

Интерпретируемость NLP-моделей

Дмитрий Колодезев, Промсофт

NLP-митап в Омске

15.10.2022

Как бы план

- NLP и его друзья
- Что и зачем интерпретируем
- Трансформеры
- Fine-tuning
- Few-shot learning
- Чеклист

Однажды на Датафесте

Вопрос из зала: Трансформеры NLP, трансформеры неинтересны для визуализации?

Дмитрий Колодезев: Они интересны для визуализации, просто я не понимаю, какую я пользу могу из визуализации извлечь. У меня есть большая модель. Я смотрю в них и вижу, что как-то она не так визуализирует мои слова, мне не нравится. Что я с этим сделаю?

Если у меня есть маленькая табличная модель, я пошел ее доучил, посмотрел, сравнил.

А языковые модели столь сложны, что по большому счету их уже не столько учат, а в лучшем случае чуть-чуть доучивают. Или вообще подбирают подводки. То есть с языковыми моделями проблема, если вы не живете языковыми моделями, то, наверное, смотреть их лишний раз не надо, чтобы не расстраиваться.

Я не понимаю, какие выводы можно сделать посмотрев Bert, поэтому и не пользуюсь визуализацией, а так-то трансформеры интересны.

Про attention есть история, что attention – это та же самая saliency map в нейронках, то есть он подсвечивает не то, почему модель принимает решение, а то, на что она обращает внимание. Это могут быть просто необычные для набора данных вещи. То есть внезапно что-то странное появилось в наборе данных, модель обратит на это внимание. Необязательно, что она будет учитывать это для принятия решения.

Natural Language Processing

- Классификация текста
- Генерация текста
- Разбор текста на кусочки (NER)
- Заполнение пропусков в тексте
- Преобразование текста в текст
- Суммаризация
- ...

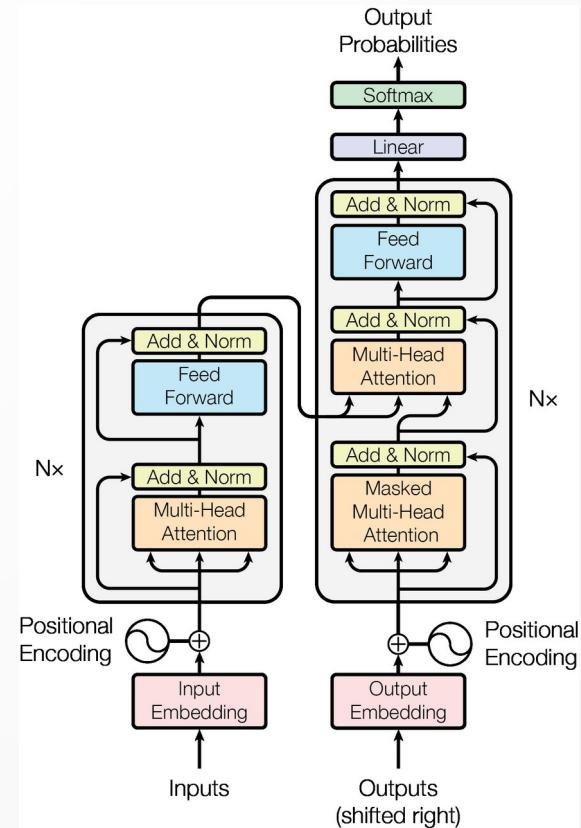
XAI: eXplainable AI

- Имеем право знать
- Отладка модели и данных
- Приемочное тестирование
- Анализ ошибок и разбор полетов
- Выявление уязвимостей
- Проверка очевидных зависимостей
- Социализация модели

https://en.wikipedia.org/wiki/Explainable_artificial_intelligence

В чем, собственно, проблема?

- Преобразования простые
- Но их очень много
- Очень-очень-очень-очень много
 - **YaLM** – 10^{11} параметров
- Не собираются в голове
- Как будто черный ящик какой-то



Или, говоря другими словами

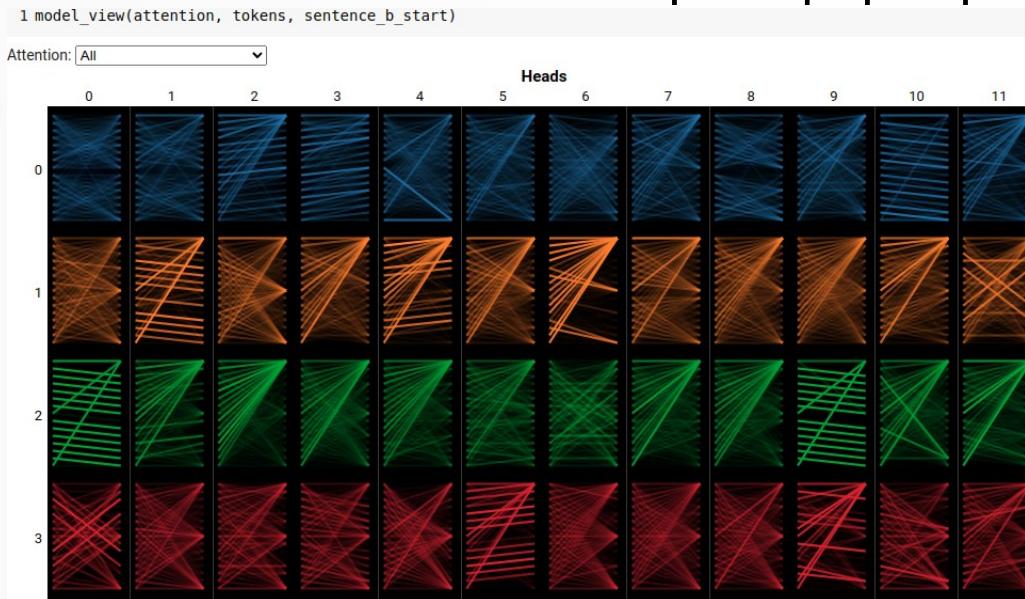
- «Despite constant advances and seemingly super-human performance on constrained domains, state-of-the-art models for NLP are imperfect. These imperfections, coupled with today's advances being driven by (seemingly black-box) neural models, leave researchers and practitioners scratching their heads asking, why did my model make this prediction?»

Как объясняют люди

- Вот по этому слову я сразу узнал в Вас интеллигента
- Он научился этой фразе у одноклассников
- Только не говорите ему про деньги, его это бесит
- Логичным продолжением сказанного будет ...
- Это все равно что сказать «я не буду вам помогать»

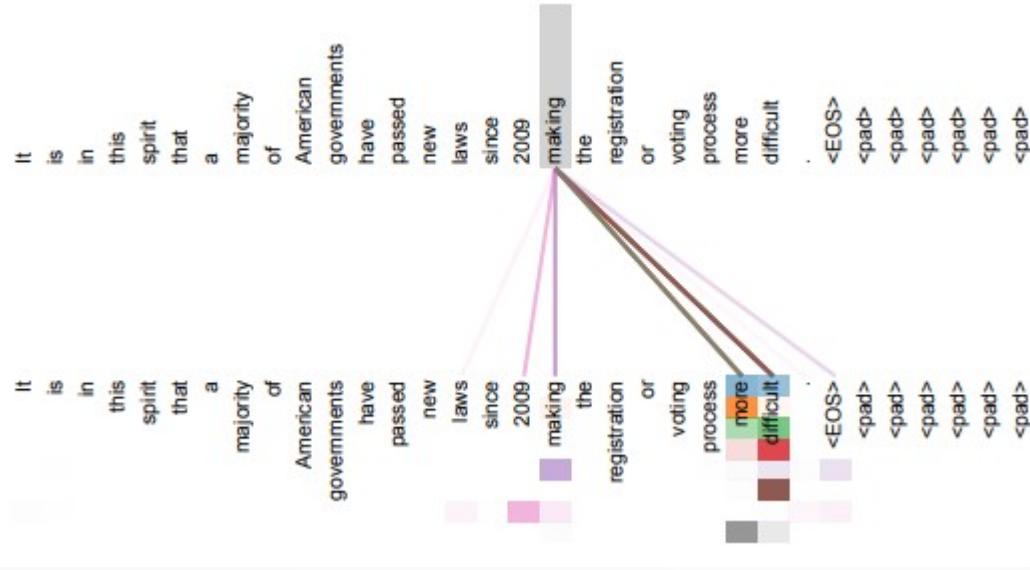
Начнем с трансформаторов трансформеров

- Внимание: главное в трансформерах — внимание!



<https://github.com/jessevig/bertviz>

Attention Is All You've Got



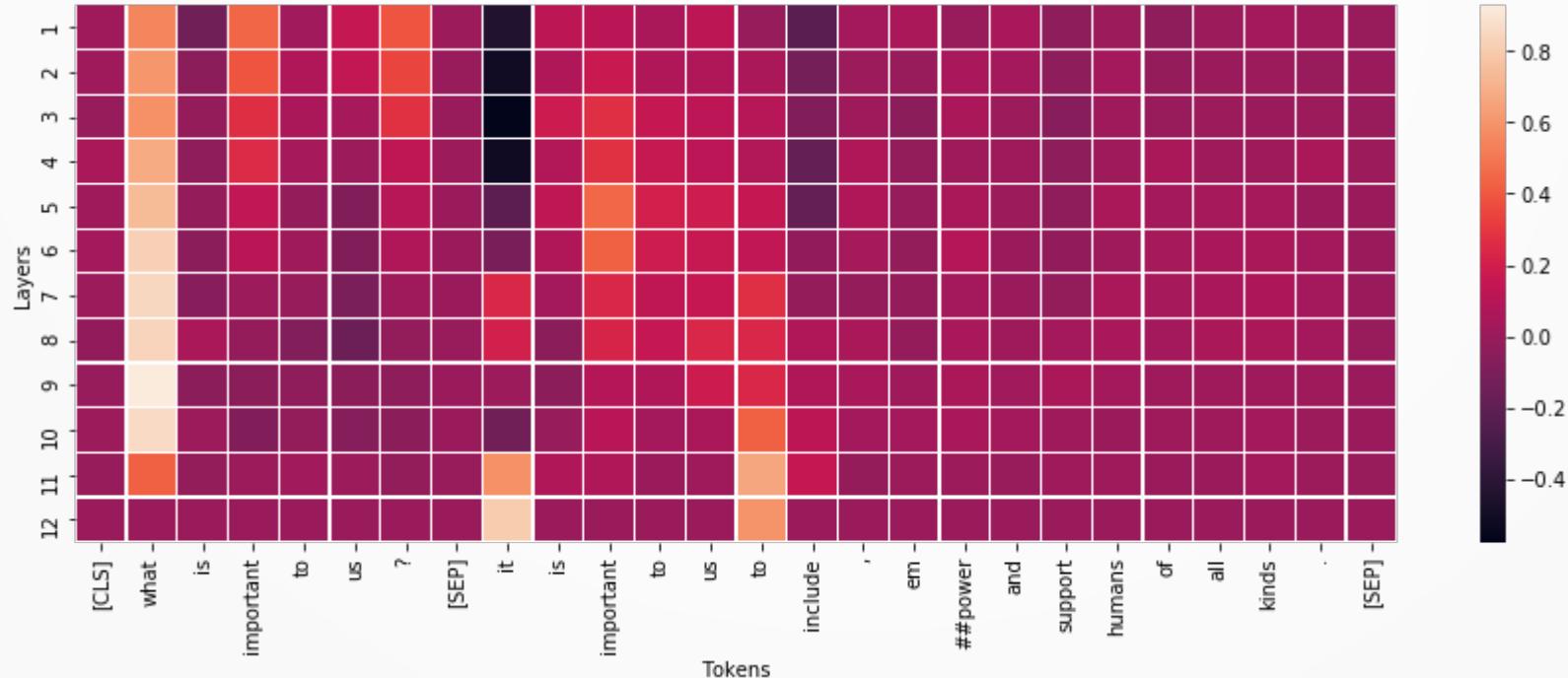
Attention Is All You Need

Может, градиент нарисуем?

Legend: ■ Negative □ Neutral ■ Positive

| True Label | Predicted Label | Attribution Label | Attribution Score | Word Importance |
|------------|-----------------|-------------------|-------------------|--|
| pos | pos (0.96) | pos | 1.29 | it was a fantastic performance ! #pad |
| pos | pos (0.87) | pos | 1.56 | best film ever #pad #pad #pad #pad |
| pos | pos (0.92) | pos | 1.14 | such a great show ! #pad #pad |
| neg | neg (0.29) | pos | -1.11 | it was a horrible movie #pad #pad |
| neg | neg (0.22) | pos | -1.03 | i 've never watched something as bad |
| neg | neg (0.07) | pos | -0.84 | that is a terrible movie . #pad |

Норма вектора имеет значение



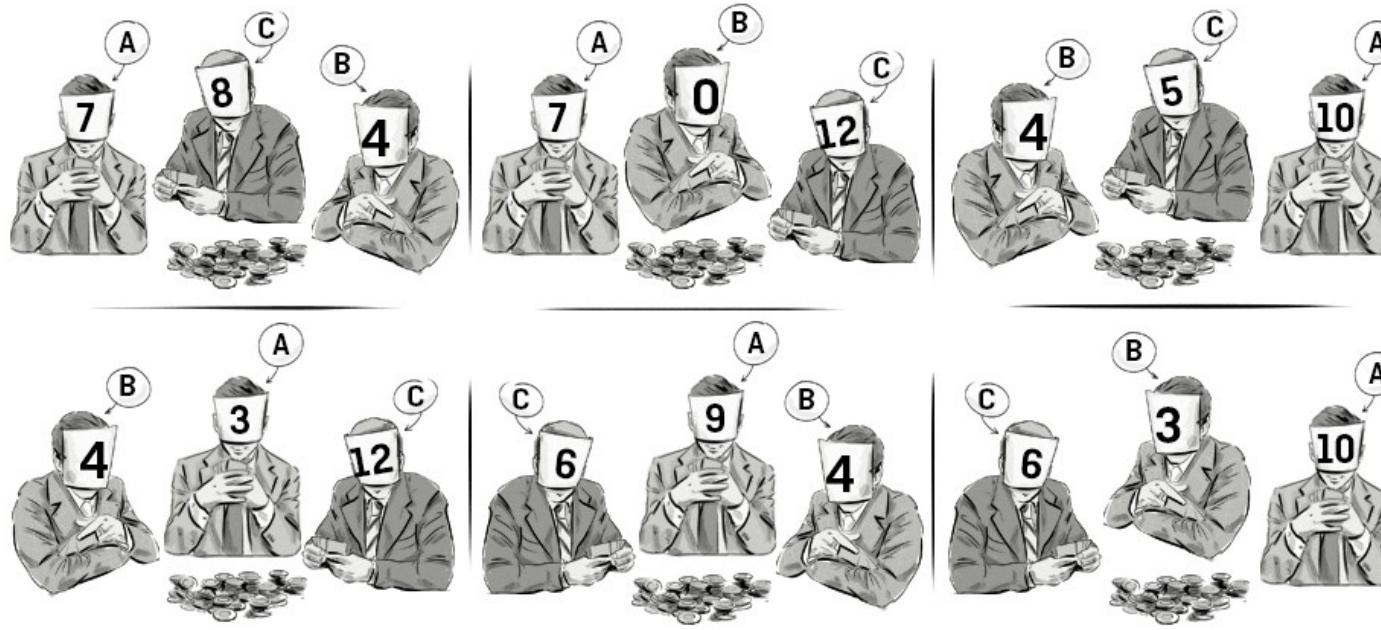
Attention is Not Only a Weight: Analyzing Transformers with Vector Norms

https://captum.ai/tutorials/Bert_SQuAD_Interpret2

Shapley Values

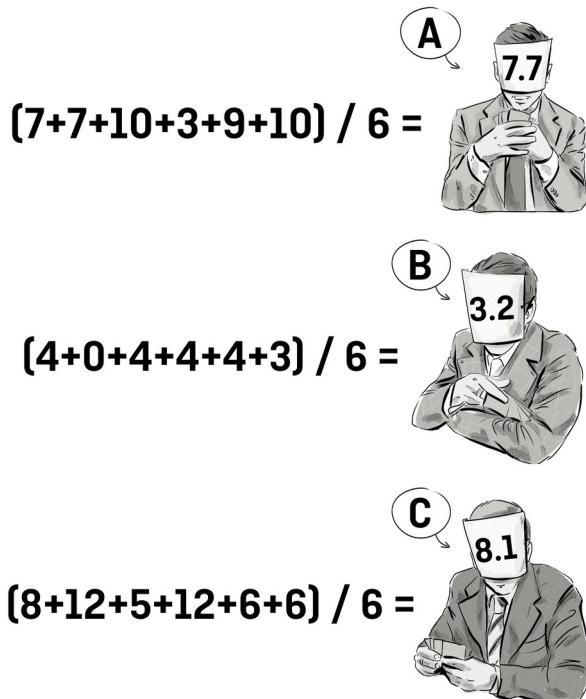


Shapley Values



<https://clearcode.cc/blog/game-theory-attribution/>

Shapley Values vs SHAP



A Unified Approach to Interpreting Model Predictions

Scott M. Lundberg
Paul G. Allen School of Computer Science
University of Washington
Seattle, WA 98105
slundi@cs.washington.edu

Su-In Lee
Paul G. Allen School of Computer Science
Department of Genome Sciences
University of Washington
Seattle, WA 98105
suinlee@cs.washington.edu

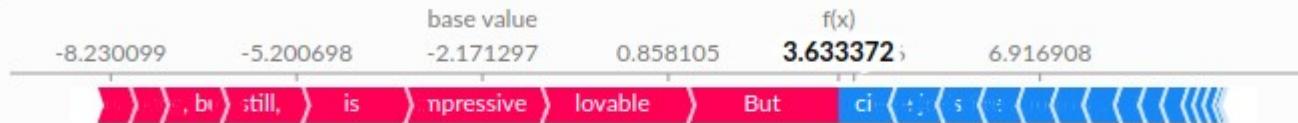
Abstract

Understanding why a model makes a certain prediction can be as crucial as the prediction's accuracy in many applications. However, the highest accuracy for large modern datasets is often achieved by complex models that even experts struggle to interpret, such as ensemble or deep learning models, creating a tension between *accuracy* and *interpretability*. In response, various methods have recently been proposed to help users interpret the predictions of complex models, but it is often unclear how these methods are related and when one method is preferable over another. To address this problem, we present a unified framework for interpreting predictions, SHAP (SHapley Additive exPlanations). SHAP assigns each feature an importance value for a particular prediction. Its novel components include: (1) the identification of a new class of additive feature importance measures, and (2) theoretical results showing there is a unique solution in this class with a set of desirable properties. The new class unifies six existing methods, notable because several recent methods in the class lack the proposed desirable properties. Based on insights from this unification, we present new methods that show improved computational performance and/or better consistency with human intuition than previous approaches.

<https://clearcode.cc/blog/game-theory-attribution/> <https://arxiv.org/pdf/1705.07874.pdf>

Нормально объясним

```
[10]: # plot the first sentence's explanation  
shap.plots.text(shap_values[3])
```



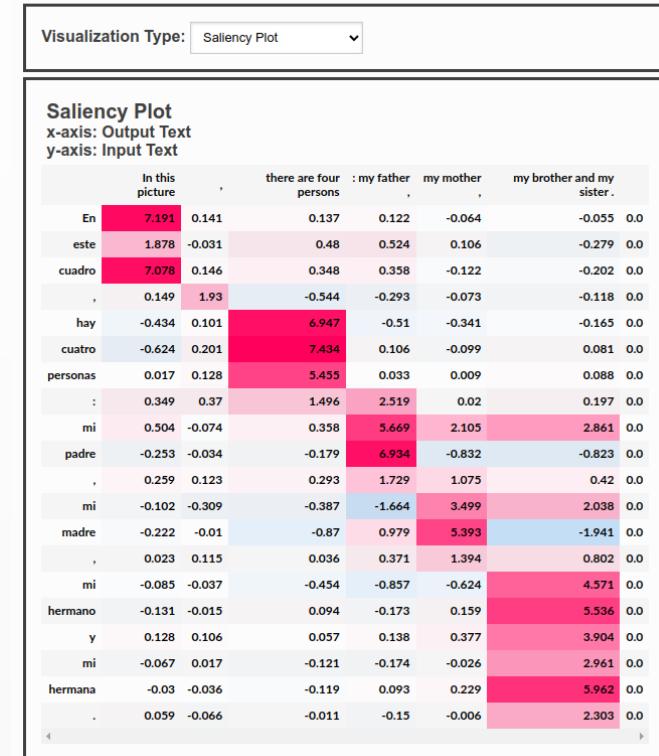
This is easily the most underrated film inn the Brooks cannon. Sure, its flawed. It does not give a realistic view of homelessness (unlike, say, how Citizen Kane gave a realistic view of lounge singers, or Titanic gave a realistic view of Italians YOU IDIOTS). Many of the jokes fall flat. But still, this film is very lovable in a way many comedies are not, and to pull that off in a story about some of the most traditionally reviled members of society is truly impressive. Its not The Fisher King, but its not crap, either. My only complaint is that Brooks should have cast someone else in the lead (I love Mel as a Director and Writer, not so much as a lead).

Понятно покажем

- In this picture , there are four persons : my father , my mother , my brother and my sister .
- En este cuadro , hay cuatro personas : mi padre , mi madre , mi hermano y mi hermana .

```
shap.plots.text(shap_values)
```

0th instance:



https://shap.readthedocs.io/en/latest/example_notebooks/api_examples/plots/text.html#

Проблема post hoc объяснений

- Люди верят красивым картинкам [Interpreting Interpretability](#)
- «Эффект Расемон» [Leo Breiman](#)
- Еще и патологии эти ваши

SQuAD

Context

In 1899, John Jacob Astor IV invested \$100,000 for Tesla to further develop and produce a new lighting system. Instead, Tesla used the money to fund his Colorado Springs experiments.

Original

What did Tesla spend Astor's money on ?

Reduced

did

Confidence

0.78 → 0.91

SNLI

Premise

Well dressed man and woman dancing in the street

Original

Two man is dancing on the street dancing

Reduced

Contradiction

Answer

0.977 → 0.706

VQA



Original

What color is the flower ?

Reduced

flower ?

Answer

yellow

Confidence

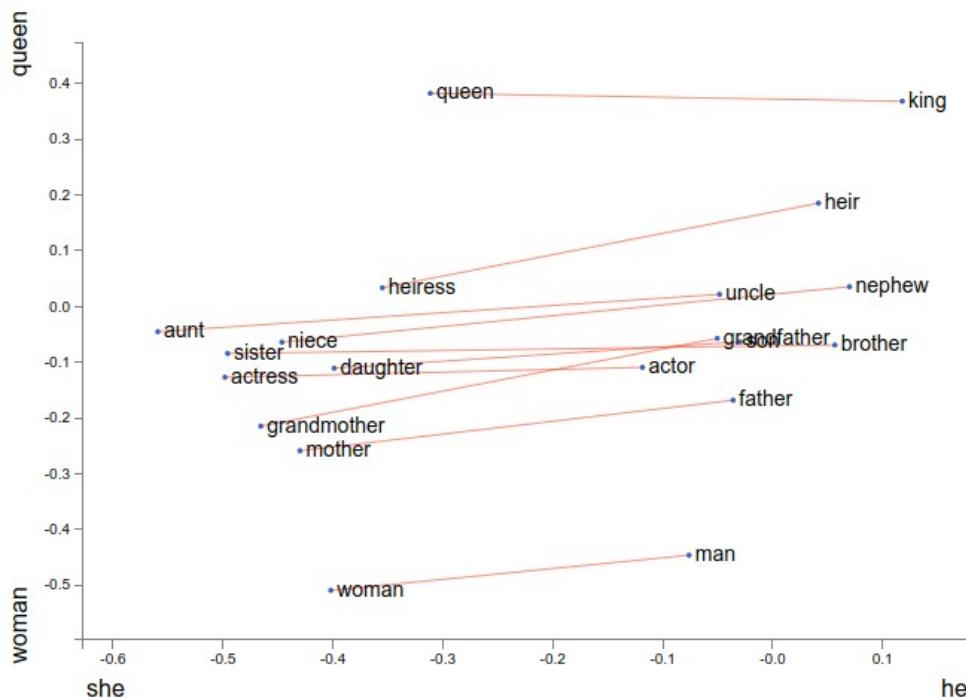
0.827 → 0.819

Сесть и во всем разобраться

The screenshot displays the Language Interpretability Tool (LIT) interface, which provides visualizations and explanations for NLP model decisions. The interface is divided into several sections:

- Top Left (Panel 1):** A 2D scatter plot titled "Embeddings" showing a cluster of points. It includes dropdown menus for "projector: UMAP" and "Embedding: sst2-tiny.cls_emb". A green line highlights a specific point from the Data Table below.
- Top Right (Panel 2):** A "Data Table" section showing a list of sentences with their corresponding labels (0 or 1). One row is highlighted in green, corresponding to the point in the Embedding plot. The table also includes a "label" column and a "sentence" column containing the text of the sentences.
- Right Side (Panels 3 and 4):** A "Datapoint Editor" section with two tabs: "TextSegment" and "CategoryLabel". The "TextSegment" tab shows a sentence and its production details. The "CategoryLabel" tab shows a label of "1".
- Middle Left (Panel 4):** A "Classification Results" section comparing "sst2-tiny - Reference" and "sst2-tiny - Main". Both show a confusion matrix with columns "Class", "Label", "Predicted", and "Score". The scores for predicted class 1 are 0.033 and 0.967 respectively.
- Middle Right (Panels 3 and 4):** "Attention" sections for "sst2-tiny - Reference" and "sst2-tiny - Main". These sections show attention matrices between tokens in the input sentence and the reference sentence. The matrices are visualized as grids where higher values are shown in darker shades of blue.

Король + Королева – Арифметика = ?



Explore word analogies

What do you want to see?

Gender analogies

Modify words

Type a new word...

Add

Type a new word...

Type a new word...

Add pair

X axis:

she

he

Y axis:

woman

queen

Change axes labels

Чеклист

| Capability | Min Func Test | INVariance | DIRectional |
|------------|------------------|------------|-------------|
| Vocabulary | Fail. rate=15.0% | 16.2% | C 34.6% |
| NER | 0.0% | B 20.8% | N/A |
| Negation | A 76.4% | N/A | N/A |
| ... | | | |

| Test case | Expected | Predicted | Pass? |
|---|----------|-----------|-------|
| A Testing Negation with MFT Labels: negative, positive, neutral | | | |
| A Testing Negation with MFT Labels: negative, positive, neutral | | | |
| Template: I {NEGATION} {POS_VERB} the {THING} . | | | |
| I can't say I recommend the food. | neg | pos | X |
| I didn't love the flight. | neg | neutral | X |
| ... | | | |
| Failure rate = 76.4% | | | |

Чеклист

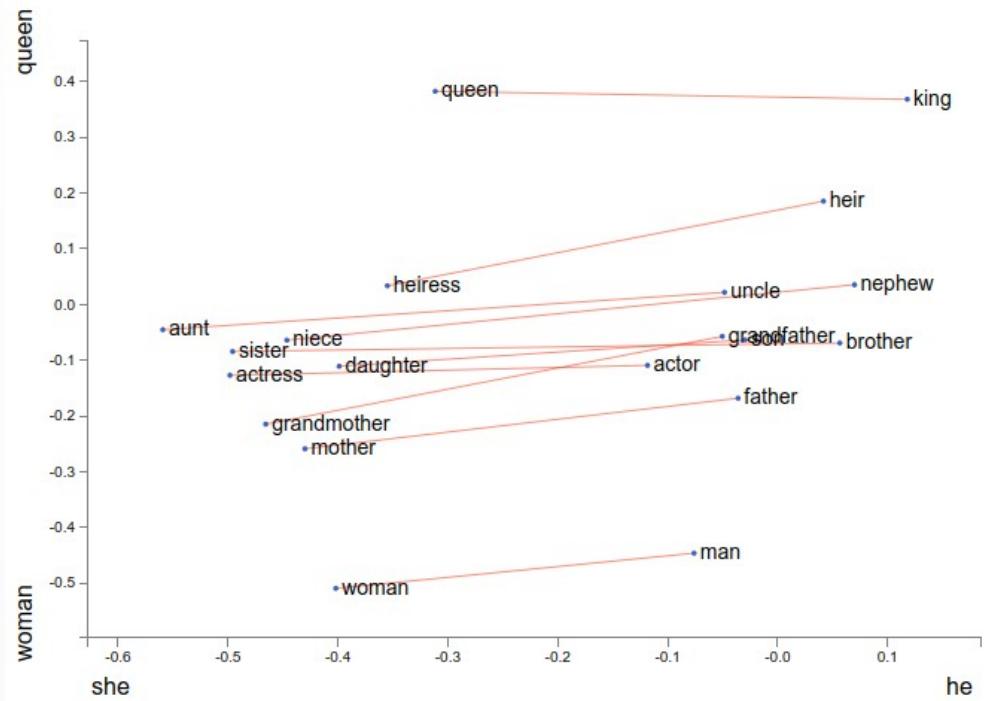
| B | Testing NER with <i>INV</i> | Same pred. (<i>inv</i>) after removals / additions |
|----------------------|--|--|
| | @AmericanAir thank you we got on a different flight to [Chicago → Dallas]. | |
| | @VirginAmerica I can't lose my luggage, moving to [Brazil → Turkey] soon, ugh. | |
| ... | | |
| Failure rate = 20.8% | | |
| C | Testing Vocabulary with <i>DIR</i> | Sentiment monotonic decreasing (↓) |
| | @AmericanAir service wasn't great. You are lame. | |
| | @JetBlue why won't YOU help them?! Ugh. I dread you. | |
| ... | | |
| Failure rate = 34.6% | | |

Секретный слайд №1

- Fine-tuning
 - Изменение весов
 - Интерполяция весов
- Экстраполяция весов

Секретный слайд №2

- Опорные фразы
- Граф



Вопросы

Слайды тут



dkolodezev



dmitry_kolodezev



promsoft



d_key

https://kolodezev.ru/download/omsk_nlp_2022.pdf