres.ru > Историческая-информа...   
Более 1 000 000 книг в форматах FB2, EPUB, TXT, PDF, Аудиокниги. Выберите и читайте! - Без подписок. Книга ваша навсегда. Все аудиокниги. Без скрытых платежей
- Историческая информатика** — Википедия  
ru.wikipedia.org > Историческая информатика   
Историческая информатика — междисциплинарная область исторических исследований, целью которой является расширение информационного...
- Журнал "Историческая информатика"**  
kleio.asu.ru   
Историческая информатика. Информационные технологии и математические методы в исторических исследованиях и образовании. Читать ещё >
- Методологические проблемы исторической информатики  
nbpublish.com > e\_jstinf/   
Ключевые слова: виртуальные исторические реконструкции, историческая информатика, источниковедение, методология, исторические источники, классификация, научно-техническая документация, электронные... Читать ещё >
- Историческая информатика.**  
ost-talent.org > 40526-istoricheskae-informatika- =



$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

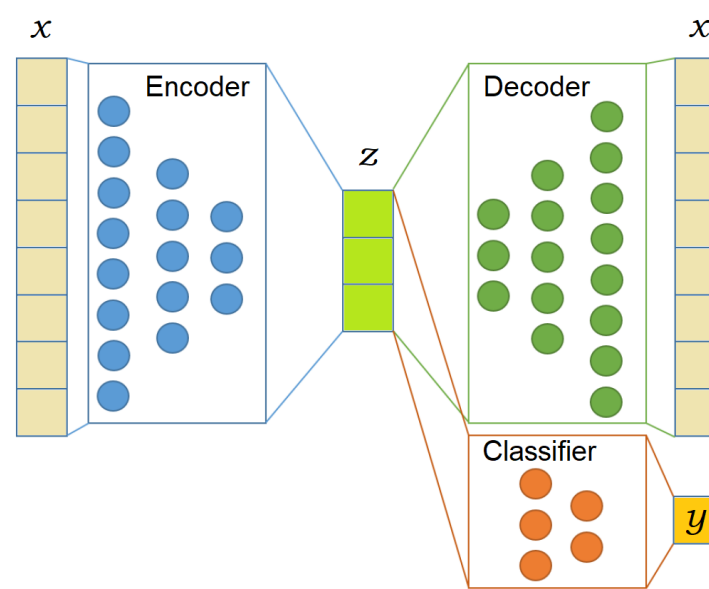
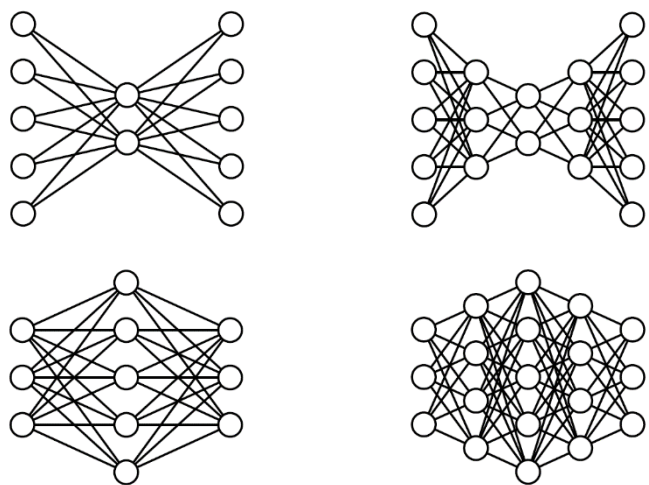
# Векторизация объектов (autoencoders)

$x$  – вектор объекта обучающей выборки, ответов не дано

$z = f(x, w)$  – модель кодирования  $x$  в векторное представление  $z$

$x' = g(z, w)$  – модель декодирования  $z$  в реконструкцию  $x'$

$\text{Loss}(x, w) = \|x'(w) - x\|$  – точность реконструкции объекта



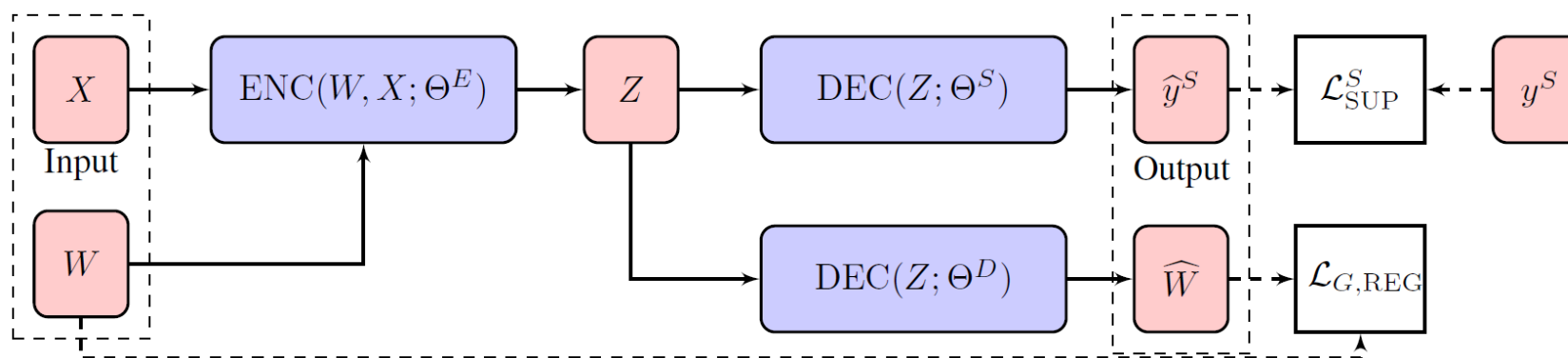
# Векторизация графов (graph embeddings)

$x; (x, x')$  – данные об объектах и взаимодействиях между объектами

$z = f(x, w)$  – модель кодирования  $x$  в векторное представление  $z$

$x' = g(z, w)$  – модель декодирования  $z$  в реконструкцию  $x'$

$\text{Loss}(x, w) = \|x'(w) - x\|$  – точность реконструкции объекта



*T.Mikolov et al.* Efficient estimation of word representations in vector space, 2013.

*I.Chami et al.* Machine learning on graphs: a model and comprehensive taxonomy. 2020.



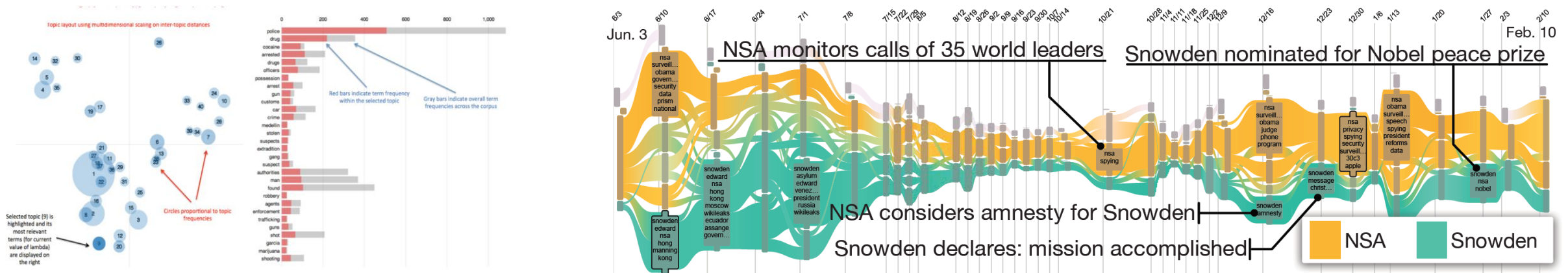
# Тематизация текстов (topic modelling)

$x$  – текстовый объект, последовательность слов или «мешок слов»

$z = f(x, w)$  – модель распределения текста  $x$  по темам  $z = p(t|x)$

$x' = g(z, w)$  – модель декодирования  $z$  в реконструкцию текста  $x'$

$\text{Loss}(x, w) = \|x'(w) - x\|_{KL}$  – точность реконструкции текста



# Тематический поиск (exploratory search)

## Схема эксперимента:

- Длинные запросы (1 стр. А4)
- 100 запросов на коллекцию
- 3 ассессора на каждый запрос
- 30 минут в среднем на запрос
- Разметка на Яндекс.Толока
- Две коллекции техно-новостей:

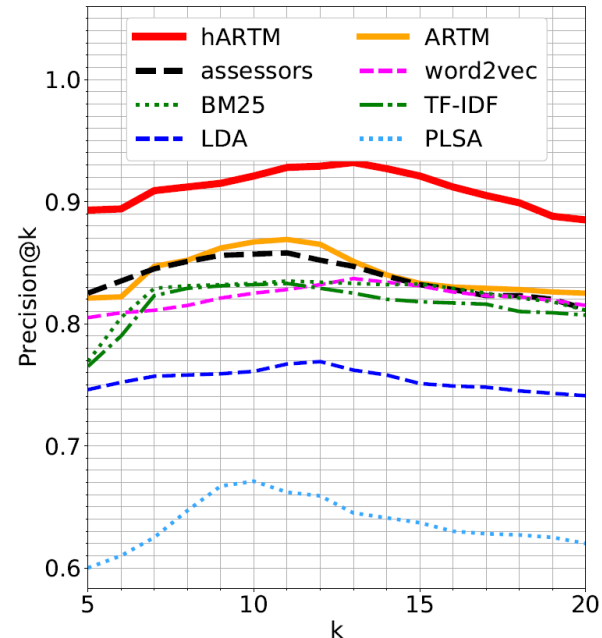


(170К текстов на русском) (750К текстов на английском)

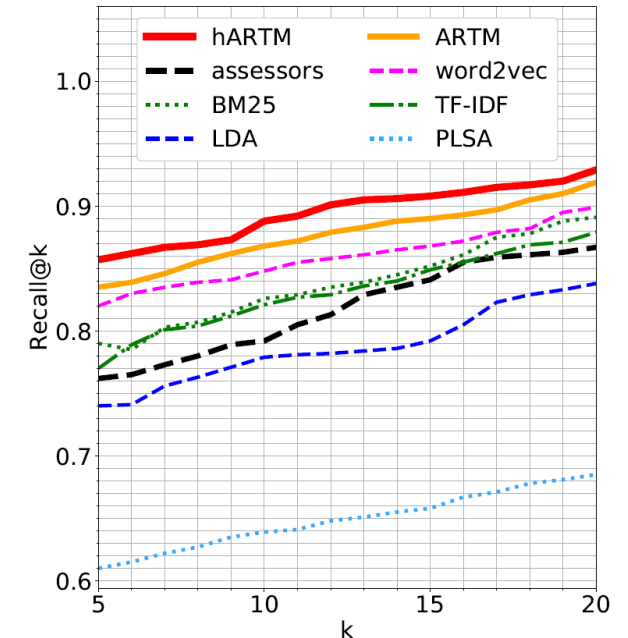


## Оценки качества поиска:

точность (precision@k)



полнота (recall@k)



Ianina A., Vorontsov K. [Regularized Multimodal Hierarchical Topic Model for Document-by-Documents Exploratory Search](#), 2019.

Ianina A., Golitsyn L., Vorontsov K. [Multi-objective topic modeling for exploratory search in tech news](#). 2017.

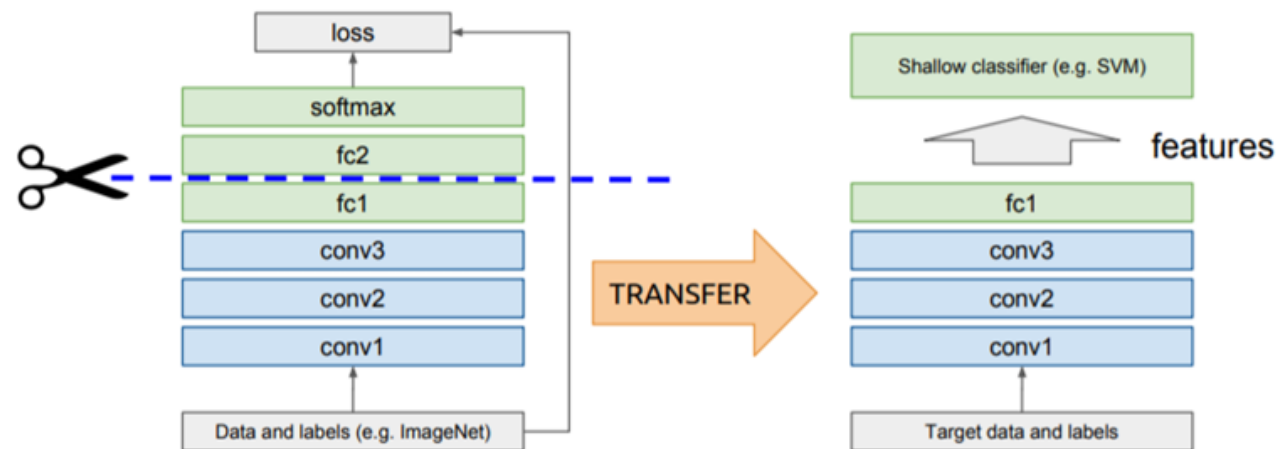
# Перенос обучения (transfer learning)

$f(x, w)$  – часть модели, универсальная для всех задач

$g(x, w')$  – часть модели, специфичная для каждой конкретной задач

$\min_{w, w'} \sum_x \text{Loss}_1(f(x, w), g_1(x, w'))$  – обучение по большим данным

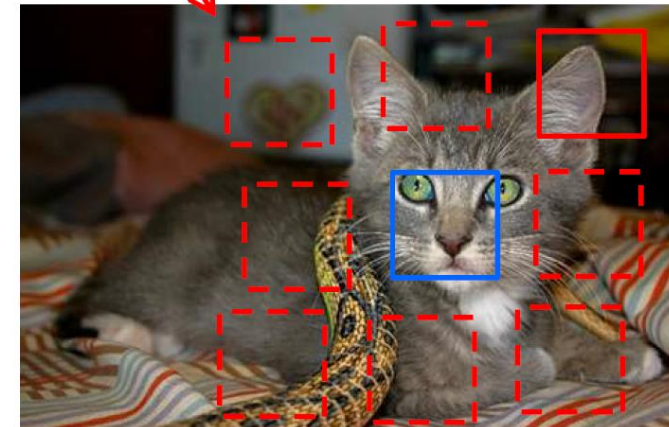
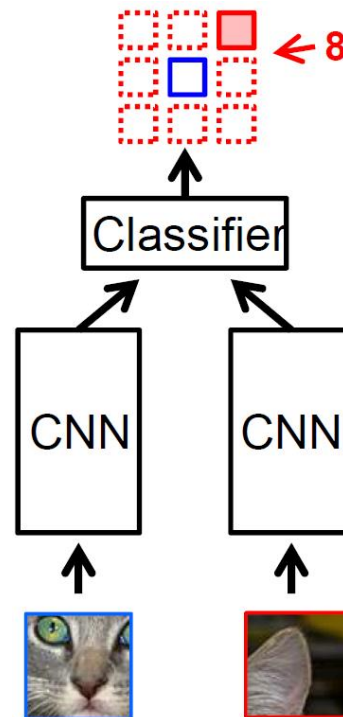
$\min_{w'} \sum_{x'} \text{Loss}_2(f(x', w^*), g_2(x', w'))$  – обучение по своим данным



# Самостоятельное обучение (self-supervised)

Модель векторизации  $z = f(x, w)$  обучается предсказывать взаимное расположение пар фрагментов одного изображения

**Преимущество:** сеть выучивает векторные представления объектов без размеченной обучающей выборки



Randomly Sample Patch  
Sample Second Patch

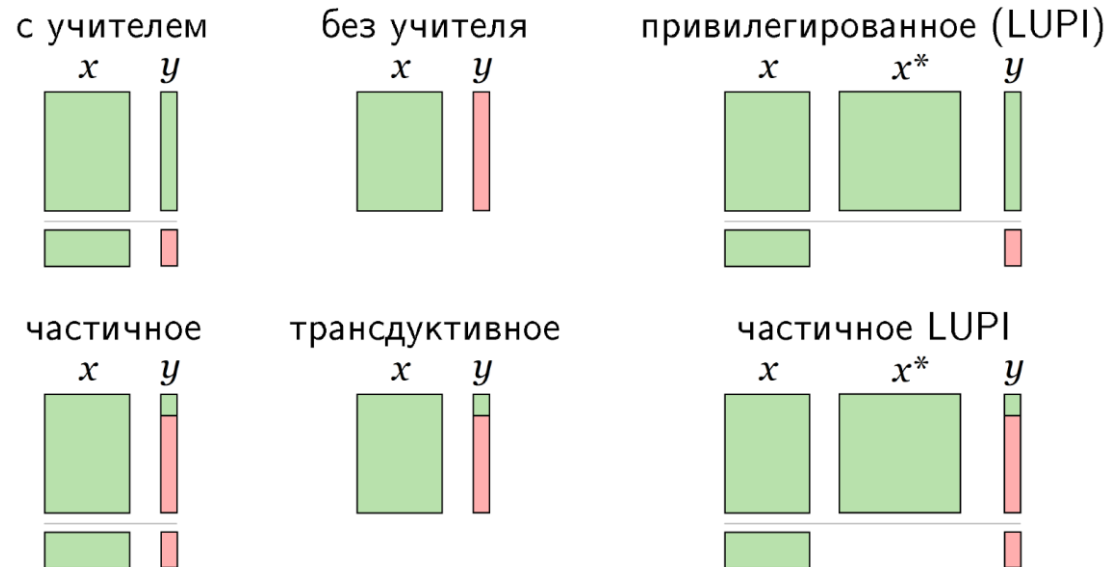
Unsupervised visual representation learning by context prediction,  
Carl Doersch, Abhinav Gupta, Alexei A. Efros, ICCV 2015

# Обучение с привилегированной информацией

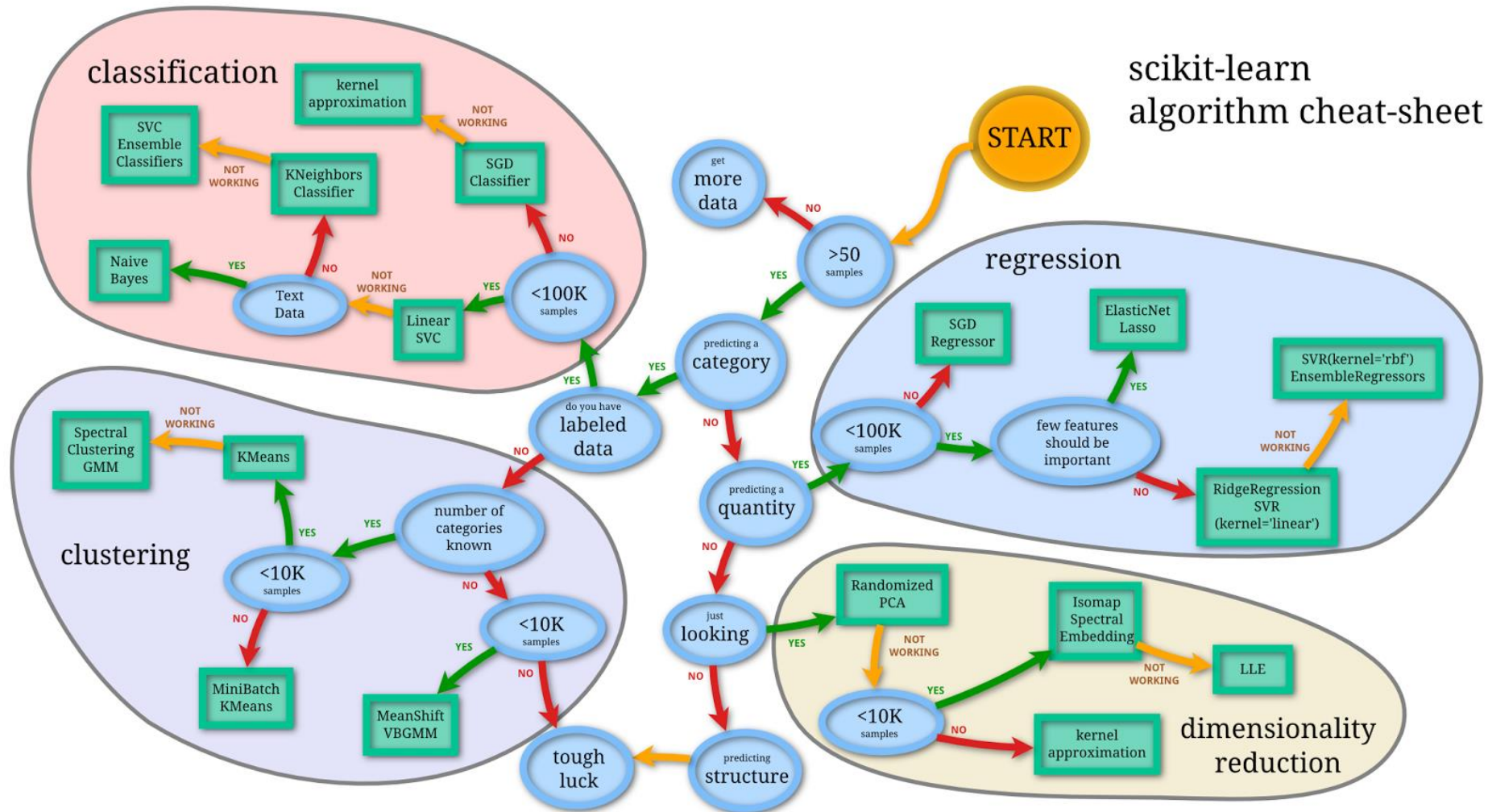
$g(x, x^*, w')$  – модель-учитель, имеет доступ к  $x^*$

$f(x, w)$  – модель-ученик, учится повторять ошибки учителя

$$\min_{w, w'} \sum_x \text{Loss}(f(x, w), y) + \text{Loss}(g(x, x^*, w'), y) + \text{Loss}(f(x, w), g(x, x^*, w'))$$



# Задачи и методы машинного обучения







# О методологии машинного обучения

## 1. Задачи машинного обучения

- Бум искусственного интеллекта и нейронных сетей
- Постановки задач и терминология машинного обучения
- Примеры задач машинного обучения

## 2. Методология машинного обучения

- Нейронные сети и глубокое обучение
- Обучение как задача оптимизации
- Типология и методология машинного обучения

## 3. Проблемы и перспективы применения

- Решение прикладных задач на практике
- Необходимые условия применения DS/AI/ML
- Мифы об искусственном интеллекте

# Особенности реальных данных

## В реальных приложениях данные бывают ...

- разнородные (признаки измерены в разных шкалах)
- неполные (признаки измерены не все, имеются пропуски)
- неточные (признаки измерены с погрешностями)
- противоречивые (объекты одинаковые, ответы разные)
- избыточные (сверхбольшие, не помещаются в память)
- недостаточные (объектов меньше, чем признаков)
- неструктурированные (нет признаковых описаний)
- «грязные» (ошибочные, грубо не соответствующие истине)

*со всем этим  
можно  
работать*

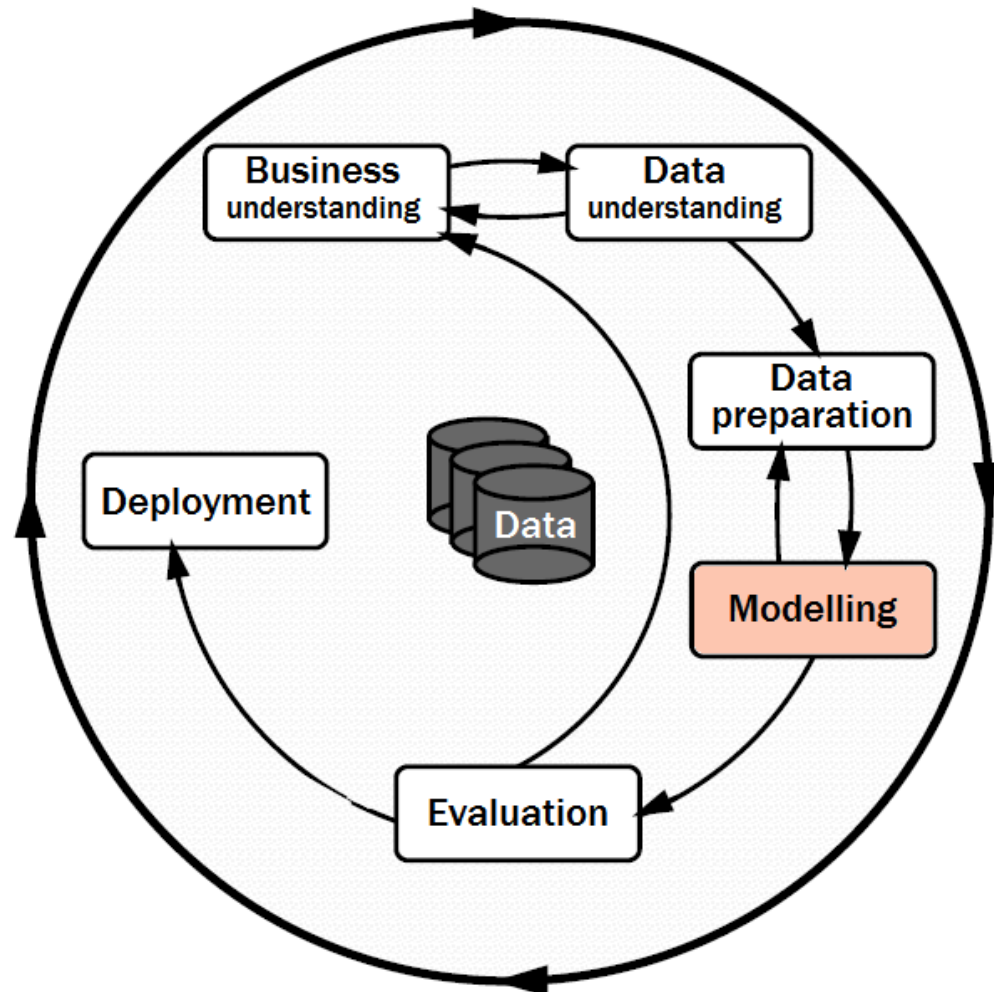


*но только не  
с грязными  
данными!*



# Этапы решения задач ML/DS/AI

CRISP-DM: Cross Industry Standard Process for Data Mining (1999)



- понимание прикладной задачи
- понимание данных
- предобработка данных
- инженерия признаков
- построение моделей
- оптимизация параметров
- контроль переобучения
- (кросс-)валидация решения
- внедрение и эксплуатация

# Необходимые условия применения ИИ

- **Полнота, чистота, достоверность данных**
  - Автоматизация и цифровизация процессов, порождающих данные
  - Контроль качества данных (цифровой двойник или «цифровое чучело»?)
  - Трудовая и технологическая дисциплина при работе с данными
- **Культура постановки задач**
  - Предметная экспертиза вместо «абстрактной веры во всемогущий ИИ»
  - Понимание целей анализа и их формализация в критериях качества
  - Готовность пилотировать новые технологии (кто на деле «data-driven»?)
- **Культура анализа данных**
  - Владение средствами визуализации и понимания данных
  - Грамотный и тщательный анализ ошибок при выборе моделей
  - Умение находить «простые но гениальные» решения

# Рынок труда в области анализа данных

## ***Инженер по данным (Data Engineer)***

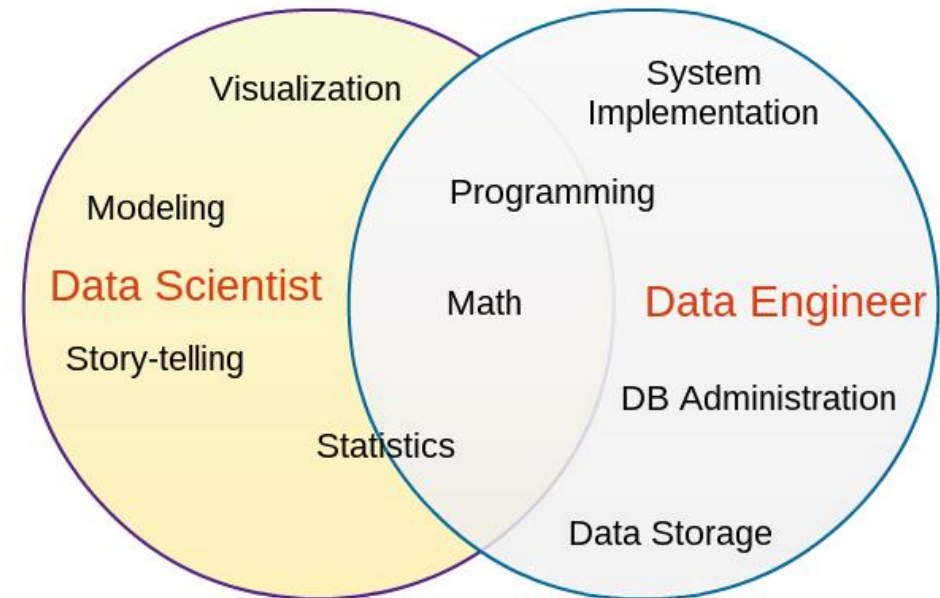
- Понимает бизнес-процессы, порождающие данные
- Работает с данными в различных форматах
- Визуализирует, понимает, очищает, готовит данные

## ***Исследователь данных (Data Scientist)***

- Моделирует, строит признаки (feature engineering)
- Выбирает модели и методы, оценивает решения
- Ходит по кругу CRISP-DM

## ***Менеджер проектов по анализу данных***

- Организует бизнес-процессы сбора и очистки данных
- Видит бизнес задачи и формализует их в терминах «Дано-Найти-Критерий»
- Организует открытые конкурсы и пилотные проекты
- Реалистично оценивает сложность задач и трудозатраты



# Миф №1

**«Сейчас наблюдается прорыв в области Искусственного Интеллекта»**

- Нет, это лишь прорыв в технологиях глубоких нейронных сетей
- Точнее, в технологиях эффективного решения задач численной оптимизации больших размерностей
- Прикладной потенциал этих технологий намного шире «ИИ»



# Миф №2

## «Скоро будет создан Общий Искусственный Интеллект (AGI)»

- Пока создаётся лишь функциональный **ИИ = Имитация Интеллекта**
- Artificial Intelligence – это мечта учёных, поэтичное название перспективного научного направления, придуманное в 1955 г.

### Отличия ИИ от естественного биологического интеллекта:

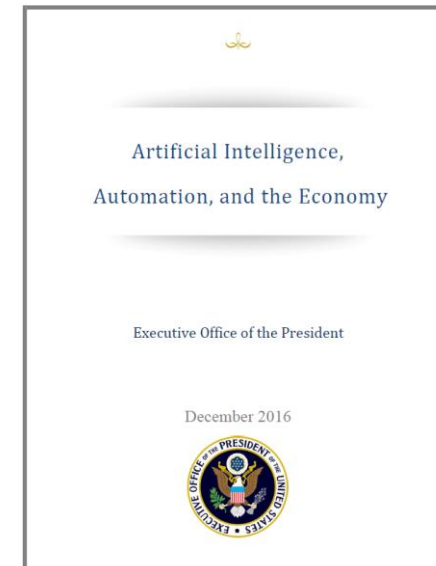
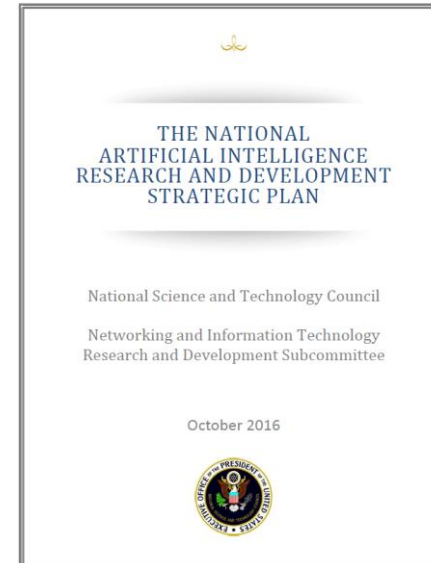
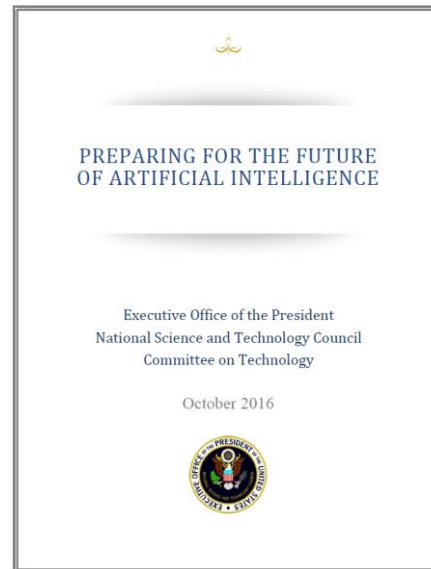
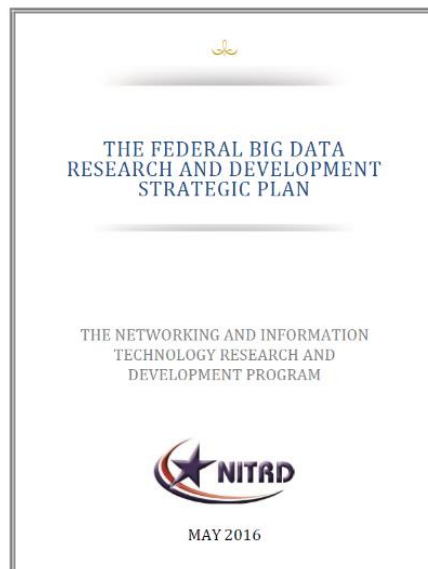
- Мы обучаемся не по выборкам, а на основе объяснений учителей, воспитания, опыта, коммуникации, то есть в естественной среде
- Мы строим картину мира и имеем целеполагание
- У нас 80 млрд. нейронов, и они устроены намного сложнее

# Миф №3

**«Тот, кто станет лидером в сфере ИИ, будет властелином мира»**

Отчёты Белого дома США, май-октябрь 2016:

**«Nations with the strongest presence in AI R&D will establish leading positions in the automation of the future»**



## Основные выгоды ИИ

- **Сокращение издержек и повышение производительности труда**
- Автоматизация банковских и финансовых услуг (FinTech)
- Автоматизация юридических услуг (LegalTech)
- Автоматизация посреднической деятельности, распределённая экономика
- Роботизация производств, автономный транспорт
- Оптимизация логистики и цепей поставок
- Оптимизация энергетических и транспортных сетей
- Сенсорные сети, мониторинг сельского хозяйства
- Персональная медицина, улучшение клинических практик
- Персональные образовательные траектории, социальная инженерия
- Автономные системы вооружений

## Некоторые из 23 рекомендаций

- #1. Организации должны активно развивать партнёрство с научными коллективами для эффективного использования данных.
- #2. В приоритетном порядке развивать стандарты *открытых данных* для привлечения научного сообщества к решению задач.
- #8. Инвестировать в разработку систем автоматического управления воздушным трафиком.
- #11. Вести постоянный мониторинг развития ИИ в других странах.
- #13. Приоритетно поддерживать фундаментальные и долгосрочные исследования в области искусственного интеллекта.
- #14. Развивать образовательные программы по ИИ и курсы повышения квалификации для прикладных специалистов.
- #20. Развивать международную кооперацию по ИИ.
- #22. Учитывать взаимовлияние ИИ и кибербезопасности.

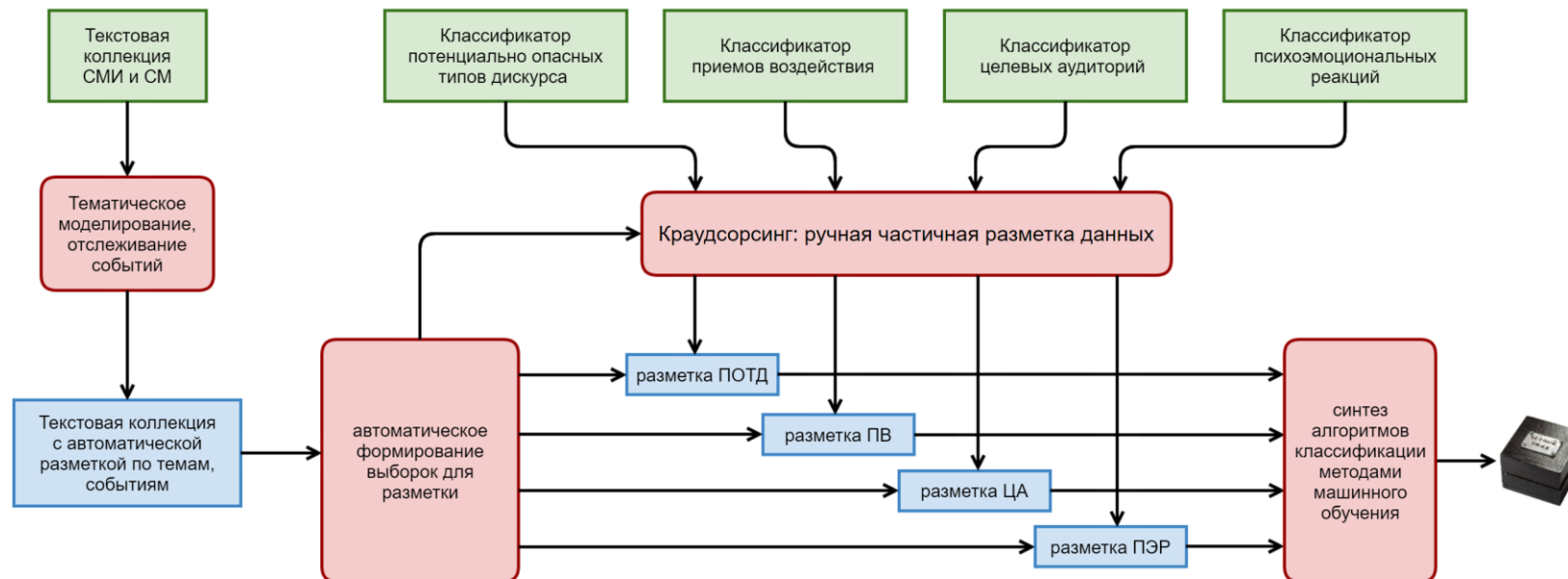
# Миф №4

## «ИИ является альтернативой классическому моделированию»

- То и другое – восстановление законов природы по данным
- ML: много данных, мало предметных знаний
- ML: универсальные аппроксиматоры для решения разных задач
- ML: модель может работать хорошо, но не понятно, почему
- ML: данные могут быть неполные, неточные, разнородные, ...
  
- **Нет чётких различий! Граница размыта!**

# Наши проекты в области анализа текстов

- Инструменты тематического моделирования BigARTM, TopicNet
- Поисково-рекомендательная система [<https://arxiv-search.mipt.ru>]
- Выявление потенциально опасных типов дискурса, приемов воздействия, целевых аудиторий и их психоэмоциональных реакций



# Рекомендуемая литература

- Бенджио И., Гудфеллоу Я., Курвилль А. Глубокое обучение. ДМК-Пресс, 2018.
- Николенко С. И., Кадурын А. А., Архангельская Е. О. Глубокое обучение. Питер, 2018.
- Воронцов К. В. Лекции по машинному обучению. [www.MachineLearning.ru](http://www.MachineLearning.ru), 2004-2020.
- Коэльо Л. П., Ричарт В. Построение систем машинного обучения на языке Python. 2016.
- Мерков А. Б. Распознавание образов. Введение в методы статистического обучения. 2011.
- Мерков А. Б. Распознавание образов. Построение и обучение вероятностных моделей. 2014.
- Hastie T., Tibshirani R., Friedman J. The Elements of Statistical Learning. Springer, 2014.
- Bishop C. M. Pattern Recognition and Machine Learning. - Springer, 2006.
- Домингос П. Верховный алгоритм. 2016.