

KOCAELİ ÜNİVERSİTESİ

BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

YAZILIM LAB.II – 1.PROJE:

Web İndeksleme Uygulaması

MÜRVET NUR ŞEN- ÜMMÜHAN TEPEBAŞ
190201097 – 180201088

ÖZET:

Bizden istenen; verilen bir URL'deki web sayfa içeriğine göre diğer birden fazla web sayfasını benzerlik bakımından indeksleyip sıralayan web tabanlı bir uygulama geliştirmektir.

Projemizi geliştirirken Python-Flask Web Framework' ü ve Javascript dilini ve Visual Studio Code geliştirme ortamını kullandık. Projemiz İngilizce dilinde web siteleri kullanılarak geliştirilmiştir.

GENEL BİLGİLER :

Proje dokümanında belirtilen “Aşama 1: Sayfada geçen kelimelerin frekanslarını hesaplama” adımı için BeautifulSoup modülü kullanılmıştır. Http requests ataması için JQuery- Axios kullanılmıştır.

Web scraping işlemlerinde Python – Natural Language Toolkit kullanılmıştır. NLTK kütüphanesi için detaylı bilgiler şöyledir:

NLTK, doğal dil araç takımı anlamına gelir. Natural Language Toolkit; insan dili verileriyle çalışmak için Python programlama dili ile geliştirilmiş ve geliştirilmekte olan 50'nin üzerinde derleme(corpus) ve sözcük kaynağı(lexical resources) ile oluşturulmuş açık kaynaklı bir kütüphanedir.

NLTK, bir veri setinin ön işlemlerini(preprocessing), yani veriyi makinenin anlayacağı hale getireceğimiz zaman bizi gereksiz kelimelerle uğraştırmaktan da kurtarıyor. Biz de bu özelliği sayfalardaki (İngilizce dilindeki) bağlaç ve gereksiz kelimeleri çıkarmak için kullandık .

<http://127.0.0.1:5000/frekans> Sayfası
Aşama 1 için verilen sitedeki tüm kelimeleri frekansları ile birlikte başarılı olarak döndürmektedir.

“Aşama 2: Anahtar kelime çıkarma ve Aşama 3: Benzerlik skorlaması” adımlarında girilen her iki Url için 15 anahtar kelimeyi frekansları ile birlikte döndürmektedir. Ve Aşama 3 için benzerlik oranını her iki Url'de dönen aynı anahtar kelimelerin frekanslarının, toplam anahtar kelime frekansına yüzdesi ile hesapladık.

<http://127.0.0.1:5000/scoring> Sayfasında Bootstrap progress bar kullanarak yüzdeyi ve anahtar kelimeleri başarılı bir şekilde ekrana bastırdık.

“Aşama 4: Site indeksleme ve sıralama”
<http://127.0.0.1:5000/indexing> Sayfasında

Verilen Url içindeki Url'lerde gezerek kelimeleri başarılı bir şekilde çıkarır.

“Aşama 5: Semantik Analiz” adımımda NLTK “from nltk.corpus import wordnet” kütüphanesini kullanarak girilen Url deki anahtar kelimelerin benzer, yakın anlamlı kelimeleri sayfaya döndürüyor. Anahtar kelimeleri kullanırken NLTK - “from nltk.stem import PorterStemmer” kütüphanesini kullanarak kelimeyi kök halinde işleme alır. Bu kelimelerin anlamsal olarak alakalı, benzer olma oranını inisiyatif olarak 0.65% den büyük olan kelimeleri kullanmak üzere sayfaya yazdırır.

<http://127.0.0.1:5000/semantik> Sayfasında girilen ilk Url’deki anahtar kelimeleri ikinci Url’deki bunlarla yakın anlamlı kelimelerle birlikte başarılı olarak yazdırmaktadır.

KULLANILAN KÜTÜPHANELER:

- ➔ `scraper_controller.py`
 - `from flask import Flask, render_template, jsonify`
 - `from flask import request`
 - `from utils import scrape_url, scrape_url2, recursive_helper, semantik`
- ➔ `utils.py`
 - `import io`
 - `import re`
 - `import requests`
 - `import operator`
 - `from nltk.corpus import wordnet`
 - `from nltk.corpus import stopwords`
 - `from snowballstemmer import stemmer`
 - `from nltk.stem import PorterStemmer`
 - `from bs4 import BeautifulSoup, Comment`
 - `from nltk.tokenize import word_tokenize`

KULLANILAN FONKSİYONLAR:

➔ `def scrape_url(url):`
`def sozlukolustur(tumkelimeler):`

`def sembolleritemizle(tumkelimeler):`

Fonksiyona alınan web sitesi URL'indeki kelimeleri alıp boşluklarına göre ayırıp gereksiz sembollerden ayırmak için kullanılan bir fonksiyondur.

➔ `def scrape_url2(url):`

Yukarıda belirttiğim "scrape_url" fonksiyonun işlemlerine ek olarak Aşama 2 ve 3 için anahtar kelimeleri return eder.

➔ `def create_nested_list(limit):`

Verilen web sitesi kümesini 3 derinlikli olarak alır

➔ `def get_all_links(soup, current_url):`

Verilen sitenin içindeki "href" etiketli linkleri alır

➔ `def`

`recursive_scrapper(input_url_list):`

➔ `def scrape_url3(url):`

Girilen Urlde aşağıdaki fonksiyonlarla işlem yapıyor.

➔ `def sozlukolustur(tumkelimeler):`

➔ `def`

`sembolleritemizle(tumkelimeler):`

➔ `def get_all_links(soup):`

➔ `def semantik(first_url,second_url):`

Fonksiyonu ile sayfadan 2 farklı Url alır. “scrape_url2” fonksiyonu çağırarak her Url'e ait ayrı semantik_list return eder

➔ `def find_similarity(first_list, second_list):`

Semantik_listlerdeki benzer, yakın anlamlı kelimeleri return eder.

Önemli bilgi: İşlem yoğunluğunun fazla zaman alması sebebiyle verilen her iki URL'deki tüm kelimelerin semantik karşılığını bulmak yerine kısa süreli sonuç

için ilk 15 anahtar kelimeyi kütüphane içerisinde işledik.

PROJE GENEL GÖRÜNÜM

Our Search Engine

SKORLAMA

İNDEKSLEME

SEMANTİK

MURVET&UMMUHAN

Hoşgeldiniz...

Bu sitede url girişi ile arama sağlanacak.

Search

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

#	Keyword
1	computer
2	computers
3	first
4	memory
5	machine

Page: <http://127.0.0.1:5000/frekans>

Our Search Engine

SKORLAMA

İNDEKSLEME

SEMANTİK

MURVET&UMMUHAN

Skorlama

Bu aşamada her iki URL metninde geçen kelimelerden en önemli kelimelerin belirlenerek anahtar kelimelerin çıkartılması işlemi gerçekleştirilecek ve ilk URL için elde edilen anahtar kelimelerin 2. URL'nin içerisinde yer alma sayısına dayalı bir benzerlik skor formülü uygulanacak.

Search

Search

Benzerlik Oranı Hesapla

1#	Keyword	Value
1	computer	255
2	computers	132

6	program	50
7	devices	49
2#	Keyword	Value
1	computer	138
2	engineering	95
3	systems	35
4	software	29
5	engineers	28
6	design	25
7	degree	21
GRAM-37,861271676300575%		

Page: <http://127.0.0.1:5000/scoring>

İndeksleme

Bu aşamada projenin web sitesindeki bir sayfada URL girilecek bir alan oluşturulacaktır. Girilen bu URL'nin içeriği ile web site kümesindeki her bir web sayfasının içeriklerinin benzerlik skorları ayrı ayrı hesaplanacaktır. Ancak bu sefer skor hesaplaması yaparken bu web site kümesinde bulunan web sayfalarının içeriğine ilaveten yine bu sayfalarda bulunan tüm alt URL'leri de dikkate alınacaktır. Alt URL'lerindeki anahtar kelimelerin yer alma sayılarına dayalı olarak skor formülünü yeniden geliştirilecek.

Search

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

1#	Keyword	Value
1	https://en.wikipedia.org/wiki/Computer	computer,computers,first,memory,machine,program,devices,retrieved,used,may,electronic,computing,digital,programs,instructions,isbn,integrated,data,system,history,modern,main,software,circuit,lan,guages,hardware,logic,unit,control,p,device,s,programming,syste

		l,processes,culture,study,man,illustrations,plain,mesopotamia,president,day,iraq,tetrahedron,labortexas,center,trussed,duck,ox,vessel,fruit,pomegranate,textile,honey,jar,oil,fleece,wool,ingot,metal,content,tablet,featuring,godin,t,cuyler,young,royal,ontario,museum,canada
4	https://www.facebook.com/sharer/sharer.php?u=https%3A%2F%2Fwww.etymonline.com%2Fword%2Fcomputer	facebook,facebookpaylaşmak için,hesabına,giriş,yaptelefon,numarası,veya,e-postafacebook,çiftresişifreni,mıunuttunveyatürkçearabieespañolkurdikurmanclenglish,uk,inc
5	https://twitter.com/share?url=https%3A%2F%2Fwww.etymonline.com%2Fword%2Fcomputer&text=computer%20%7C%20Origin%20and%20meaning%20of%20computer%20by%20Online%20Etymology%20Dictionary	javascript,browser,supported,help,center,policy,available,we've,detected,disabled,please,enable,switch,continue,using,twittercom,select,list,browsers,terms,service,privacy,cookie,imprint,ads,info,twitter,inc,something,went,wrong,don't,fret,—let's,give,another,shot
6	https://www.reddit.com/submit?url=https%3A%2F%2Fwww.etymonline.com%2Fword%2Fcomputer&title=computer%20%7C%20Origin%20and%20meaning%20of%20computer%20by%20Online%20Etymology%20Dictionary	press,original,content,t,jump,feed,question,mark,learn,rest,keyboards,shortcutslog,insign,upuser,account,menucreate,post,draft,postimages,&,videolinkpollbolditalicslinkstrikethroughline,code,super,scriptspoil,erheading,bulleted,list,numbered,list,quote,block,code,blockquote,add,image,add,videomarkdown,mode,ec,this,community,allow,tags,spoiler,marks,spoilersfwmarks,safe,work,fair,select,subredditlenable,flair,postsave,draft,send,post,reply,notification,posting,redditlemember,human,behave,like,would,real,lifelook,source,content,search,duplicates,posting,read,community's,rules,please,mindful,rel

Page: <http://127.0.0.1:5000/indexing>

Semantik Analiz

Burada verilen web siteleri içerisinde anahtar kelimelerle semantik olarak alakalı kelimeler ayıklanır.

1#	Keyword	Value
1	computer	computer,systems,hardware
2	computers	computer,systems,hardware
3	memory	science
4	machine	computer,systems,hardware

1#	Keyword	Value
1	computer	computer,systems,hardware
2	computers	computer,systems,hardware
3	memory	science
4	machine	computer,systems,hardware
5	program	science
6	retrieved	retrieved
7	computing	science
8	programs	science

Page: <http://127.0.0.1:5000/semantik>

KAYNAKÇA:

- Kaynakça 1:
 - https://towardsdatascience.com/synonyms-and-antonyms-in-python-a865a5e14ce8#_
- Kaynakça 2:
 - <https://www.veribilimiokulu.com/natural-language-toolkitnltk/>
- Kaynakça 3:
 - <https://www.nltk.org/data.html>
- Kaynakça 4:
 - <https://www.youtube.com/watch?v=II9mlqPmgII&t=1s>
- Kaynakça 5:
 - Udemy - Mustafa Murat Coşkun (40+ Saat) Python | Sıfırdan İleri Seviye Programlama (2020)

