Домашнее задание №2 Асланов А.Б., ИУ9-21М Задание 2 Запросы - это проанализированные факты из Википедии. Документы - это предложения из статей Википедии, указанных в этих фактах. Обработать морфологическим анализатором и найти наиболее релевантные предложения: По векторной модели без idf; • По TF-IDF ( df в данном случае - количество предложений, в которых встретилось слово). Нормировать запросы и предложения, то есть выстроить все предложения из статей по мере сходства с запросом по векторной модели. В отчёте должны быть показаны веса выдаваемых предложений. **Запросы** (== факты): 1) В медленном движении медленно читают; 2) Петру I доставляли из Китая посылки с золотым песком; 3) Зоны отдыха экипажей дальнемагистральных авиарейсов находятся над или под пассажирским салоном. Документы: К факту 1: Д1.1: Медленное чтение как часть медленного движения; Д1.2: Однако в последнее время наблюдается повышенный интерес к медленному чтению как направлению в рамках медленного движения. К факту 2 прямо относящихся предложений из статьи не было. Наиболее подходящие косвенные: Д2.1: Коробчатое золото — золотой песок, который поставлялся в Российскую империю из Китая для чеканки червонцев. Д2.2: Коробчатое золото использовалось для чеканки червонцев в 1701 и 1704 годах. Д3.1: Отсеки для отдыха экипажа обычно присутствуют на дальнемагистральных лайнерах и могут располагаться над пассажирским салоном (в этом случае в них нужно подниматься по лестнице) или рядом с ним. In [1]: from pymystem3 import Mystem m = Mystem() from collections import Counter import re import string import pandas as pd import numpy as np In [2]: def make morph analysis(fact): """Проводим морфологический анализ для последующей обработки""" # Удлить все знаки препинания fact = ''.join(ch for ch in fact if ch not in "[«»;,()'`--.?!:']") # Сделать все слова предложения в нижнем регистре fact = fact.lower() # Лемматизация fact = m.lemmatize(fact) fact = ''.join(fact) return fact In [3]: docs = [ 'Отсеки для отдыха экипажа обычно присутствуют на дальнемагистральных лай нерах и могут располагаться над пассажирским салоном (в этом случае в них нужно по дниматься по лестнице) или рядом с ним', 'Зона отдыха экипажа – отсек или специально отведённое место внутри самол ёта, предназначенное для отдыха членов экипажа при выполнении дальних рейсов. ', 'Члены экипажа могут воспользоваться зонами отдыха в нерабочее время', 'Модуль либо имеет форму стандартного грузового контейнера и помещается в грузовой отсек самолёта, расположенный под пассажирским салоном (вход в него осущ ествляется через специальный люк, скрытый от пассажиров), либо предназначен для уст ановки внутри пассажирского салона', 'Так, на дальнемагистральных самолётах Аэрофлота обособленный отсек отдых а есть только у лётного экипажа, бортпроводники могут отдыхать только на задних ряд ах салонов, если там нет пассажиров', 'Отсеки отдыха также могут присутствовать на транспортных и военно-трансп ортных самолётах', 'Секции для пилотов, расположенные в передней части самолёта, отделены от расположенных сзади секций для бортпроводников. ', 'Исключения могут быть сделаны для отсеков, оборудованных креслами с рем нями'] query = ['Зоны отдыха экипажей дальнемагистральных авиарейсов находятся над или п од пассажирским салоном'] In [4]: **def** extract all docs words(docs): """Извлекаем все слова из документа или запроса и прогоняем их через морфолог ический анализатор""" all docs words = [] for doc in docs: all\_docs\_words += doc.split() all docs words = ' '.join(all docs words) all\_docs\_words = make\_morph\_analysis(all\_docs\_words) all docs words = list(set(all docs words.split())) return all docs words all docs words = extract all docs words(docs) def check the same words(query, doc, counter): In [5]: """Находим количество совпадений слов в запросе и документе""" # Морфологическая обработка поданного на вход документа doc = make morph analysis(doc) # Разбиваем предложение формата <str> на список doc = doc.split() # Словари с частотами слов в запросе и в документе d doc = Counter(doc) # Если і-ое слово из запроса есть в документе, то считаем частотность этог о і-го слова в документе. # Если же такого слова нет, ставим 0. doc vec = []for word in all docs words: if word in doc: doc vec.append(d doc[word]) else: doc\_vec.append(0) d = {f'd {counter}':doc vec} dataframe = pd.DataFrame(data=d, index=all docs words) return dataframe def make termdoc matrix(query, docs): """Создаём матрицы терм-документ""" counter = 0 # счётчик документов dataframe = pd.DataFrame() for doc in docs: counter += 1 dataframe curr = check the same words(query, doc, counter) dataframe = pd.concat([dataframe, dataframe\_curr], axis=1) # Считаем df как кол-во предложений, в которых встретилось i-ое слово запр oca. # Его проще вычислить по датафрейму: df будет равен количеству ненулевых с толбцов для і-го слова. dataframe['df'] = pd.Series((dataframe != 0).astype(int).sum(axis= 1)) # Добавляем в датафрейм столбец с запросом query = make\_morph\_analysis(query) query = query.split() d\_query = Counter(query) query vec = [] for word in all\_docs\_words: if word in query: query vec.append(d query[word]) else: query\_vec.append(0) dataframe['q'] = query\_vec return dataframe dataframe = make\_termdoc\_matrix(query, docs) dataframe Out[5]: d\_1 d\_2 d\_3 d\_4 d\_5 d\_6 d\_7 d\_8 df 0 0 0 1 0 0 2 0 2 0 располагать 0 1 0 0 1 0 0 0 0 0 ряд 0 0 0 0 0 0 0 при 0 0 0 0 1 0 0 0 1 0 ОНО 1 0 0 0 0 0 0 0 0 1 иметь 0 0 0 0 1 0 0 0 1 0 аэрофлот 0 0 0 0 0 0 0 там 0 0 1 0 0 0 0 1 0 0 люк 1 0 0 1 0 0 0 0 1 0 под 2 0 1 0 1 0 0 0 0 0 предназначать 0 0 1 0 0 0 0 0 0 помещаться 0 0 0 0 0 1 0 0 0 1 также 0 0 1 0 0 0 0 0 0 располагаться 0 0 0 1 0 0 0 0 0 вход 1 0 0 0 0 0 0 контейнер 0 0 0 0 1 0 0 0 0 летный 0 0 0 0 0 1 0 0 0 1 ПО 0 0 0 0 1 0 0 0 0 1 воспользоваться 0 0 0 0 0 0 0 0 если 0 0 0 1 0 0 0 0 1 0 отделять 1 0 0 0 2 1 0 0 3 0 на 0 0 1 0 0 0 0 0 1 0 ОН 1 1 0 0 0 0 0 0 2 0 ОТ 0 0 1 0 0 0 0 0 1 0 установка 1 0 0 0 0 0 0 0 0 1 стандартный 0 1 0 1 1 1 1 0 5 0 самолет 1 0 0 0 0 0 0 0 0 специальный 0 2 0 1 0 1 0 0 0 0 внутри 0 1 0 0 0 0 0 0 1 0 они 0 0 2 0 0 0 0 0 1 0 секция ... ... ... 0 1 0 0 0 0 0 0 1 0 этот 1 0 0 0 0 0 0 0 0 1 осуществляться 2 0 0 4 1 1 0 0 1 1 экипаж 0 0 0 0 0 0 0 0 нужно 1 0 1 0 1 1 0 5 0 1 мочь 0 0 1 0 0 0 0 1 0 0 задний 0 0 2 0 1 0 0 0 0 1 член время 0 0 0 0 0 0 0 0 1 0 0 0 1 0 0 0 1 2 0 быть 0 0 0 0 0 0 0 1 0 1 кресло 0 0 0 1 0 1 0 0 0 0 рядом 0 1 0 0 0 0 0 0 0 место 2 0 0 0 0 0 0 0 1 0 либо 0 2 0 1 2 1 0 0 4 0 В 0 0 0 1 0 0 0 0 1 0 через 2 0 0 0 5 1 1 отдых 0 0 0 0 1 0 0 0 1 0 У 3 салон 0 0 военнотранспортный 0 0 0 1 0 0 0 0 1 0 нерабочий 0 0 0 0 0 1 1 0 транспортный 0 0 0 1 0 0 1 0 отдыхать 0 0 0 0 0 2 или 1 0 0 0 0 0 0 1 0 0 дальний 1 0 0 0 0 0 0 1 0 часть 0 1 0 0 1 0 0 0 0 0 так 0 0 0 отводить 0 0 0 0 1 0 0 0 1 0 нет 1 0 0 0 0 0 0 0 1 0 лайнер 90 rows × 10 columns 1. Векторная модель (без учета idf) Вес определяется как функция от количества вхождений терма в документ. In [6]: import operator from scipy.spatial import distance def calculate cosine similarity(dataframe, doc num): """Расчёт косинусного расстояния для векторной модели""" document = np.array(dataframe[f'd {doc\_num}']) query = np.array(dataframe['q']) return 1 - distance.cosine(document, query) # Запись в словарь всех косинусных расстояний для каждого конкретного документа answers =  $\{\}$ for i, doc num in zip(range(len(docs)), range(1, len(docs)+1)): answers[docs[i]] = calculate cosine similarity(dataframe, doc num) # Сортировка по значениям словаря для упорядочения порядка ранжирования def sort answers(answers, query): print('Sampoc: ', query, '\n'\*3) for ans in sorted(answers, key=answers.get, reverse=True): print(ans, answers[ans]) final\_answers = sort\_answers(answers, query) [ 'Зоны отдыха экипажей дальнемагистральных авиарейсов находятся над или под пассажирским салоном'] Отсеки для отдыха экипажа обычно присутствуют на дальнемагистральных лайнерах и могут располагаться над пассажирским салоном (в этом случае в них нужно поднимать ся по лестнице) или рядом с ним 0.42600643361512924 Зона отдыха экипажа — отсек или специально отведённое место внутри самолёта, пред назначенное для отдыха членов экипажа при выполнении дальних рейсов. 0.41702882 811414954 Члены экипажа могут воспользоваться зонами отдыха в нерабочее время 0.333333333 33333326 Модуль либо имеет форму стандартного грузового контейнера и помещается в грузовой отсек самолёта, расположенный под пассажирским салоном (вход в него осуществляетс я через специальный люк, скрытый от пассажиров), либо предназначен для установки внутри пассажирского салона 0.25125945381480297 Так, на дальнемагистральных самолётах Аэрофлота обособленный отсек отдыха есть то лько у лётного экипажа, бортпроводники могут отдыхать только на задних рядах салон ов, если там нет пассажиров 0.24759378423606915 Отсеки отдыха также могут присутствовать на транспортных и военно-транспортных са молётах 0.10540925533894596 Секции для пилотов, расположенные в передней части самолёта, отделены от располож енных сзади секций для бортпроводников. 0.0 Исключения могут быть сделаны для отсеков, оборудованных креслами с ремнями 0.0 2. Модель TF - IDFВес определяется как произведение функции от количества вхождений терма в документ и функции от величины, обратной количеству документов коллекции, в которых встречается этот терм. Расчёт косинусного расстояния для векторной модели Косинусное расстояние позволяет вычислить меру сходства двух документов как  $cos(\vec{q},\vec{d}) = \frac{\vec{q}\vec{d}}{||\vec{q}||*||\vec{d}||}$ , где q - вектор запроса, а d - вектор документа. In [11]: from scipy.spatial import distance import math def calculate\_cosine\_similarity\_tfidf(dataframe, doc\_num, N=len(docs )): tf = np.array(dataframe[f'd\_{doc\_num}']) df = np.array(dataframe['df']) idf = np.log10(N / df)document = tf \* idf query = np.array(dataframe['q']) return 1 - distance.cosine(document, query) # Запись в словарь всех косинусных расстояний для каждого конкретного документа answers tfidf = {} for i, doc num in zip(range(len(docs)), range(1, len(docs)+1)): answers tfidf[docs[i]] = calculate cosine similarity tfidf(datafra me, doc num) sort answers(answers tfidf, query) Запрос: [ 'Зоны отдыха экипажей дальнемагистральных авиарейсов находятся над или под пассажирским салоном'] Отсеки для отдыха экипажа обычно присутствуют на дальнемагистральных лайнерах и могут располагаться над пассажирским салоном (в этом случае в них нужно поднимать ся по лестнице) или рядом с ним 0.32961809135077114 Зона отдыха экипажа — отсек или специально отведённое место внутри самолёта, пред назначенное для отдыха членов экипажа при выполнении дальних рейсов. 0.25856776 2497498 Модуль либо имеет форму стандартного грузового контейнера и помещается в грузовой отсек самолёта, расположенный под пассажирским салоном (вход в него осуществляетс я через специальный люк, скрытый от пассажиров), либо предназначен для установки внутри пассажирского салона 0.2039570167126472 Члены экипажа могут воспользоваться зонами отдыха в нерабочее время 0.199098183 92541212 Так, на дальнемагистральных самолётах Аэрофлота обособленный отсек отдыха есть то лько у лётного экипажа, бортпроводники могут отдыхать только на задних рядах салон ов, если там нет пассажиров 0.1327452154007247 Отсеки отдыха также могут присутствовать на транспортных и военно-транспортных са молётах 0.037382952568958516 Секции для пилотов, расположенные в передней части самолёта, отделены от располож енных сзади секций для бортпроводников. 0.0 Исключения могут быть сделаны для отсеков, оборудованных креслами с ремнями 0.0 Результаты Жирным шрифтом выделены предложения, которые, как ожидал при выделении вручную, должны быть наиболее релевантны запросу. Однако сейчас соглашусь с моделями, которые утверждают, что предложение Коробчатое золото использовалось для чеканки червонцев в 1701 и 1704 годах не является настолько подходящим, как предложения, которым присвоены бОльшие вероятности. В целом системы на данных примерах сопоставимы по качеству, однако TF-IDF присваивает вариантам в целом меньшие веса, то есть более "не уверен" в верности выдачи. Запрос: ['В медленном движении медленно читают'] Векторная модель 0.5Однако в последнее время наблюдается повышенный интерес к медленному чтению как направлению в рамках медленного движения. 0.4743416490252569 0.33541019662496846 Однако теоретики метода утверждают, что идея медленного чтения зародилась гораздо раньше появления медленного движения. 0.33541019662496846 Создатели сайта утверждают, что лучший способ быть частью медленного движения — это читать хорошие книги 0.3113995776646092 В настоящее время существует сайт медленного движения, который предоставляет доступ к информации, ресурсам, услугам и сетевым возможностям всем, кто 0.28284271247461895 Последователи методики медленного чтения не являются организованными в своем движении: «нет ни одного бланка, ни совета директоров, ни, о ужас, центра. 0.10540925533894598 Метод медленного чтения позволяет более полно понять и оценить сложный текст при изучении философии и литературы Модель TF-IDF 0.22677868380553634 Прежде всего, будем говорить медленно 0.18527268850463563 Однако в последнее время наблюдается повышенный интерес к медленному чтению как направлению в рамках медленного движения. 0.161980178 Создатели сайта утверждают, что лучший способ быть частью медленного движения — это читать хорошие книги 0,137428803 0,107582101 В настоящее время существует сайт медленного движения, который предоставляет доступ к информации, ресурсам, услугам и сетевым возможностям всем, кто 0,090019083 Последователи методики медленного чтения не являются организованными в своем движении: «нет ни одного бланка, ни совета директоров, ни, о ужас, центра. 0,042542798 Однако теоретики метода утверждают, что идея медленного чтения зародилась гораздо раньше появления медленного движения 0,008051486 Метод медленного чтения позволяет более полно понять и оценить сложный текст при изучении философии и литературы Запрос: ['Петру I доставляли из Китая посылки с волотым песком'] Векторная модель 0.5735393346764044 Из Китая поступал в основном золотой песок, а золотые слитки привозили в намного меньшем количестве 0.5345224838248488 Коробчатое золото — золотой песок, который поставлялся в Российскую империю из Китая для чеканки червонцев 0.3779644730092272 Золотой песок упаковывали в коробочки небольших размеров 0.32025630761017426 Золотые двухрублёвики, которые также называли андреевскими золотыми из-за размещения на одной из сторон изображения апостола Андрея Первозвани 0.15075567228888181 По одним данным, коробчатое золото, поступавшее на Монетный двор через Сибирский приказ, использовалось иногда для чеканки двухрублёвиков, а по ; 0.0 Коробчатое золото использовалось для чеканки червонцев в 1701 и 1704 годах 0.0 Уже в конце XVII века через китайскую границу коробчатое золото стало поставляться в значительном количестве 0.0 Деньги пошли на оплату работы мастеров, стоимость самого золота, потери при плавке Модель TF-IDF 0,39168 Коробчатое золото — золотой песок, который поставлялся в Российскую империю из Китая для чеканки червонцев 0,353772759 Из Китая поступал в основном золотой песок, а золотые слитки привозили в намного меньшем количестве 0,168373998 Золотой песок упаковывали в коробочки небольших размеров 0,123042338 Золотые двухрублёвики, которые также называли андреевскими золотыми из-за размещения на одной из сторон изображения апостола Андрея Первозванн 0,054296564 По одним данным, коробчатое золото, поступавшее на Монетный двор через Сибирский приказ, использовалось иногда для чеканки двухрублёвиков, а по ; Коробчатое золото использовалось для чеканки червонцев в 1701 и 1704 годах 0.00.0 Уже в конце XVII века через китайскую границу коробчатое золото стало поставляться в значительном количестве 0.0 Деньги пошли на оплату работы мастеров, стоимость самого золота, потери при плавке 0.0 Запрос: ['Зоны отдыха экипажей дальнемагистральных авиарейсов находятся над или под пассажирским салоном'] Векторная модель 0.42600643361512924 Отсеки для отдыха экипажа обычно присутствуют на дальнемагистральных лайнерах и могут располагаться над пассажирским салоном (в этом случае в них нужно подниматься по лестнице) и 0.41702882811414954 Зона отдыха экипажа — отсек или специально отведённое место внутри самолёта, предназначенное для отдыха членов экипажа при выполнении дальних рейсон 0.3333333333333333 Члены экипажа могут воспользоваться зонами отдыха в нерабочее время 0.251259453814803 Модуль либо имеет форму стандартного грузового контейнера и помещается в грузовой отсек самолёта, расположенный под пассажирским салоном (вход в него осуществляется через специальный люк, скры 0.24759378423606918 Так, на дальнемагистральных самолётах Аэрофлота обособленный отсех отдыха есть только у лётного экипажа, бортпроводники могут отдыхать только на задних рядах салонов, если там нет пассажирог 0.10540925533894598 Отсеки отдыха также могут присутствовать на транспортных и военно-транспортных самолётах 0.0 Исключения могут быть сделаны для отсеков, оборудованных креслами с ремнями Модель TF-IDF 0.329618091 Отсеки для отдыха экипажа обычно присутствуют на дальнемагистральных лайнерах и могут располагаться над пассажирским салоном (в этом случае в них нужно подниматься по лестнице) и 0,258567762 Зона отдыха экипажа — отсек или специально отведённое место внутри самолёта, предназначенное для отдыха членов экипажа при выполнении дальних рейсов. 0,203957017 Модуль либо имеет форму стандартного грузового контейнера и помещается в грузовой отсек самолёта, расположенный под пассажирским салоном (вход в него осуществляется через специальный люк, скры 0,199098184 Члены экипажа могут воспользоваться зонами отдыха в нерабочее время 0.132745215 Так, на дальнемагистральных самолётах Аэрофлота обособленный отсек отдыха есть только у лётного экипажа, бортпроводники могут отдыхать только на задних рядах салонов, если там нет пассажирог 0,037382953 Отсеки отдыха также могут присутствовать на транспортных и военно-транспортных самолётах Секции для пилотов, расположенные в передней части самолёта, отделены от расположенных сзади секций для бортпроводниког 0.0 Исключения могут быть сделаны для отсеков, оборудованных креслами с ремнями Перепроверка Возьмём результаты векторной модели и посмотрим, правильно ли было вычислено косинусное расстояние между векторами слов. Здесь запрос q: "В медленном движении медленно читают", а документ d: "Прежде всего, будем говорить медленно" (предложения лемматизированы). In [8]: from scipy.spatial import distance q = [1, 1, 1, 1, 1, 0, 0, 0, 0] # в медленный движение медленно читать d = [0, 0, 0, 1, 0, 1, 1, 1, 1] # прежде всего быть говорить медленно print(1 – distance.cosine(d, q)) # б/м недоставание до 0,2 в этом случае - издержка интерпретатора python'a 0.1999999999999996 Lemma d q 0 1 1 0 медленный 1 движение 1 медленно 1 читать 0 прежде всего 0 0 быть 0 говорить  $cos(\theta) = \frac{q*d}{||q||*||d||} = \frac{0+0+0+1+0+0+0+0}{\sqrt{1^2+1^2+1^2+1^2+1^2+0^2+0^2+0^2+0^2}\sqrt{0^2+0^2+0^2+1^2+0^2+1^2+1^2+1^2+1^2}} = \frac{1}{\sqrt{5}\sqrt{5}}$ = 0, 2Перепроверка 2 Проверка вычисления векторной модели между запросом  $q = \mathbf{B}$  медленном движении медленно читают и документом d =Медленное движение как часть медленного движения. Lemma d медленный 1 движение медленно 1 читать 1 0 как 0 1 часть чтение  $cos(\theta) = \frac{\vec{q} * \vec{d}}{||\vec{q}|| * ||\vec{d}||} = \frac{0 + 2 + 1 + 0 + 0 + 0 + 0}{\sqrt{0^2 + 2^2 + 1^2 + 0^2 + 0^2 + 1^2 + 1^2} \sqrt{1^2 + 1^2 + 1^2 + 1^2 + 1^2 + 0^2 + 0^2 + 0^2}} = \frac{3}{\sqrt{8}\sqrt{5}} \approx 0,474$ In [ ]: