

WEB İNDEKSLEME UYGULAMASI

Alparslan Beraat ÖZDEMİR – 170202045
Bilgisayar Mühendisliği Bölümü
Kocaeli Üniversitesi
Kocaeli, Türkçe
beraat78@gmail.com

ÖZET

Bu proje, 5 farklı PHP programından oluşmaktadır. Her aşama, farklı işlevler yerine getirir.

İlk aşama, verilen bir web sayfasındaki kelimelerin frekanslarını bulur. Bunları ekrana yazdırır.

İkinci aşama, verilen bir web sayfasındaki anahtar kelimeleri ekrana yazdırır.

Üçüncü aşama, iki web sayfası arasındaki benzerlik oranını bulur ve ekrana yazdırır.

Dördüncü aşama, verilen bir URL kümesini benzerliklerine göre sıralar.

Beşinci aşama, verilen bir web sayfasının anahtar kelimelerinin eş anlamlılarını bularak ekrana yazdırır.

I. GİRİŞ

Web indeksleme uygulaması, 5 adet farklı php dosyası ve 1 adet html2text.php kütüphanesinden oluşmuştur

Bir apache sunucusunda dosyalar tarayıcı yardımı ile ayrı ayrı çalıştırılabilir.

Her sayfa önce html kısmıyla input alanlarını barındırır, daha sonra PHP kodu kısmı yer alır.

Bütün aşamalarda string temizleyici ayıkla fonksiyonu standart olarak yer almaktadır.

Temizle fonksiyonu da aynı şekilde, özel karakter içeren kelimeleri temizlemek için standart olarak yer almaktadır.

II. TEMEL BİLGİLER

Proje gelişiminde;

Tümleşik geliştirme ortamı olarak PHPStorm kullanılmıştır.

Program Windows 10 İşletim sisteminde test edilmiştir. Program PHP dilinde 8.0 versiyonunda geliştirilmiştir.

III. TASARIM

Web indeksleme projesinin programlanma aşamaları altta belirtilen başlıklar altında açıklanmıştır.

A. Kullanıcı arayüzünü oluşturmak.

Temel HTML ile her sayfada amaca uygun, çok sade şekilde tasarlanmış formlar yer almaktadır. Tasarım kısmına önem verilmemiş, sadece görevini yerine getirecek kadar bir tasarım yapılmıştır.

B. Fonksiyonlar

Main Fonksiyonu:

Php dil yapısı gereği özel bir main fonksiyonu tanımlanmamıştır. Programın php taglerinin başladığı yerden itibaren main fonksiyonu olarak sayılabilir, HTML formlarından verilerin alınması genelde bu aşamada yer almaktadır. Ayrıca, bazı aşamalarda gerekli işlemlerin temel kısmı bu kısımda yapılmıştır.

Ayıkla Fonksiyonu:

Bu fonksiyon, bütün aşamalarda standart olarak yer almaktadır. İşlem yapılacak URL leri parametre olarak alan bu fonksiyon, sayfa kaynak kodunu en ideal hale gelene kadar onlarca filtrasyona tabi tutar. Bu filtrasyonların büyük çoğunluğu, Html2Text kütüphanesi kullanılarak yapılmıştır. Tüm HTML taglerini büyük oranda temizleyen bu kütüphane, sayfa içeriğini büyük oranda temizlenmiş bir string olarak return eder ve diğer filtreler de uygulandıktan sonra, sadece sayfada “anlam ifade eden” kelimeler kümesi kalır.

Temizle fonksiyonu:

Bu fonksiyon, içerisinde özel karakter yer alan bütün kelimeleri yok eder. Diğer filtrasyonlar uygulanırken

bozulmuş olarak kalan kelimeleri bu fonksiyon aracılığı ile regexler kullanılarak yok edilir.

Esanlambul fonksiyonu:

Bu fonksiyon, sadece asama5.php dosyasında, son aşamada yer alır. Sayfanın anahtar kelimelerini parametre olarak alan bu fonksiyon, online bir sözlük hizmetinin API sini kullanarak parametre olarak aldığı kelimelerin eş anlamlılarını uzak sunucuda sorgular. API, bu kelimenin eş anlamlılarını response olarak döndürür, ve fonksiyon tarafından anahtar kelimelerin eş anlamlıları eğer bulunabildiyse ekrana yazdırılır.

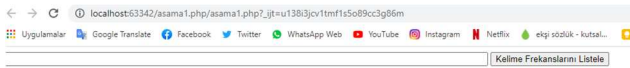
C. Algoritma

Program, birden fazla aşamadan oluştuğundan dolayı, her aşamada farklı bir algoritma kullanılmıştır. Genel olarak program akışı şu şekildedir:

- 1- Başla
- 2- Kullanıcıdan gerekli inputları form üzerinden al.
- 3- POST isteği ile PHP kodundaki değişkenlere ata.
- 4- Hedef adreslere istek atarak sayfa kaynaklarını değişkene ata.
- 5- Sayfa kaynaklarına gerekli filtreleri uygula.
- 6- Eğer gerekiyorsa, sonuçları birbiriyle kıyasla.
- 7- İsterleri kaydet.
- 8- Ekrana çıktı olarak web sayfasında yazdır.

IV. SONUÇLAR VE EKRAN ÇIKTILARI

Aşama1'in çalışması ve sonucu şekil 1 ve 2'de gösterilmiştir.



Şekil1.

Kalime Frekanslarını Listele	
Array ([array] => 54 [echo] => 21 [badges] => 21 [stack] => 20 [overflow] => 16 [string] => 16 [stuff] => 15 [teams] => 14 [post] => 13 [notice] => 13 [code] => 11 [value] => 9 [using] => 9 [language] => 9 [any] => 8 [conversion] => 8 [html] => 8 [input] => 8 [output] => 8 [convert] => 8 [sign] => 7 [when] => 7 [gold] => 7 [prints] => 7 [into] => 6 [success] => 6 [success] => 6 [policy] => 6 [success] => 6 [found] => 5 [exchange] => 5 [here] => 5 [done] => 5 [initial] => 5 [answered] => 5 [answered] => 4 [answered] => 4 [result] => 4 [data] => 4 [undefined] => 4 [about] => 3 [developer] => 3 [meta] => 3 [blog] => 3 [join] => 3 [create] => 3 [active] => 3 [edited] => 3 [previous] => 3 [comment] => 3 [access] => 3 [format] => 3 [error] => 3 [structure] => 3 [development] => 3 [handle] => 3 [valid] => 3 [application] => 3 [object] => 3 [item] => 3 [cookies] => 3 [products] => 2 [public] => 2 [technology] => 2 [private] => 2 [friend] => 2 [career] => 2 [build] => 2 [feature] => 2 [feature] => 2 [unit] => 2 [asked] => 2 [month] => 2 [body] => 2 [type] => 2 [like] => 2 [instead] => 2 [phrase] => 2 [text] => 2 [text] => 2 [turn] => 2 [male] => 2 [female] => 2 [named] => 2 [cell] => 2 [join] => 2 [input] => 2 [plot] => 2 [function] => 2 [team] => 2 [content] => 2 [purpose] => 2 [most] => 2 [them] => 2 [being] => 2 [example] => 2 [level] => 2 [prevent] => 2 [might] => 2 [cause] => 2 [possible] => 2 [nikola] => 2 [haric] => 2 [badge] => 2 [right] => 2 [away] => 2 [variable] => 2 [index] => 2 [diagonal] => 2 [rep] => 2 [something] => 2 [item] => 2 [specific] => 2 [index] => 2 [reference] => 2 [amodulus] => 2 [article] => 2 [box] => 2 [skiller] => 2 [truck] => 2 [police] => 2 [spectra] => 2 [p] => 2 [english] => 2 [gamer] => 2 [game] => 2 [accept] => 2 [cooperate] => 1 [programming] => 1 [technical] => 1 [opportunity] => 1 [career] => 1 [tech] => 1 [empty] => 1 [sites] => 1 [script] => 1 [hub] => 1 [taxi] => 1 [find] => 1 [pages] => 1 [collaborate] => 1 [group] => 1 [owned] => 1 [connect] => 1 [location] => 1 [structure] => 1 [es]	

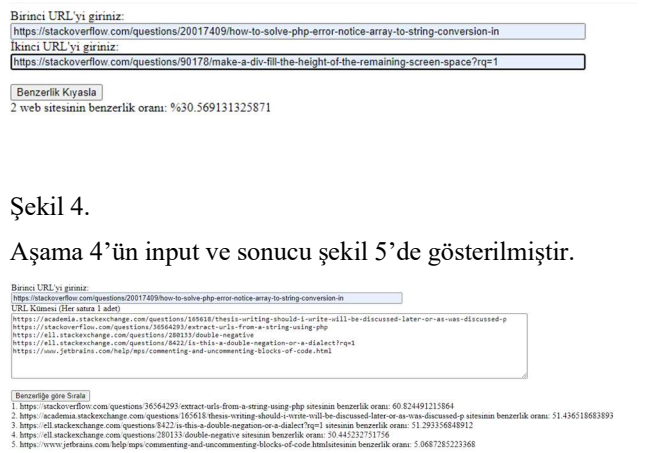
Şekil 2.

Aşama2'nin çalışması ve sonucu şekil 3'de gösterilmiştir.



Şekil 3.

Aşama3'ün çalışması ve sonucu şekil 4'de gösterilmiştir.



Şekil 5.

Aşama5'in Çalıştırılması ve çıktısı şekil6'da gösterilmiştir.



KAYNAKÇA

- [1] <https://github.com/mtibben/html2text>
- [2] <https://www.wordsapi.com/#try>
- [3] <https://stackoverflow.com/questions/9442249/php-regex-remove-words-from-string-which-contain-non-letters-numbers>