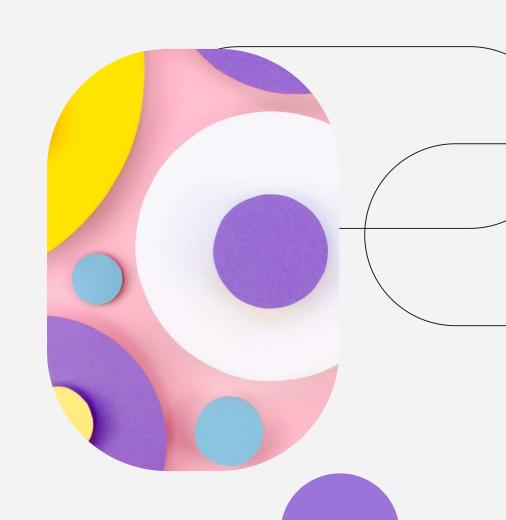


FastAPI

Application Programming Interface

O1Jupyter



Подключение

NLP

Подключение библиотек

[1]: import warnings
 warnings.filterwarnings('ignore')

[2]: import pandas as pd import nltk import matplotlib.pyplot as plt

 $from \ sklearn.feature_extraction.text \ import \ CountVectorizer \\ from \ sklearn.feature_extraction.text \ import \ TfidfVectorizer \\$

from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize

Загрузка данных

[3]: df = pd.read_csv('dataset03.csv')

[4]: **df**

[4]:		namecompany	description	rating	activity	date_publish	text_article	prep_text	tokenize_text	text_stem
	0	Timeweb Cloud	Облачная платформа для разработчиков и бизнеса	1572.77	Веб-разработка, Домены и хостинг, Веб-сервисы	2023-03-27	В своём блоге я рассмотрел устройство и назнач	в своём блоге я рассмотрел устройство и назнач	своём блоге рассмотрел устройство назначение к	сво блог рассмотрел устройств назначен контрол
	1	OTUS	Цифровые навыки от ведущих экспертов	795.43	Консалтинг и поддержка, Рекрутинг и НR, Произв	2023-03-27	Автор статьи: Роман КозловРуководитель курса В	автор статьи роман козловруководитель курса bi	автор статьи роман козловруководитель курса bi	автор стат рома козловруководител курс bi анал
	2	Онлайн Патент	Ваш личный патентный офис	220.74	Консалтинг и поддержка, Веб- сервисы	2023-03-27	Рынок биологически активных добавок(БАД) с каж	рынок биологически активных добавок бад с кажд	рынок биологически активных добавок бад каждым	рынок биологическ активн добавок бад кажд год

Векторизация

Векторизация текстовых данных

```
[4]: # sazpyska cmucka cmon-cnoß

nltk.download('stopwords')

[nltk_data] Downloading package stopwords to

[nltk_data] C:\Users\HOME\AppData\Roaming\nltk_data...

[nltk_data] Package stopwords is already up-to-date!

[4]: True

[5]: # sazpyska cmon-cnoß на русском языке

russian_stopwords = stopwords\words("russian")
```

TF-IDF

TfidfVectorizer - это инструмент векторизации текста в машинном обучении, который используется для преобразования текстовых данных в числовые признаки. Он основан на понятии TF-IDF (Term Frequency-Inverse Document Frequency), который позволяет оценить важность слова в документе.

TF-IDF вычисляется для каждого слова в документе путем умножения его частоты встречаемости (ТF) на обратную частоту документа (IDF). ТF оценивает, насколько часто слово появляется в документе, в то время как IDF оценивает, насколько информативно слово является путем снижения веса слова, которое часто встречается во всей коллекции документов.

TfidfVectorizer векторизует текст, создавая матрицу, в которой строки представляют документы, а столбцы - каждое уникальное слово или термин. Значения в ячейках этой матрицы представляют собой значения ТF-IDF для каждого термина в каждом документе. Таким образом, каждый документ представлен в виде вектора ТF-IDF значений, который может быть использован в алгоритмах машинного обучения.

- [7]: tfidf_matrix.shape
- [7]: (438, 10000)
- [8]: print(tfidf_matrix)
 - (0, 8947) 0.017091908228614894 (0, 2203) 0.018917742841662737 (0, 1831) 0.01187840743372064

Кластеризация

```
[12]: pip install threadpoolctl==3.1.0
       Requirement already satisfied: threadpoolctl==3.1.0 in c:\users\home\anaconda3\lib\site-packages (3.1.0)
       Note: you may need to restart the kernel to use updated packages.
      Метод к-средних - KMeans
[13]: from sklearn.cluster import KMeans
[14]: # ввод количества кластеров
       num clusters = 5
       # обучение модели к-средних
       km = KMeans(n clusters=num clusters)
       idx = km.fit(tfidf_matrix)
      # список с значениями кластеров
       clusters = km.labels .tolist()
[15]: frame = pd.DataFrame()
       #k-means
      out = { 'text': df['text_article'], 'cluster': clusters }
      frame = pd.DataFrame(out, columns = ['text', 'cluster'])
[16]: frame
                                                     text cluster
         В своём блоге я рассмотрел устройство и назнач...
        1 Автор статьи: Роман КозловРуководитель курса В...
         2 Рынок биологически активных добавок(БАД) с каж...
        3 Автор статьи: Александр КолесниковВирусный ана...
         4 Всем привет! Меня зовут Артём Пузанков, я DevS...
```



Анализ кластеров

Облако слов

Облако слов (Word Cloud) - это визуализация, которая отображает наиболее часто встречающиеся слова в наборе текстовых данных. В облаке слов размер и положение каждого слова зависят от его частоты в тексте: чем чаще слово встречается, тем больше его размер и тем ближе оно к центру.

Облака слов часто применяются для визуализации результатов анализа текста, выявления ключевых тем и идей, а также для быстрого восприятия наиболее релевантных слов в текстовых данных.

[21]: !pip install WordCloud

Анализ кластеров

```
from wordcloud import (
    WordCloud.
    STOPWORDS
# получение текстовой строки из списка слов
def str_corpus(corpus):
   str corpus = ''
   for i in corpus:
       str corpus += ' ' + i
    str corpus = str corpus.strip()
    return str corpus
# получение списка всех слов в корпусе
def get corpus(data):
    corpus = []
    for phrase in data:
        for word in phrase.split():
            corpus.append(word)
    return corpus
# получение облака слов
def get wordCloud(corpus):
    wordCloud = WordCloud(background color='white',
                              stopwords=russian stopwords,
                              width=3000,
                              height=2500,
                              max words=200,
                              random state=42
                         ).generate(str corpus(corpus))
    return wordCloud
```

```
num = 0
for i in frame[frame['cluster']==0]['text_lemm3']:
    num += 1
    corpus = get_corpus(i)
    procWordCloud = get_wordCloud(corpus)

fig = plt.figure(figsize=(17, 6))
    plt.subplot(1, 2, 1)
    plt.imshow(procWordCloud)
    plt.axis('off')
    plt.show()
    if num == 5:
        break
```

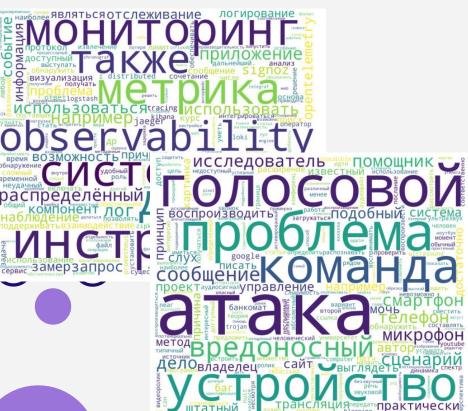
ВВИХОДНО И ВОСПРИЯТЬ В ОБОРМЕНТ В ОБОРМЕНТ

Статьи, связанные с исследованием каких либо изобретений, моделей.





Статьи, связанные с системами, их использованием. Рассказываются какие либо проблемы и как их решить





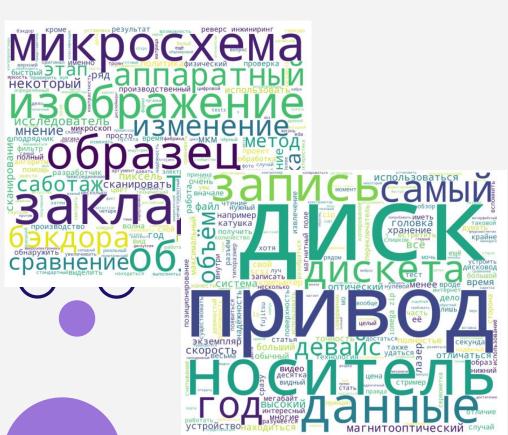


Статьи, связанные с обзором данных и с разработкой проектов для анализа данных





Статьи, связанные с устройствами, контролерами такие как ноутбук, смартфон, привод, микрофон. Исследования данных устройств.







Статья содержит много кода либо англоязычных слов.



Классификация

```
from sklearn.model selection import train test split
from sklearn.metrics import accuracy score
from sklearn.linear model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.neighbors import KNeighborsClassifier
# разделение на две выборки датафрейма с кластерами
X_train, X_test, y_train, y_test = train_test_split(frame['text'], frame['cluster'], test_size=0.25, random_state=42)
# преобразование текстовых данных в числовые признаки с помощью TF-IDF
vectorizer = TfidfVectorizer(max_features=10000, ngram_range=(1, 2), stop_words=russian_stopwords)
X train tfidf = vectorizer.fit transform(X train)
X test tfidf = vectorizer.transform(X test)
LogisticRegression
# обучение модели логистической регрессии
model lr = LogisticRegression()
model lr.fit(X train tfidf, y train)
# предсказание на тестовом наборе
y pred = model lr.predict(X test tfidf)
# оценка качества модели
accuracy lr = accuracy score(y test, y pred)
print("Accuracy:", accuracy_lr)
Accuracy: 0.7636363636363637
```

RandomForestClassifier

```
[30]: # обучение модели случайного леса
      model rf = RandomForestClassifier()
      model rf.fit(X train tfidf, y train)
      # предсказание на тестовом наборе
      y_pred = model_rf.predict(X_test_tfidf)
       # оценка качества модели
      accuracy_rf = accuracy_score(y_test, y_pred)
      print("Accuracy:", accuracy_rf)
      Accuracy: 0.7363636363636363
      KNeighborsClassifier
[32]: # обучение модели k-ближайших соседей
      model knn = KNeighborsClassifier()
      model knn.fit(X train tfidf, y train)
      # предсказание на тестовом наборе
      y pred = model knn.predict(X test tfidf)
       # оценка качества модели
      accuracy knn = accuracy score(y test, y pred)
      print("Accuracy:", accuracy_knn)
      Accuracy: 0.72727272727273
[33]: print("Точность логистической реграссии -", accuracy_lr)
      print("Точность случайного леса -", accuracy_rf)
      print("Точность дерева решений -", accuracy dt)
      print("Точность k-ближайших соседей -", accuracy knn)
      Точность логистической реграссии - 0.7636363636363637
      Точность случайного леса - 0.7363636363636363
      Точность дерева решений - 0.51818181818182
       Точность k-ближайших соседей - 0.72727272727273
```

Сохранение моделей

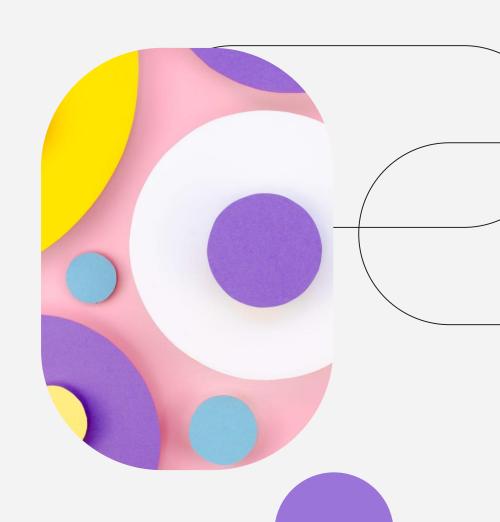
```
import pickle
with open('model_rf.pkl', 'wb') as file:
    pickle.dump(model_rf, file)

with open('vectorizer.pkl', 'wb') as f:
    pickle.dump(vectorizer, f)

frame.to_csv("Data_NLP7.csv")
```



FastAPI



Kmo mακοῦ FastAPI?

API (Application Programming Interface) — это набор способов и правил, по которым различные программы общаются между собой и обмениваются данными.

FastAPI — веб-фреймворк для создания API, написанный на Python. Один из самых быстрых[2] и популярных (после Django и Flask)[3] веб-фреймворков, написанных на Python (на 2023 год).



Подключение FastAPI

```
(venv) C:\Users\Admin\PycharmProjects\api_text_clust
pip install fastapi[all]
```

```
from fastapi import FastAPI
from fastapi.middleware.cors import CORSMiddleware
```

Библиотеки и создание АРІ

Middleware (промежуточное или связующее программное обеспечение) — это фрагмент кода в конвейере приложения, используемый для обработки запросов и ответов.

```
from fastapi import FastAPI

from fastapi.middleware.cors import CORSMiddleware
import pickle
import numpy as np
import pandas as pd
from pydantic import BaseModel
import string
import re
import nltk
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize

import pymorphy2
```

Загрузка моделей

```
# загрузка моделей
with open('model_rf.pkl', 'rb') as file:
    model = pickle.load(file)

with open('vectorizer.pkl', 'rb') as file:
    vectorizer = pickle.load(file)

with open('text2.txt', 'r', encoding='utf-8') as file:
    text2 = file.read()
```

```
def fun_punctuation_text(text):
   text = ''.join([ch for ch in text if ch not in string.punctuation])
   text = re.sub(r'\s+', ' ', text, flags=re.I)
   text = re.sub('[a-z]', '', text, flags=re.I)
   return text
def fun_tokenize_text(text):
   nltk.download('stopwords')
   stopword_eng = stopwords.words("english")
   stopword_ru = nltk.corpus.stopwords.words('russian')
   nltk.download('word tokenize')
   t = word_tokenize(text)
   text = [token for token in t if token not in stopword_eng and token not in stopword_ru]
   return text
```

Функции по обработке текста

```
# предобработка теста.

def fun_pred_text(text):
    text = fun_punctuation_text(text)
    text = fun_tokenize_text(text)
    text = fun_lemmatizing_text(text)
    return text
```

```
# лемматизация текста

def fun_lemmatizing_text(text):
    morph = pymorphy2.MorphAnalyzer()
    res = list()

for word in text:
    p = morph.parse(word)[0]
    res.append(p.normal_form)
    return res
```

Функции для предсказания

```
# TF-IDF

def fun_TfidfVectorizer(text):
    text = fun_pred_text(text)

df = pd.DataFrame({'text': str_corpus(text)}, index=range(1))
    new_row = pd.DataFrame({'text': text2}, index=range(1))
    df = pd.concat([df, new_row])

tfidf_matrix = vectorizer.transform(df['text'])
return tfidf_matrix
```

```
# преобразование списка в текст

def str_corpus(corpus):

   str_corpus = ''

   for i in corpus:

       str_corpus += ' ' + i

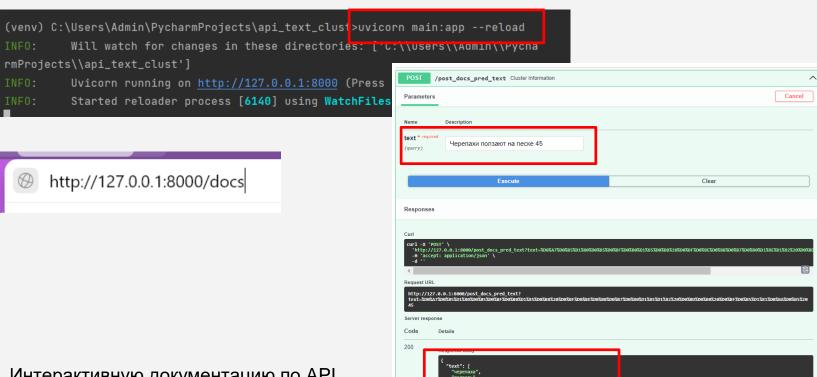
       str_corpus = str_corpus.strip()

   return str_corpus
```

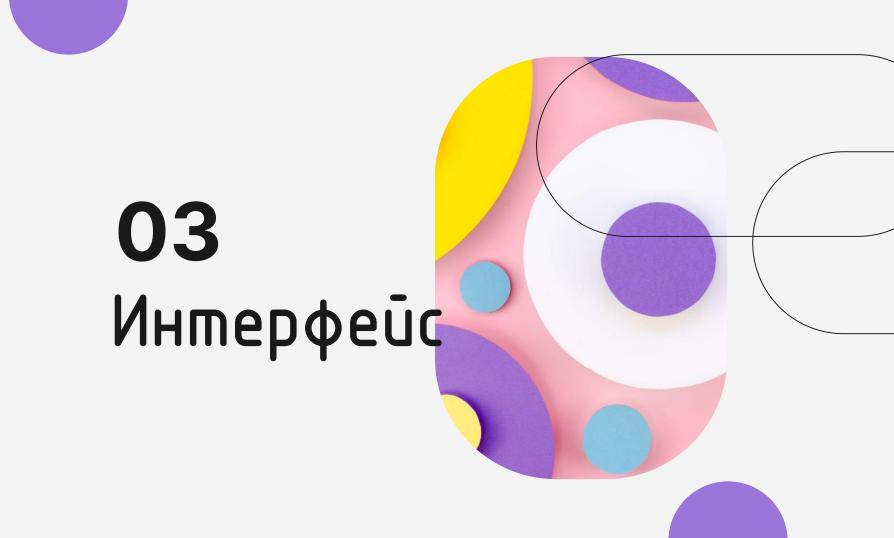
Получение и отправка данных

```
@app.post("/post_docs_pred_text")
|def cluster_information(text):
    return {'text': fun_pred_text(text)}
class Item(BaseModel):
    text: str
# метод для подключения к интерфейсу
@app.post("/post_cluster")
def post_pred_text(item: Item):
    return {'cluster': fun_predict_cluster(item.text)}
```

3anyck API



Интерактивную документацию по API



PyQt

Предсказание кластера \times Современные бизнес приложения имеют сложную структуру, состоят из множества независимых компонентов, часто распределенных по разным узлам, контейнерам или виртуальным машинам. В связи с этим, поиск неисправностей в таких серьезных приложениях становится тоже непростой задачей, зачастую превращающейся в нетривиальный квест. При этом не стоит забывать, каждая минута простоя приложения в продакшене стоит денег, поэтому выявлять причины сбоев необходимо как можно быстрее. В этой и последующих статьях мы поговорим о том, как производить мониторинг работы приложений, осуществлять сбор событий и трейсинг для решения конкретных проблем в работе приложения. Эти статьи ориентированы на специалистов, занимающихся администрированием Linux, хотя о Windows мы тоже будем иногда упоминать.\xa0\xa0 Эта стать Предсказать Предсказанный кластер: 2 – статьи, связанные с обзором данных и с разработкой проектов для анализа данных.

Streamlit

Гадалка онлайн

Введите свой текст тут

Черепаха – воплощение мудрости и спокойствия в мире животных. Ее панцирь, словно древний щит, защищает от жизненных невзгод, символизируя стойкость и выносливость. Медленные движения черепахи скрывают в себе глубокую философию – она учит нас ценить каждый момент, не спеша идти к своей цели. Ее глаза, полные мудрости и таинственности, словно окно в древние времена, напоминают о важности уважения к прошлому и умении жить в гармонии с окружающим миром. Черепаха – символ долголетия и стойкости, напоминание о том, что и в медленном темпе можно достичь великих высот.

Предсказать

Предсказанный кластер

Кластер 3 – статьи, связанные с работниками, участниками проектов или компаний.

PyQt umnopm

```
import sys
from PyQt5.QtWidgets import QApplication, QMainWindow, QVBoxLayout, \
    QHBoxLayout, QLineEdit, QPushButton, QLabel, \
    QWidget, QTextEdit
import requests
```

Создание окна

```
class PredictionApp(QMainWindow):
       super().__init__()
       self.setWindowTitle("Предсказание кластера")
       self.setGeometry(100, 100, 600, 400)
       main_layout = QVBoxLayout()
       self.input_field = QTextEdit()
       self.input_field.setPlaceholderText("Ввод текста")
       main_layout.addWidget(self.input_field)
       self.predict_button = QPushButton("Предсказать")
       self.predict_button.clicked.connect(self.predict_cluster)
       main_layout.addWidget(self.predict_button)
```

```
# Cosdanue лейбла для вывода ответа
self.output_label = QLabel("Нажмите предсказать для отображения кластера")
self.output_label.setWordWrap(True)
main_layout.addWidget(self.output_label)

# Cosdanue главного лейаута для приложения
central_widget = QWidget()
central_widget.setLayout(main_layout)
self.setCentralWidget(central_widget)
```

Предсказание кластера

```
def predict_cluster(self):
    # Получаем текст из нашего поля
    text = self.input_field.toPlainText()
    data = {
        "text": text
   url = "http://127.0.0.1:8000/post_cluster"
    response = requests.post(url, json=data)
    result = response.json()
    clust = result.get("cluster")
    self.output_label.setText(f"Предсказанный кластер: {clust}")
```

Ну и главный метод который создаёт само приложение

```
# Ну и главный метод который создаёт само приложение

pif __name__ == "__main__":
    app = QApplication(sys.argv)
    prediction_app = PredictionApp()
    prediction_app.show()
    sys.exit(app.exec_())
```

Streamlit

```
# Загрузка библиотек
import streamlit as st
import requests
import os
# Обнуление статусов прокси для корректного подключения
os.environ['HTTP_PROXY'] = ''
os.environ['HTTPS_PROXY'] = ''
```

```
(venv) C:\Users\Admin\PycharmProjects\api_text_clust>
pip install streamlit
```

всеѕо Создание

```
if __name__ == "__main__":
main(<u>)</u>
```

```
def main():
    input_text = st.text_area("Введите свой текст тут", height=200)
        if input_text == "":
            st.write(f"Введите текст")
           data = {
                "text": input_text
           url = "http://127.0.0.1:8000/post_cluster"
           response = requests.post(url, json=data)
           result = response.json()
            clust = result.qet("cluster")
            st.markdown(f"""
                    111111
            st.write(f"Кластер {clust}")
```

3anyck streamlit

```
(venv) C:\Users\Admin\PycharmProjects\api_text_clust;
streamlit run app_streamlit.py

You can now view your Streamlit app in your browser.

Local URL: http://localhost:8501
Network URL: http://192.168.98.249:8501
```

PyQt

Предсказание кластера \times Современные бизнес приложения имеют сложную структуру, состоят из множества независимых компонентов, часто распределенных по разным узлам, контейнерам или виртуальным машинам. В связи с этим, поиск неисправностей в таких серьезных приложениях становится тоже непростой задачей, зачастую превращающейся в нетривиальный квест. При этом не стоит забывать, каждая минута простоя приложения в продакшене стоит денег, поэтому выявлять причины сбоев необходимо как можно быстрее. В этой и последующих статьях мы поговорим о том, как производить мониторинг работы приложений, осуществлять сбор событий и трейсинг для решения конкретных проблем в работе приложения. Эти статьи ориентированы на специалистов, занимающихся администрированием Linux, хотя о Windows мы тоже будем иногда упоминать.\xa0\xa0 Эта стать Предсказать Предсказанный кластер: 2 – статьи, связанные с обзором данных и с разработкой проектов для анализа данных.

Streamlit

Гадалка онлайн

Введите свой текст тут

Черепаха – воплощение мудрости и спокойствия в мире животных. Ее панцирь, словно древний щит, защищает от жизненных невзгод, символизируя стойкость и выносливость. Медленные движения черепахи скрывают в себе глубокую философию – она учит нас ценить каждый момент, не спеша идти к своей цели. Ее глаза, полные мудрости и таинственности, словно окно в древние времена, напоминают о важности уважения к прошлому и умении жить в гармонии с окружающим миром. Черепаха – символ долголетия и стойкости, напоминание о том, что и в медленном темпе можно достичь великих высот.

Предсказать

Предсказанный кластер

Кластер 3 – статьи, связанные с работниками, участниками проектов или компаний.

