







DST-OFF

Задача 01

Сервис проверки поддельных новостей (fake news) в сфере технологий и инноваций



DST-OFF







НИКОЛАЙ ГАНИБАЕВ Капитан команды

Telegram: nganibaev Email: ganibaev@gmail.com Тел. +7 903 851 5919

Функции: аналитика (фреймворки fasttext, natasha), формирование базы данных новостей



АЛЕКСАНДР ГАНИБАЕВ Аналитик, разработчик

Telegram: aganibaev Email: <u>acleriot@gmail.com</u> Тел. +7 960 111 4192

Функции: программирование apiсервиса (фреймворк fastapi) и работа с github, автоматическое разворачивание в прод



АЛЕКСАНДР ГАНИБАЕВ Аналитик, разработчик

Telegram: valthazari Email: <u>vasif.faradzhov@yandex.ru</u> Тел. +7 906 062 9089

Функции: описание концепции презентация и дизайн



Подход к решению задачи*







Работает постоянно в фоне. Требует список авторитетных источников и алгоритм разбора страниц для сбора текста новостей. В целях хакатона база данных новостей будет собрана вручную с mos.ru.



Работает по запросу На входе текст, на выходе оценка достоверности новости (1-100 баллов)

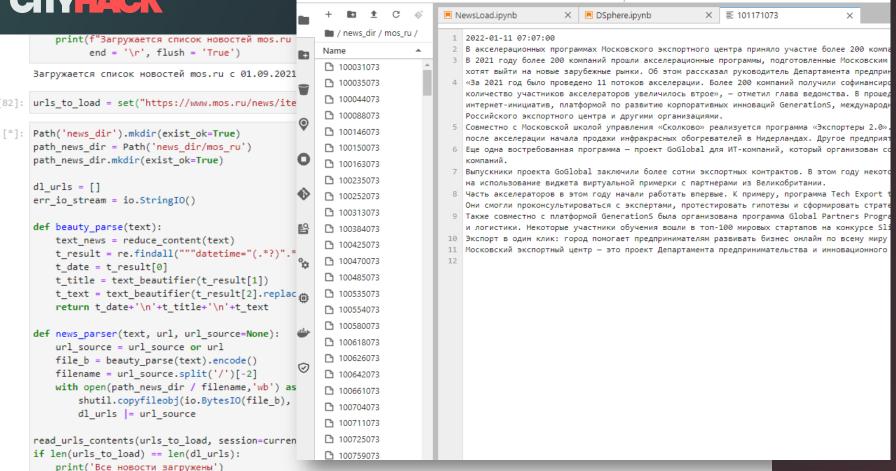


для демонстрации работы АРІ

*Базовый стек: Python, Natasha, FastText, FastApi







File Edit View Run Kernel Git Snippets Tabs



df_news[0] = """Москва заняла первое место среди европейских городов в рейтинге инноваций, помогающих в борьбе с COVID-19 В мире Москва занимает третье место, уступая лишь Нью-Йорку и Сан-Франциско.

Москва признана первой среди европейских городов в рейтинге инноваций, помогающих в формировании устойчивости коронавирусу Среди мировых мегаполисов российская столица занимает третью строчку — после Сан-Франциско и Нью-Йорка. Пятерку замыкают В

добиться высоких показателей Москве помогло п Среди них алгоритмы компьютерного зрения на о Еще одно инновационное решение — облачная пла Способствовали высоким результатам и технолог Эксперты агентства StartupBlink оценивали при Московский опыт

В борьбе с коронавирусом Москва отказалась от Московская система здравоохраненя за время па Столица поддерживает бизнес, выделяя субсидии Как составляется рейтинг

Рейтинг составляется на базе глобальной карты Алгоритм рейтинга учитывает количество и тип

df_news[1] = """Москва заняла первое место ср
В мире Москва занимает третье место, уступая .
Москва признана первой среди европейских горо,
Среди мировых мегаполисов российская столица :
Добиться высоких показателей Москве помогло п

```
from transformers import pipeline

p = pipeline(
    task='zero-shot-classification',
    model='cointegrated/rubert-base-cased-nli-threeway'
)

p(
    sequences=df_news[8],
    candidate_labels="Xopowo, Noxo",
    hypothesis_template="{}."
)
```

{'sequence': 'Проектный офис ФЭСН РАНХиГС завершает 30 проектов. И начинает новые.\пЕжегодно в мае-июне Проектный офис Факультета экономических и зультаты заказчикам — российским и международным компаниям.\пСреди компаний-заказчиков были крупные российские компании и организации: Сбер, РЖД, ого, холдинг САВООВ ФУДС, Альфа-Банк, ВТБ, Дом РФ, Очаково, Сегежа-Гряпп а также представительства зарубежных компаний: ВМW, DeLonghi, L'Oreal, РТ ФЭСН является крупнейшим университетским проектным центром не только по количеству и сложности проектов, и не только по количеству студентов-участ ы. Так, еще пять лет назад ФЭСН вовлек в проектную деятельность университеты пяти стран Европы и стал разрабатывать с ними полугодичные проекты дл делались усилиями проектных групп нексольких стран; например, в проекте от генерального директора ВНW Russi по «Разработке маркетинговой стратеги ии» бакалавры ФЭСН работали вместе с магистрами из Германии и Бельгии. А в одном из проектов участвовали команды из Германии, Бельгии, Италии, 4 п , 'labels': ['Хорошо', 'Плохо'], 'scores': [0.7546549439480237, 0.2453451007604599]}

Среди них алгоритмы компьютерного зрения на основе искусственного интеллекта. Это методика уже помогла рентгенологам проак Еще одно инновационное решение — облачная платформа, которая объединяет пациентов, врачей, медицинские организации, страхо Способствовали высоким результатам и технологии, которые помогают адаптировать жизнь горожан во время пандемии. Это проект Эксперты агентства StartupBlink оценивали принятые в Москве меры с точки зрения эпидемиологических показателей и влияния н

df_news[2] = """Москва стала лидером в Европе в рейтинге инноваций, помогающих в борьбе с COVID-19

В мире российская столица заняла третье место, обогнав Лондон и Барселону.

Москва заняла первое место среди европейских городов в рейтинге инноваций, помогающих в борьбе с COVID-19, опередив Лондов ""В российской столице применяются почти 160 передовых решений для борьбы с распространением коронавируса. Среди них алгор Рейтинг составляется на базе глобальной карты инновационных решений по борьбе с коронавирусом и оценивает около 100 ведущи

df news[3] = """Москва заняла первое место в Европе по инновациям в борьбе с COVID-19

Москва обошла европейские столицы в рейтинге инноваций по устойчивости к COVID-19, опередив Лондон и Барселону, сообщается Российская столица также заняла третье место среди мировых мегаполисов. В пятерку лидеров вошли Бостон и Лондон.

Занять лидирующие позиции в рейтинге Москве помогли около 50 передовых решений, которые применяются для борьбы с распростр Одно из таких решений - алгоритмы компьютерного зрения на основе искусственного интеллекта, которые уже помогли рентгеноло Также высоким результатам способствовали технологии, помогающие адаптировать жизнь москвичей во время пандемии. Среди них







```
from transformers import pipeline
p = pipeline(
   task='zero-shot-classification',
   model='cointegrated/rubert-base-cased-nli-threeway'
)
p(
  sequences=df_news[8],
  candidate_labels="Xopowo, Плохо",
  hypothesis_template="{}."
)
```

{'sequence': 'Проектный офис ФЭСН РАНХиГС завершает 30 проектов. И начинает новые.\nЕжегодно в мае-июне Проектный офис Факультета экономических и зультаты заказчикам – российским и международным компаниям.\nСреди компаний-заказчиков были крупные российские компании и организации: Сбер, РЖД, ого, холдинг САВООВ ФУДС, Альфа-Банк, ВТБ, Дом РФ, Очаково, Сегежа-Групп а также представительства зарубежных компаний: ВМW, DeLonghi, L'Oreal, Pf ФЭСН является крупнейшим университетским проектным центром не только по количеству и сложности проектов, и не только по количеству студентов-участ ы. Так, еще пять лет назад ФЭСН вовлек в проектную деятельность университеты пяти стран Европы и стал разрабатывать с ними полугодичные проекты дл делались усилиями проектных групп нескольких стран; например, в проекте от генерального директора ВМW Russia по «Разработке маркетинговой стратеги ии» бакалавры ФЭСН работали вместе с магистрами из Германии и Бельгии. А в одном из проектов участвовали команды из Германии, Бельгии, Италии, 4 п , 'labels': ['Хорошо', 'Плохо'], , 'scores': [0.7546549439430237, 0.2453451007604599]}







```
[236]: from sentence_transformers import SentenceTransformer, util
       model = SentenceTransformer('sentence-transformers/LaBSE')
[241]: import numpy as np
       # Two lists of sentences
       sentences1 = list(df_news.values())
       sentences2 = list(df news.values())
       #Compute embedding for both lists
       embeddings1 = model.encode(sentences1, convert to tensor=True)
       embeddings2 = model.encode(sentences2, convert to tensor=True)
       #Compute cosine-similarits
       cosine scores = util.cos sim(embeddings1, embeddings2).numpy().mean()
       cosine_scores
       0.7916351
[244]: cosine_scores = util.cos_sim(embeddings1, embeddings2)
       #Output the pairs with their score
       for i in range(1):
           for j in range(len(sentences2)):
               print("{} \t\t {} \t\t Score: {:.4f}".format(i, j, cosine_scores[i][j]))
                                        Score: 1.0000
                                        Score: 1.0000
                                        Score: 0.9246
                                        Score: 0.9011
                                        Score: 0.8411
                                        Score: 0.9015
                                        Score: 0.9513
                                        Score: 0.7405
                                        Score: 0.5089
```







```
import torch
from transformers import AutoTokenizer, AutoModel
tokenizer = AutoTokenizer.from pretrained("cointegrated/LaBSE-en-ru")
model = AutoModel.from pretrained("cointegrated/LaBSE-en-ru")
sentences = list(df news.values())
encoded_input = tokenizer(sentences, padding=True, truncation=True, return_tensors='pt')
with torch.no grad():
    model output = model(**encoded input)
embeddings = model output.pooler output
embeddings = torch.nn.functional.normalize(embeddings)
#Compute cosine-similarits
cosine scores = util.cos sim(embeddings, embeddings).numpy().mean()
cosine_scores
Some weights of the model checkpoint at cointegrated/LaBSE-en-ru were not used when initializing BertModel: ['cls.predictions.transform.dense.bias
ls.predictions.transform.dense.weight', 'cls.seq relationship.weight', 'cls.predictions.bias', 'cls.seq relationship.bias', 'cls.predictions.trans
- This IS expected if you are initializing BertModel from the checkpoint of a model trained on another task or with another architecture (e.g. ini
ng model).
- This IS NOT expected if you are initializing BertModel from the checkpoint of a model that you expect to be exactly identical (initializing a Be
model).
0.7849529
cosine scores = util.cos sim(embeddings, embeddings)
#Output the pairs with their score
for i in range(1):
    for j in range(len(sentences)):
        print("{} \t\t {} \t\t Score: {:.4f}".format(i, j, cosine scores[i][j]))
                                 Score: 1,0000
a
                                 Score: 0.9433
                                 Score: 0.8958
0
                                 Score: 0.8647
                                 Score: 0.8255
                                 Score: 0.9211
                                 Score: 0.9210
                                 Score: 0.6971
                                 Score: 0.5612
```



```
1 df_news[1].splitlines()[1]
Ввод [11]:
            executed in 22ms, finished 21:08:46 2022-06-11
 Out[11]: 'В мире Москва занимает третье место, уступая лишь Нью-Йорку и Сан-Франциско.'
Ввод [28]:
               1 import pandas as pd
               2 from dostoevsky.tokenization import RegexTokenizer
               3 from dostoevsky.models import FastTextSocialNetworkModel
                4 import fasttext
                6 fasttext.FastText.eprint = lambda x: None
                7 tokenizer = RegexTokenizer()
               9 model = FastTextSocialNetworkModel(tokenizer=tokenizer)
              10
              11 messages = df_news.values()
              12 # messages = df topics
              13 # messages = [x.splitlines()[1] for x in df_news.values()]
              14
              15 results = model.predict(messages)
              16
              17 | p = pd.DataFrame(messages, columns=['text'])
              18 p['positive']=''
              19 p['negative']=''
              20 p['rezult']=''
              21
              22 for id, sentiment in enumerate(results):
                       p['positive'][id] = f"{sentiment['positive']:.2f}"
                       p['negative'][id] = f"{-sentiment['negative']:.2f}"
              25
                       p['rezult'][id] = f"{sentiment['positive']-sentiment['negative']:.2f}"
              26
              27 p
            executed in 851ms, finished 21:25:32 2022-08-11
  Out[28]:
                                                         text positive negative rezult
                   Москва заняла первое место среди европейских г...
                                                                        -0.19 -0.12
                                                                0.06
                   Москва заняла первое место среди европейских г...
                                                                0.08
                                                                        -0.18 -0.10
                   Москва стала лидером в Европе в рейтинге иннов...
                                                                        -0.17 -0.08
                   Москва заняла пеовое место в Европе по инновац...
                                                                        -0.15 -0.07
                   Москва заняла первое место в Европе по инновац...
                                                                        -0.17 -0.10
                   Москва стала первой в Европе среди городов с и...
                                                                        -0.22 -0.14
                                                                0.08
```

-0.16 -0.08

-0.22 -0.14

-0.20 -0.13

Москва заняла восьмое место среди европейских ...

Бостон занял первое место среди европейских го...

8 Проектный офис ФЭСН РАНХиГС завершает 30 проек...







```
Ввод [29]:
                  1 from autocorrect import Speller
                  2 spell = Speller(lang='ru', fast=True, only replacements=True)
              executed in 8.94s, finished 21:25:42 2022-06-11
Ввод [30]: -
                  1 # Проходим спелл-чеккером по всем текстам
                  2 p["spell correct"] = p['text'].apply(lambda x: spell(str(x)))
                  3 D
              executed in 119ms, finished 21:25:42 2022-08-11
Ввод [36]:
              executed in 44ms, finished 21:25:53 2022-08-11
  Out[36]:
                                                                text positive negative rezult
                                                                                                                                      spell correct err count
                      Москва заняла первое место среди европейских г...
                                                                                  -0.19 -0.12
                                                                                                    Москва заняла первое место среди европейских г...
                                                                         0.06
                      Москва заняла первое место среди европейских г...
                                                                         0.08
                                                                                  -0.18
                                                                                        -0.10
                                                                                                    Москва заняла первое место среди европейских г...
                     Москва стала лидером в Европе в рейтинге иннов...
                                                                                                   Москва стала лидером в Европе в рейтинге иннов...
              2
                                                                         0.09
                                                                                  -0.17 -0.08
                     Москва заняла первое место в Европе по инновац...
                                                                                  -0.15 -0.07
                                                                                                   Москва заняла первое место в Европе по инновац...
                                                                         0.08
                     Москва заняла первое место в Европе по инновац...
                                                                                  -0.17 -0.10
                                                                                                   Москва заняла первое место в Европе по инновац...
                                                                         0.07
                      Москва стала первой в Европе среди городов с и...
                                                                                  -0.22 -0.14
                                                                                                    Москва стала первой в Европе среди городов с и...
                                                                         0.08
                                                                                                                                                            0
                                                                                                   Москва заняла восьмое место среди европейских ...
                     Москва заняла восьмое место среди европейских ...
                                                                         80.0
                                                                                  -0.16 -0.08
                      Бостон занял первое место среди европейских го...
                                                                                  -0.22 -0.14
                                                                                                    Бостон занял первое место среди европейских го...
                                                                         0.08
              8 Проектный офис ФЭСН РАНХиГС завершает 30 проек...
                                                                                  -0.20 -0.13 Проектный офис ФЭСН РАНХиГС завершает 30 проек...
                                                                         0.07
```



```
1 from natasha import (
Ввод [102]: "
                       Segmenter,
                       MorphVocab,
                      NewsEmbedding,
                      NewsMorphTagger,
                      NewsSyntaxParser,
                      NewsNERTagger,
                8
                       Doc,
                9 )
               10
               11 emb = NewsEmbedding()
               12 | segmenter = Segmenter()
               13 morph_vocab = MorphVocab()
               14 ner_tagger = NewsNERTagger(emb)
               15 # morph_tagger = NewsMorphTagger(emb)
               16 # syntax_parser = NewsSyntaxParser(emb)
               17
             executed in 2.44s, finished 05:47:48 2022-06-12
Ввод [114]: * 1 # text = ''.join(m.lemmatize(df_news[0]))
                2 text = df_news[0]
                3 doc = Doc(text)
                4 | doc.segment(segmenter)
                5 doc.tag_ner(ner_tagger)
                6 # doc.tag_morph(morph_tagger)
                7 # doc.parse_syntax(syntax_parser)
             executed in 104ms, finished 05:58:08 2022-06-12
               1 for span in doc.spans:
Ввод [115]: -
                      span.normalize(morph_vocab)
                   NER_extract = dict(pd.DataFrame([(x.type, x.normal) for x in doc.spans], columns=['type', 'title'])
                                      .groupby('type').title.unique())
               7 if 'ORG' in NER_extract:
                      print('Найдены организации: ', ', '.join(NER_extract['ORG']))
               9 if 'LOC' in NER_extract:
                      print('\nНайдены места: ', ', '.join(NER_extract['LOC']))
            + 11 if 'PER' in NER_extract:
                      print('\nНайдены люди: ', ', '.join(NER_extract['PER']))
             executed in 22ms, finished 05:58:12 2022-06-12
             Найдены организации: StartupBlink
             Найдены места: Москва, Нью-Йорку, Сан-Франциско, Лондон, Барселону, Нью-Йорка, Бостон, Москве
```













Заголовок

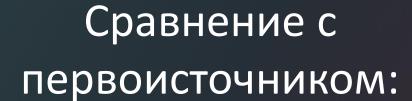
- кликбейт/желтизна заголовка
- тональность заголовка
- близость семантики заголовка содержанию новости



Новости

- тональность новости
- Охват /распространенность новости в авторитетных источниках
- проверка на опечатки

















Базовые метрики:





Формируем для каждой новости набор метрик/оценок:

семантическая близость заголовков

01



0

04

близость по тональности заголовков

семантическая близость первых абзацев (обычно суть новости)

02





05

близость по тональности первых абзацев

семантическая близость текстов без первого абзаца

03





06

близость по тональности текстов без первого абзаца



Базовые метрики:







- учитываются новости с семантической схожестью более _80_% за _2_ дня

Тональность новости

- стандартная новость нейтральна



Дополнительные метрики*







^{*}рассчитываются после выявления в тексте новости и первоисточника именованных сущностей (NER) - персоналий, организаций, мест, числительных























ДЕМО





