

КОМПЬЮТЕРНАЯ ЛИНГВИСТИКА

Маша Шеянова, masha.shejanova@gmail.com

Саша Ершова, asershova@edu.hse.ru

December 16, 2016

НИУ ВШЭ

ЧТО ТАКОЕ "КОМПЬЮТЕРНАЯ ЛИНГВИСТИКА"?

Есть лингвистика. Есть компьютеры. Что хорошего можно с этим сделать?

1. Можно делать корпуса и вспомогательные инструменты для теоретических лингвистов.
2. Computational linguistics: изучение языка при помощи формальных математических моделей, статистики и всего такого.
3. Natural language processing: автоматическое извлечение чего-нибудь из текста и автоматическое его порождение.

NB: 2 и 3 — очень разные вещи, хотя и то, и другое в русском называют "компьютерной лингвистикой"

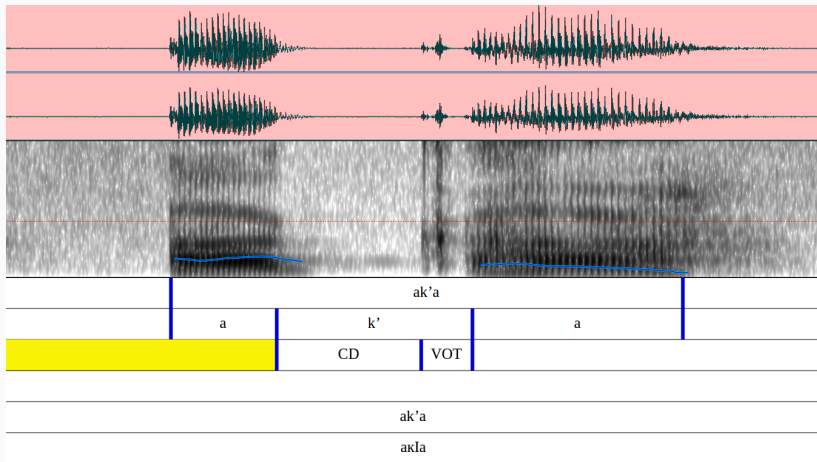
ВСПОМОГАТЕЛЬНЫЕ ИНСТРУМЕНТЫ

- Корпуса
- Словари
- Инструменты сбора данных
- Программы для анализа данных (анализ звука: Praat, анализ морфологии: Fieldworks)

сногшибательный компромат



АНАЛИЗ ДАННЫХ. PRAAT.



COMPUTATIONAL LINGUISTICS

ЧТО СЮДА ВХОДИТ?

В принципе, это любые лингвистические исследования, где нужно что-то посчитать, например:

- посмотреть, от чего возникают дырки в парадигмах
- доказать, что вид в русском — это континуум
- посмотреть, какие слова ближе по значению, а какие дальше (этим умеет заниматься дистрибутивная семантика)

Что мы хотим:

- формальный способ считать лексическую близость
- глобально: научить компьютер извлекать смыслы из текста

Как делать это автоматически?

Дистрибутивная гипотеза: значения слов полностью определяются их контекстами. Слова с похожими типичными контекстами имеют схожее значение.

Нам нужно:

- много текстов, чтобы картинка была репрезентативной
- посчитать в этих текстах взаимную встречаемость слов друг с другом
- найти слова, которые могут заменить друг друга и слова, у которых нет общих контекстов

Готово! Мы прекрасны и можем

- находить слова, близкие по значению к данному
- строить семантические пропорции
- строить семантические визуализации

ДИСТРИБУТИВНАЯ СЕМАНТИКА. ЭТО РАБОТАЕТ!

На rusvectors можно найти слова, наиболее близкие к данному, построить семантическую пропорцию и многое другое.

The screenshot displays the rusvectors web application interface. At the top, there are two input boxes: "человек_S" and "кошка_S". Below "человек_S" is a blue arrow pointing down to a box containing "нога_S". Below "кошка_S" is a blue arrow pointing down to a box containing "???", indicating a missing word in an analogy. To the right of these inputs is a search bar and a "Calculate!" button. Below the input boxes, there are two columns of results. The left column is titled "News corpus" and lists five words with their similarity scores: 1. ступня 0.430, 2. котенок 0.424, 3. кошачий 0.409, 4. пес 0.403, 5. ножка 0.388. The right column is titled "Ruscorpora" and lists five words with their similarity scores: 1. лапка 0.499, 2. ножка 0.485, 3. лапа 0.482, 4. ножища 0.482, 5. ножонка 0.479. Below these columns is a section titled "Web corpus" with five words and scores: 1. лапа 0.534, 2. ступня 0.519, 3. колено 0.508, 4. спина 0.484, 5. туловище 0.472. At the bottom, there is a "Choose the model:" section with a checked box for "Ruscorpora and Russian Wikipedia". Below that, there is a "Show only results which belong to:" section with radio buttons for "Nouns", "Verbs", and "Adverbs". At the very bottom, there is another "Choose the model:" section with checked boxes for "Ruscorpora and Russian Wikipedia", "News corpus", "Ruscorpora", and "Web corpus".

человек_S

нога_S

News corpus

1. ступня 0.430
2. котенок 0.424
3. кошачий 0.409
4. пес 0.403
5. ножка 0.388

кошка_S

???

Ruscorpora

1. лапка 0.499
2. ножка 0.485
3. лапа 0.482
4. ножища 0.482
5. ножонка 0.479

Web corpus

1. лапа 0.534
2. ступня 0.519
3. колено 0.508
4. спина 0.484
5. туловище 0.472

Choose the model:

☒ Ruscorpora and Russian Wikipedia

Show only results which belong to:

☐ Nouns ☐ Verbs ☐ Adverbs

Calculate!

Choose the model:

☒ Ruscorpora and Russian Wikipedia ☒ News corpus ☒ Ruscorpora ☒ Web corpus

NATURAL LANGUAGE PROCESSING

- Спеллчекеры
- Машинный перевод
- Text mining
- Speech recognition и OCR
- Когнитивные технологии: боты, weak AI, seq2seq-нейросети

У нас есть параллельные корпуса, то есть корпуса, где каждое предложение одного языка сопоставлено с предложением другого. С их помощью мы учим компьютер переводить предложения пользователя.

Английский	Японский
How much is that red umbrella?	Ano akai kasa wa ikura desu ka.
How much is that small camera?	Ano chiisai kamera wa ikura desu ka.

Corpus-based бывает:

- Statistical
- Example-based

Параллельные корпуса не используются. Часть информации хранится в словарях, часть прописана в правилах.

Как это работает в Apertium

- словари:
 - билингвальные: лексические соответствия
 - монолингвальные: парадигмы
- правила:
 - лексический выбор: сложно → difficult, complicated или complex?
 - разрешение морфологической омонимии
 - изменение структуры

Corpus-based:

- широко используется сейчас (Google, Яндекс)
- требует параллельные корпуса: чем больше, тем лучше
- в принципе, не требует лингвистических знаний

Rule-based:

- сейчас всё больше уступает статистическому, **НО**
- может применяться при отсутствии больших корпусов → можно работать с малыми языками!
- их можно постепенно улучшать
- требует лингвистических знаний

Автоматическое извлечение информации для:

- категоризации текстов
- информационного поиска
- извлечения информации

OCR — Optical Character Recognition — извлечение текста из картинки.

Speech recognition — извлечение текста из аудиозаписи.

Зачем нам это, если можно просто взять и послушать/почитать?

OCR — Optical Character Recognition — извлечение текста из картинки.

Speech recognition — извлечение текста из аудиозаписи.

Зачем нам это, если можно просто взять и послушать/почитать?

- **Невероятно много информации.**
- Возможность "на лету" проделывать с извлечённым текстом ещё какие-нибудь операции.

Например, машинный перевод надписей на улице.

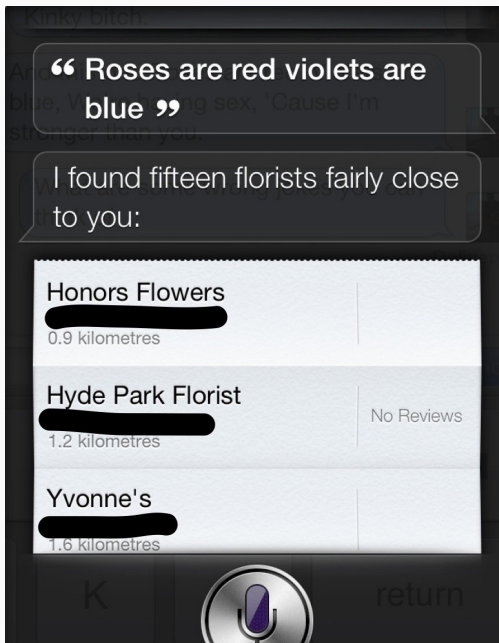


AI (Artificial Intelligence) — strong vs. weak.

strong AI — настоящий мыслящий искусственный интеллект, неотличимый от человека.

weak AI — штука, которая умеет выполнять некоторые когнитивные задачи, которыми обычно занимается человек.





Buy / Redeem Gift |
▼ | Your Account & Help

Watch Instantly
Browse DVDs
Your Queue
Movies You'll ♥
Give Netflix

Welcome,

Congratulations! Movies we think **You** will ♥

Add movies to your Queue, or **Rate** ones you've seen for even better suggestions.

The Scarlet Letter

Add

★★★★★

Not Interested

Unfaithful

Add

★★★★★

Not Interested

Two can play that game

Add

★★★★★

Not Interested

Indecent Proposal

Add

★★★★★

Not Interested

Same Time Next Year

Play Add

★★★★★

Not Interested

Whore

Add

★★★★★

Not Interested

Slutty Summer

Add

★★★★★

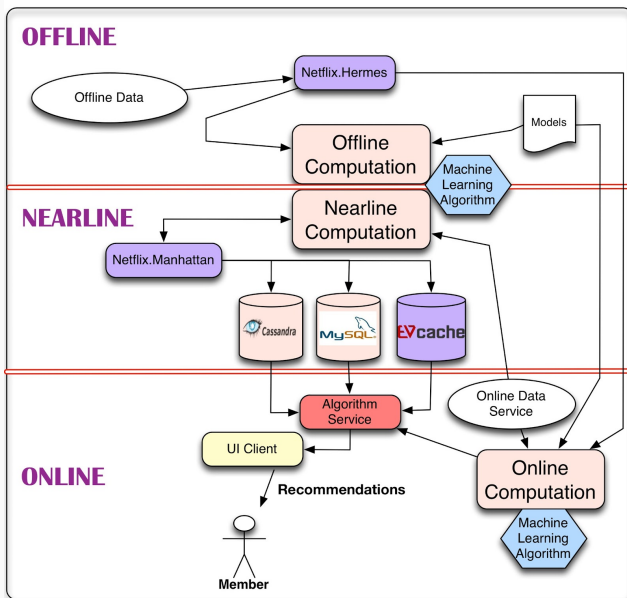
Not Interested

Bambi

Play Add

★★★★★

Not Interested



Strong AI пока не существует, но его хотят, боятся и ищут в существующих программах при помощи теста Тьюринга.

Что не так с тестом Тьюринга?

Weak AI есть вообще практически везде.

Нейросеть — это магический способ решения лингвистических (и не только) проблем. Она смотрит на данные и даёт правильные (обычно) ответы на вопросы про эти данные.

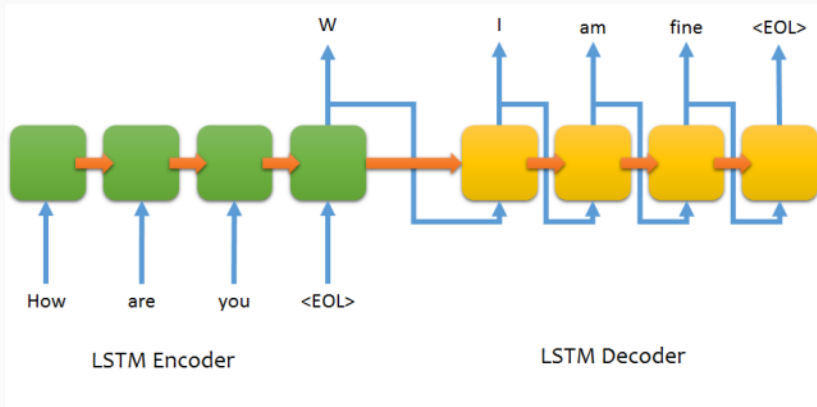
Только надо ~~выбрать правильное заклинание~~ правильно её сконфигурировать — это обычно самое сложное.

НС используются, в частности, для speech recognition и OCR.

Полным знанием о том, как работают нейросети, не обладает никто.

Нейросеть вида "sequence to sequence" принимает на вход некоторую последовательность чего угодно (чисел, пикселей, символов) и порождает другую последовательность чего угодно, соответствующую первой.

SEQ2SEQ-НЕЙРОСЕТИ



СПАСИБО ЗА ВНИМАНИЕ!