

Введение в машинное обучение

Надежда Чиркова
факультет компьютерных наук НИУ ВШЭ



Является ли сообщение спамом?

Пишите Всероссийские проверочные работы (4, 5, 10, 11 класс)? ⚠️
Тогда Вам к нам! 💙

📌 Наша группа - [https://vk.com/club\[REDACTED\]](https://vk.com/club[REDACTED]) 📌
🔥 У нас будут материалы (задания и ответы) на все классы. 🔥

Русский язык на 18 и 20 апреля для 4-х классов уже есть! 📁
Спеши купить! За день до проведения станет дороже! ⚠️

Является ли сообщение спамом?

Что подарить, чтобы не прогадать? 🤔

Стилизованные портреты на холсте - отличный выход из ситуации!

Это индивидуальный подарок, при этом подходящий всем без исключения. 😊👍

➡ Группа - [https://vk.com/а\[REDACTED\]](https://vk.com/а[REDACTED])

👉 Оформите заказ у менеджера можно прямо сейчас -
[https://vk.com/\[REDACTED\]](https://vk.com/[REDACTED])

Является ли сообщение спамом?

Как составляют прогноз погоды, почему он не всегда точен и где лучше всего смотреть погоду:

<https://postnauka.ru/faq/75484>

Спам и не спам

Как научить компьютер автоматически отличать рекламные сообщения и удалять их?

Первое, что приходит в голову

```
if num_likes > 100:
    return False
elif "недорого" in s or "поможем" in s or "в наличии" in s:
    return True
elif num_smiles > 20:
    if s.count("!") > 10:
        return True
    elif "http" not in s:
        return True
else:
```

Второе, что приходит в голову

```
super_bad_words = ["лайк", "недорого", "уникальный", "жми"]
bad_words = ["новый", "лучший", "купите"]
good_words = ["считаю", "думаю", "нет"]
super_good_words = ["понравился", "худший"]
score = 0
score += sum([5 for word in s if word in super_bad_words])
score += sum([2 for word in s if word in bad_words])
score += sum([-2 for word in s if word in good_words])
score += sum([-5 for word in s if word in super_good_words])
if score > 0:
    return True
else:
    return False
```

В чем минус подхода?

- долго
- негибко
- неоптимально
- требует опыта

Мы уже накопили много опыта!

[http://vkontakte.ru/d\[REDACTED\]](http://vkontakte.ru/d[REDACTED]) 15..

Поставьте лайк,плииииз =))))

☐ Это спам

Альтернатива

Создавать алгоритм автоматически!

Данные = много примеров:

<p>Пишите Всероссийские проверочные работы (4, 5, 10, 11 класс)? 📌 Тогда Вам к нам! 💙</p> <p>📍 Наша группа - https://vk.com/club[REDACTED] 📍 📌 У нас будут материалы (задания и ответы) на все классы. 📌</p> <p>Русский язык на 18 и 20 апреля для 4-х классов уже есть! 📌 Спеши купить! За день до проведения станет дороже! 📌</p>	1
<p>Что подарить, чтобы не прогадать? 📌</p> <p>Стилизованные портреты на холсте - отличный выход из ситуации! Это индивидуальный подарок, при этом подходящий всем без исключения. 📌</p> <p>📌 Группа - https://vk.com/[REDACTED] 📌 Оформите заказ у менеджера можно прямо сейчас - https://vk.com/[REDACTED]</p>	1
<p>http://vkontakte.ru/d[REDACTED] 15.. Поставьте лайк, п्लीиз =))) 🚫 Это спам 🚫 Ответить</p>	1
<p>А вот и среда подходит к концу. Еще один успешно пройденный день на пути к цели! Кажется, даже сама погода радует, подарила нам теплый и солнечный денек. Сам мудренее вечера, в это значит, что нужно отложить всю работу, все хаоты по учебе до утра, потому что ночь была создана для того, чтобы отбросить все заботы прошедшего дня.</p>	0
<p>Чтобы не пропустить интересные события культурной жизни школы, подписывайтесь на телеграм канал: https://t.me/sch_inf</p> <p>Будем писать о концертах, лекциях, встречах с интересными людьми. Там же будут новости науки и рецензии на книги от учителей и учеников «Интеллектуала».</p>	0

Альтернатива

- ~~долго~~ быстро
- ~~негибко~~ один код на разные данные
- ~~неоптимально~~ лучше широкого круга решений
- ~~требует опыта~~ достаточно истории примеров

Альтернатива — машинное обучение

Цель: **автоматически** принимать решения



Используя свои знания,
писать программу для конкретного случая



Используя накопленные **данные**,
компьютер сам создает **оптимальную** программу

Альтернатива — машинное обучение

Цель: **автоматически** принимать решения



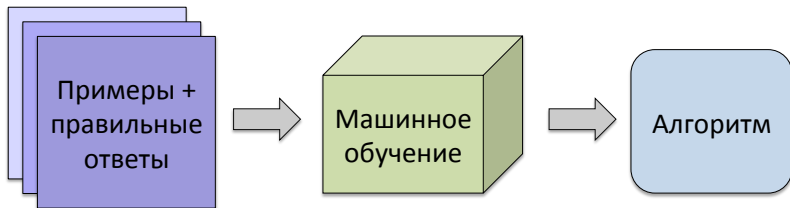
Используя свои знания,
писать программу для конкретного случая



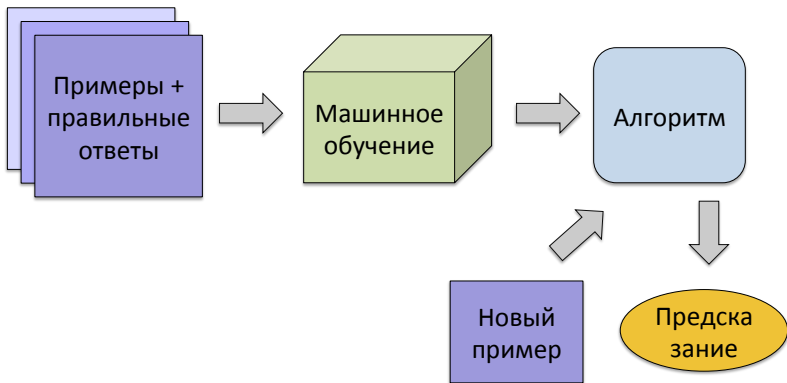
Используя накопленные **данные**,
компьютер сам создает **оптимальную** программу

Тоже программа!

Задача машинного обучения



Задача машинного обучения



Откуда данные?

Вариант А: разметка, выполняемая людьми
(машинное обучение для автоматизации того, что умеет делать человек)

- распознавание объектов на изображении, слов в речи
- перевод с одного языка на другой
- ответы на вопросы
- постановка медицинского диагноза
- интеллектуальные игры
- ...

Откуда данные?

Вариант Б: из жизни, науки, производства, бизнеса, ...
(машинное обучение для поиска зависимостей в данных, для предсказаний)

- медицина: болен ли пациент гриппом
- банки: вернет ли заемщик кредит
- торговля: купит ли покупатель товар
- Интернет: кликнет ли посетитель сайта ссылку
- экономика: обанкротится ли предприятие
- финансы: каким будет курс доллара
- ...

Медицинская диагностика

Данные — таблица с признаками и **правильными** ответами

Пол	Давление	Температура	Лейкоциты	...	Диагноз
жен	140/90	38.2	$15 \cdot 10^9$...	Грипп
муж	130/90	36.6	$7 \cdot 10^9$...	Здоров
муж	120/80	39.1	$11 \cdot 10^9$...	Ангина
		

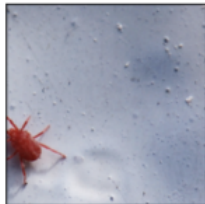
Медицинская диагностика

Данные — таблица с признаками и правильными ответами

Пол	Давление	Температура	Лейкоциты	...	Диагноз
жен	140/90	38.2	$15 \cdot 10^9$...	Грипп
муж	130/90	36.6	$7 \cdot 10^9$...	Здоров
муж	120/80	39.1	$11 \cdot 10^9$...	Ангина
		
жен	120/80	36.7	$8 \cdot 10^9$...	?

Распознавание изображений: Imagenet

10 млн изображений с ручной разметкой классов (виды животных, объекты мебели, машины, здания, ...)



mite



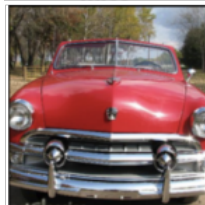
container ship



motor scooter



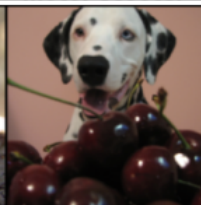
leopard



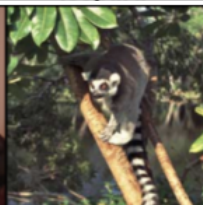
arille



mushroom



cherry



Madagascar cat

От 28.2% ошибок (2010) к 3.57% (2015)

Распознавание изображений: применение

Яндекс

мадагаскарская кошка



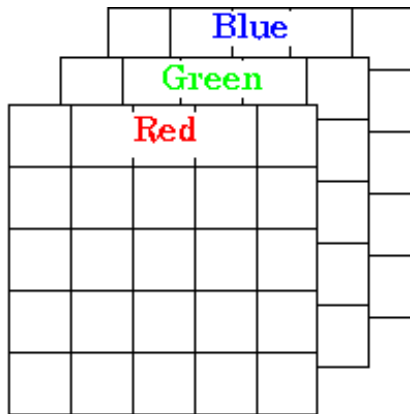
Найти

ПОИСК **КАРТИНКИ** ВИДЕО КАРТЫ МАРКЕТ НОВОСТИ ПЕРЕВОДЧИК ЕЩЁ

РАЗМЕР ▾ ОРИЕНТАЦИЯ ▾ ТИП ▾ ЦВЕТ ▾ ФАЙЛ ▾ ТОВАРЫ СВЕЖИЕ ОБОИ 1440×900 НА САЙТЕ



Представление изображений в компьютере



Предсказание возраста

У многих пользователей «В контакте» не указан возраст (или указан неверно :), хотя многим приложениям и рекламодателям эта информация нужна.

друзей	подписок	like'ов на аватарке	...	Возраст
110	20	75	...	18
7	11	3	...	45
30	4	10	...	24
		...		

Предсказание возраста

Или проще — моложе/старше 30

друзей	подписок	Like'ов на аватарке	...	Моложе 30?
110	20	75	...	да
7	11	3	...	нет
30	4	10	...	да
		...		

Два самых популярных вида задач

Классификация

предсказываем класс из конечного множества

предсказание моложе/старше 30

постановка диагноза

распознавание объектов на изображении

Регрессия

предсказываем вещественное (или целое, или натуральное) число

предсказание возраста

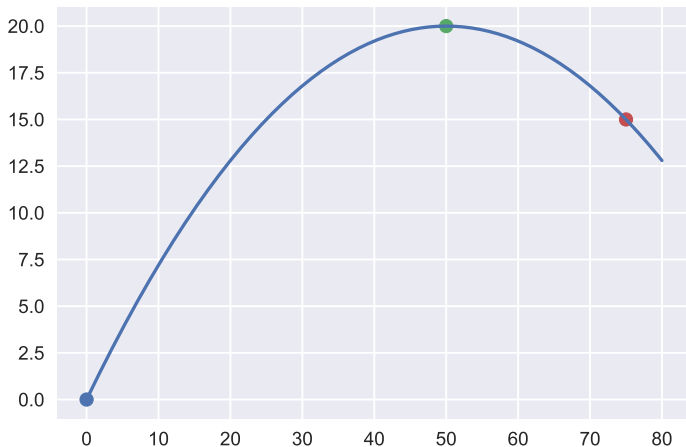
предсказание курса валют

предсказание температуры завтра

Задача про параболу

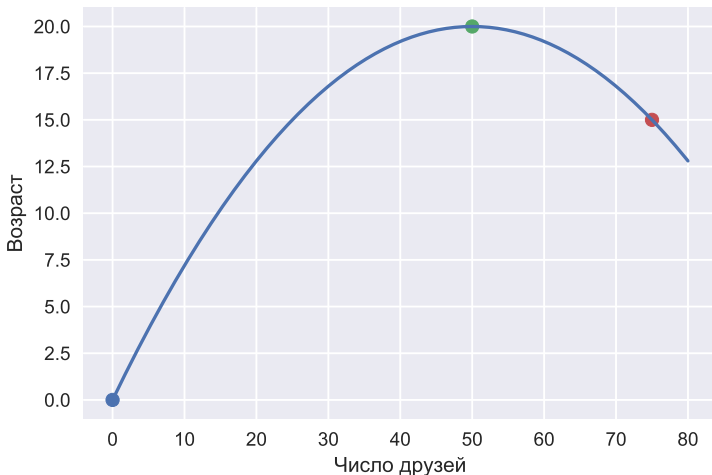
Парабола, проходящая через три точки:

$(0, 0)$, $(50, 20)$, $(75, 15)$



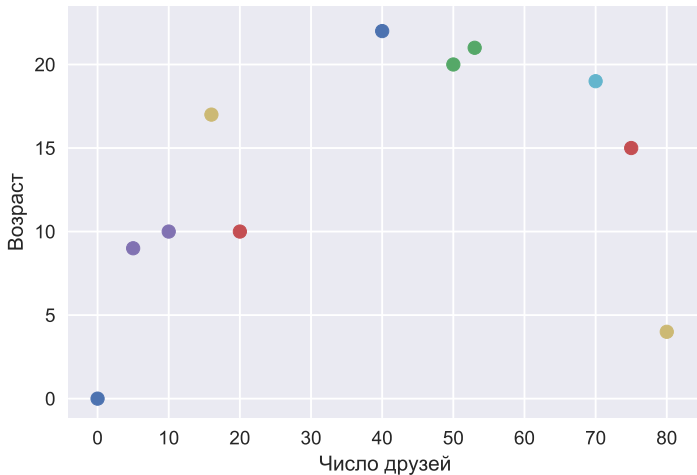
Предсказание возраста по числу друзей

Это решение задачи предсказания возраста по числу друзей, если объектов всего три!



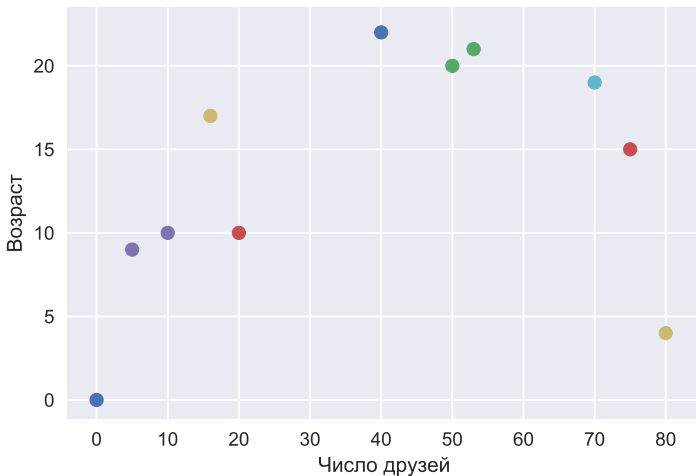
Предсказание возраста по числу друзей

А если объектов больше?



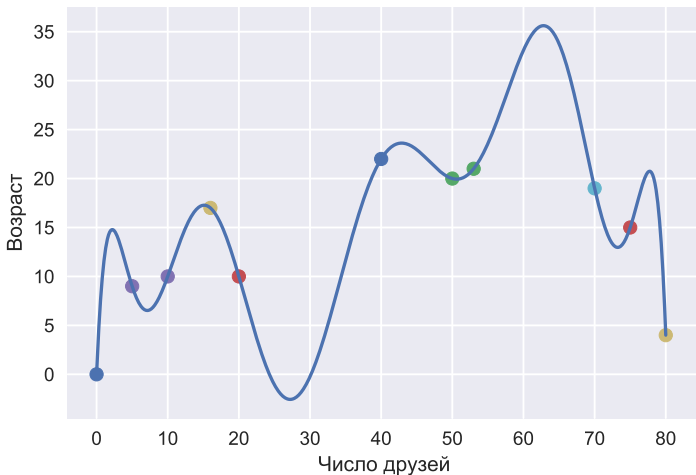
Предсказание возраста по числу друзей

А если объектов больше? Какой степени многочлен пройдет через все точки?



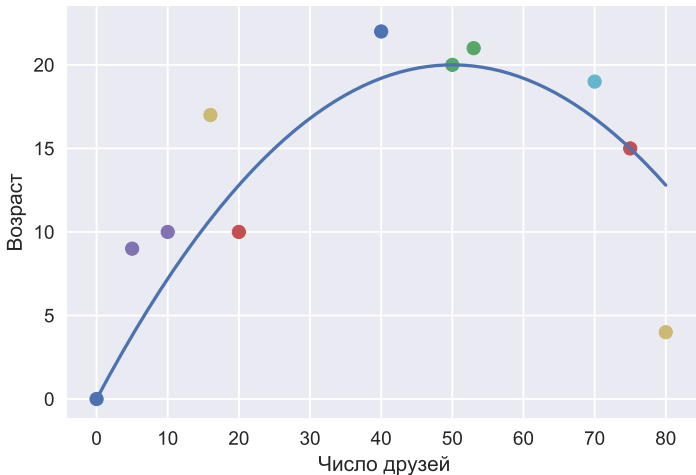
Предсказание возраста по числу друзей

Можно провести многочлен 10 степени через все точки, но если их еще больше?



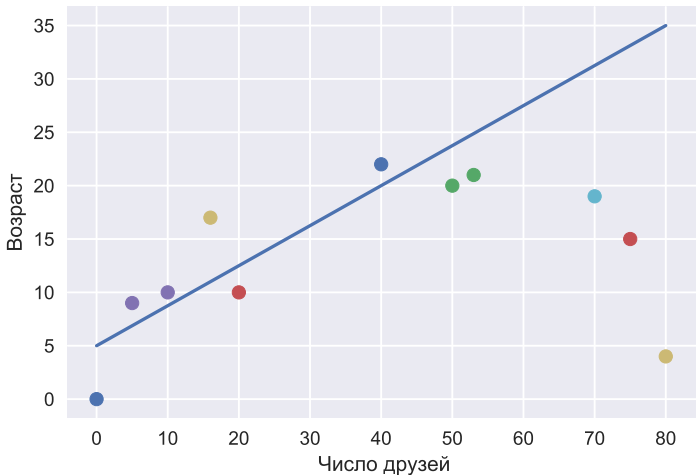
Предсказание возраста по числу друзей

Лучше приближенно, но проще, чем точно и сложно.



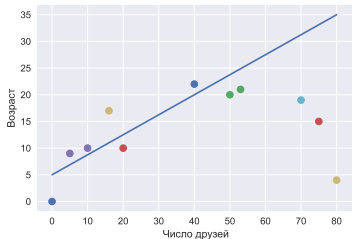
Предсказание возраста по числу друзей

Может еще проще?

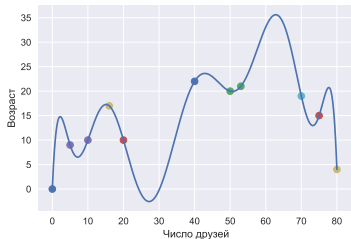


Баланс между точностью и обобщающей способностью

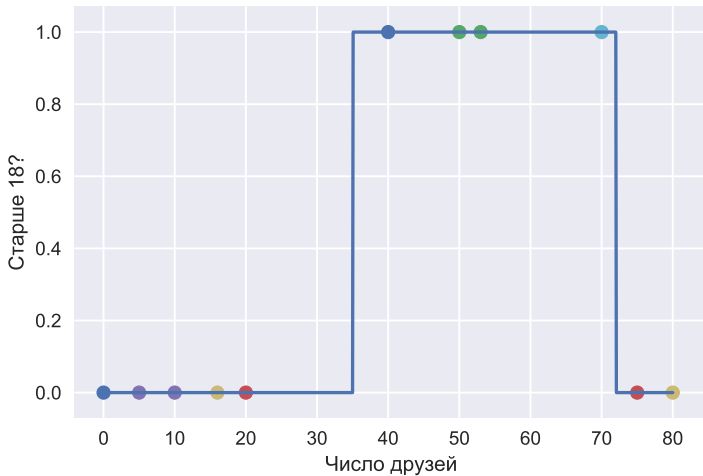
Недообучение



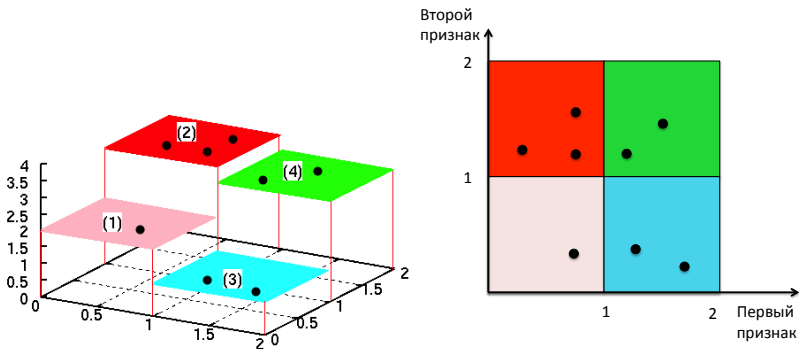
Переобучение



Задача приближения функции для классификации



Задача приближения функции для классификации

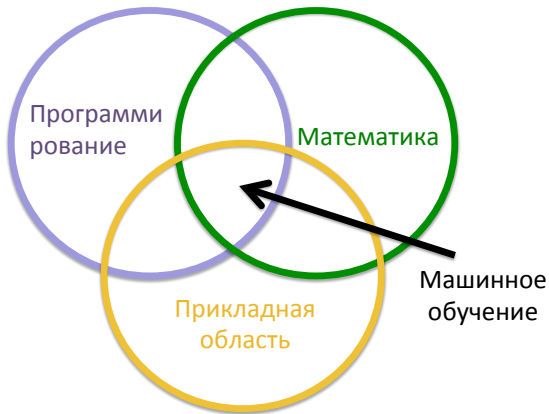


Итак,...

Машинное обучение:

- с точки зрения программирования: автоматическое создание программы, умеющей решать задачу предсказания, по данным
- с точки зрения математики: приближение функции по точкам
- находит применение во многих прикладных областях

Машинное обучение на стыке трех областей



Задачи машинного обучения

Задачи обучения с учителем: данные $X \in \mathbb{R}^{\ell \times d}$, $Y \in \mathbb{R}^{\ell}$
 ℓ – число объектов, d – число признаков

- Регрессия: $Y \in \mathbb{R}^{\ell}$
- Классификация:
 - бинарная $Y \in \{0, 1\}^{\ell}$
 - многоклассовая с непересекающимися классами
 $Y \in \{1, \dots, k\}^{\ell}$
 - многоклассовая с пересекающимися классами
 $Y \in \{0, 1\}^{\ell \times k}$
- Рекомендательные системы

Задачи обучения без учителя: данные $X \in \mathbb{R}^{\ell \times d}$

- Кластеризация: $X \in \mathbb{R}^{\ell \times d} \rightarrow Z \in \{1, \dots, c\}^{\ell}$
- Понижение размерности: $X \in \mathbb{R}^{\ell \times d} \rightarrow Z \in \mathbb{R}^{\ell \times d'}$, $d' < d$
- Визуализация: $X \in \mathbb{R}^{\ell \times d} \rightarrow Z \in \mathbb{R}^{\ell \times 2}$

Чего не умеет машинное обучение

Разумеется, есть много задач, которые решаются без данных и их анализа:

- построение оптимального маршрута на карте
- сложные физические расчеты
- рендеринг трехмерных изображений
- ...

Для их решения нужно хорошо изучать информатику, алгоритмы и средства программирования!

Поисковые сервисы

Яндекс

Новосибирские достопримечательности для |



Найти

новосибирск достопримечательности для детей

9 Новосибирск на карте России

yandex.ru/maps > Новосибирск

Новосибирск на карте России — схематической или спутниковой. Поиск на карте по адресу или названию населённого пункта.

W Новосибирск — Википедия

ru.wikipedia.org > Новосибирск

Новосибирск (произношение ; до 1925 года — Ново-Николаевск) — третий по численности населения и тринадцатый по занимаемой площади город России...

Новости Новосибирска — главные новости...

Новосибирская область

news.yandex.ru > Novosibirsk

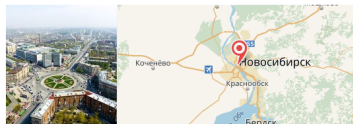
Новости Новосибирска. Выбрать другой регион. ... Завершён ремонт участка дороги на Немирова-Данченко в Новосибирске.

🖼 Новосибирск — смотрите картинки

yandex.ru/images > Новосибирск

Новосибирск

Город в России



Третий по численности населения и тринадцатый по занимаемой площади город России, имеет статус городского округа. [Википедия](#)

Погода: 13°C, Облачно

Местное время: 17 сентября, 00:15

Задачи машинного обучения:

- предсказание релевантности страницы запросу
- предсказание релевантности страницы интересам пользователя
- распознавание именованных сущностей

Поисковые сервисы

Яндекс

Новосибирские достопримечательности для |



Найти

новосибирск достопримечательности для детей

9 Новосибирск на карте России

yandex.ru/maps > Новосибирск

Новосибирск на карте России — схематической или спутниковой. Поиск на карте по адресу или названию населённого пункта.

W Новосибирск — Википедия

ru.wikipedia.org > Новосибирск

Новосибирск (произношение ; до 1925 года — Ново-Николаевск) — третий по численности населения и тринадцатый по занимаемой площади город России...

Новости Новосибирска — главные новости...

Новосибирская область

news.yandex.ru > Novosibirsk

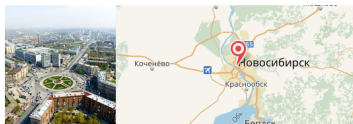
Новости Новосибирска. Выбрать другой регион. ... Завершён ремонт участка дороги на Немирова-Данченко в Новосибирске.

🖼 Новосибирск — смотрите картинки

yandex.ru/images > Новосибирск

Новосибирск

Город в России



Третий по численности населения и тринадцатый по занимаемой площади город России, имеет статус городского округа. [Википедия](#)

Погода: 13°C, Облачно

Местное время: 17 сентября, 00:15

Задачи программирования:

- быстрый отбор документов из огромной базы
- быстрое автозаполнение запроса
- выбор похожих запросов (с похожим множеством документов)

Self-driving cars



Задачи машинного обучения:

- распознавание знаков, сигналов светофора
- распознавание текстов на изображениях
- предсказание следующих действий пешеходов и других объектов
- планирование маневров
- прогнозирование времени прибытия в точку

Self-driving cars



Задачи программирования:

- поиск оптимального маршрута
- физические расчеты, моделирование движения

Машинное обучение в банках



Задачи машинного обучения:

- кредитный скоринг
- предсказание следующих транзакций пользователей
- кластеризация пользователей
- предсказание спроса на наличные деньги в банкоматах

Машинное обучение в банках



Задачи программирования:

- поддержка работы с базой клиентов банка
- написание программного обеспечения для банкоматов, онлайн-банков, создание приложений
- обеспечение информационной безопасности