

Veri Analizi ve Kullanıcı İlişkilendirmesi Raporu

Alperen Karacan
Öğrenci No: 210201096

I. ÖZET

Bu proje, JSON data seti üzerinden elde edilen kullanıcı verilerini analiz etmeyi ve bu veriler üzerinden kullanıcıların ilgi alanlarını belirlemeyi amaçlamaktadır. Proje, kullanıcıların, kullanıcı adları, adları, soyadları, bölgeleri, tweetlerini, takipçi sayılarını ve takip ettikleri hesapları içeren geniş bir veri setinden yararlanır. Hedef, bu verileri işleyerek kullanıcıların ilgi alanlarını çıkarmak ve bu bilgileri graph olarak sunmaktır. Proje, metin işleme, veri yapıları, algoritmalar ve arayüz geliştirme gibi birçok programlamayı öğrenip kullandırmayı amaçlamaktadır. Kazanım:

- Bu proje, öğrencilere veri yapıları, graflar, hash tabloları, arama algoritmaları, veri analizi gibi konularda deneyim kazanma fırsatı sunacaktır. Hem veri yapıları ve algoritmaların hem de gerçek veri analizi ve ilişkisel çıkarımların pratiğini yaparak problem çözme ve analitik düşünme becerilerimizi geliştirmeye fayda sağlayacaktır.

II. GİRİŞ

Bu projede, kullanıcıların etkileşimlerinden yola çıkarak ilgi alanlarını tespit etmektir. Bu tespit, veri toplama, işleme, analiz etme ve sonuçları görselleştirme adımlarından oluşur. Projenin temel amacı, kullanıcıların hangi konulara daha fazla ilgilendiklerini anlamaktır.

III. PROJENİN AŞAMALARI

A. 1- Veri Kaynağı

Projede 50 bin kullanıcının verilerini içeren bir JSON dosyasını veri kaynağı olarak kullandım fakat 50 bin kullanıcının verilerini işlemek çok uzun sürdüğü için ne kadar kullanıcının verilerinin alınacağını fonksiyon yazarak esnek bıraktım.

B. 2- Kullanıcı

Kullanıcı bilgileri User sınıfı aracılığıyla temsil edilir. Veri depolama ve erişim için özelleştirilmiş HashTable ve Hashmap gibi veri yapıları kullanılır.

C. 3- Analiz

Kullanıcıların tweetleri önemsiz sayılan kelimeler çıkartılarak ve kelimeleri köklerine ayırarak analiz edilir.

D. 4-İlgi Alanına Göre Eşleme

Analiz sonucunda elde edilen veriler, önceden belirlenen ilgi alanlarına göre kategorize edilir. Bu kısım kullanıcıların hangi konulara daha fazla ilgi gösterdiğini ortaya çıkarır.

- Kullanıcı ilgi alanlarını hash tablolarında tutulmalıdır. Her bir ilgi alanı için bir hash tablosu oluşturulmalı ve ilgi alanını anahtar olarak kullanarak bu ilgi alanını paylaşan kullanıcıların listesini değer olarak saklanmalıdır.
- Kullanıcıların takip ettikleri, takipçileri ve tweet içerikleri arasında benzerlikleri tespit ederek, ortak ilgi alanlarına sahip kullanıcıları eşleştirilmelidir.
- Verilen bir ilgi alanı arama algoritmaları (hashing vb.) kullanılarak bulunmalıdır. Ayrıca bu ilgi alanını paylaşan kullanıcıların listesi çıkarılmalıdır.
- Kullanıcıların ortak ilgi alanlarını belirlemek için hash tabloları ve karşılaştırma algoritmaları (arama, sıralama algoritmaları vb.) kullanılmalıdır.
- Hash tabloları ve arama algoritmaları kullanılarak elde edilen benzer ilgi alanlarına sahip kullanıcılar arasındaki ilişkiler gösterilmelidir.

- Kullanıcıların ilgi alanlarına dayalı detaylı analiz raporları oluşturularak metin dosyasına kaydedilmelidir
- Belirli bölge (konum) ve dil için trend olan hashtagler veya konular listelenmelidir.

E. 5- Graf Üzerinde Analiz

Kullanıcıdan alınan veri ile ilgili kişinin takip ettikleri kullanıcılar ile grafi çizilir.Aşağıda örnek bir graf bulunmaktadır.

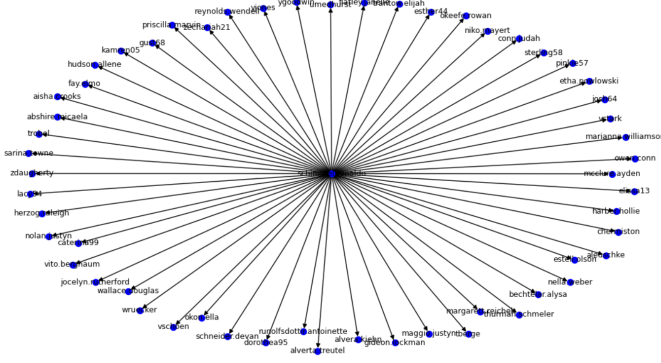


Fig. 1. Graf Örneği

IV. DENEYSEL SONUÇLAR

A. Veri Seti ve İşlenmesi

Projede kullanılan veri seti,JSON dan elde edilen 50.000 kullanıcıya ait verileri içermektedir.Bu veriler,kullanıcı bilgileri,tweet içerikleri,takipçi sayıları,takip ettikleri hesaplar ve dil,bölge gibi bilgilerdir.Verit setini daha yönetilebilir ve daha hızlı olması için boyutu indirgeyere ilk 100 kullanıcıya kadar kısıtladım,Bu kısıtı kendiniz esnek bir şekilde değiştirebilirsiniz fonksiyonunu yazdım.Bu işlem,Analiz süresini ve kaynak kullanımını olabildiğince az kullanmak için önemlidir.

B. Tweet Analizi ve Kelime Frekansı

Kullanıcıların tweetlerinin analizi,Her tweetteki kelimeleri Türkçe köklerine indirgeyerek ve önceden belirlediğim önemsiz kelimeleri hariç tutarak gerçekleştirilmiştir.Bu işlem,doğal dil işleme tekniklerini kullanarak metinden anlamlı bilgiler çıkarmayı hedefler.Analiz sonucunda,her kullanıcının tweetlerinde en sık geçen kelime ve bu kelimelerin frekansları tespit edilmiştir.Bu sonuçlar, kullanıcıların hangi konular üzerinde daha fazla konuştuğunu ve hangi kelimelerin onlar için önemli olduğunu belirlemek için kullanılır.Kullanılan algoritma bubble sort dur.Merge ve pythonun yerleşik sıralama algoritmalarıyla da denedim aralarında bariz performans farkları olmamasından dolayı bubble sort ile yaptım sizler diğer bahsettiğim algoritmaları da koda kolayca entegre ederek yapabilirsiniz.

```

Hangi kullanıcının analizini almak istersiniz? (1 - 100, Çıkmak ve kullanıcı graphını açmak için -1): 3

Kullanıcı Adı: monserrat82
Adı: Freda Abbott IV
Takipçi Sayısı: 100
Takip Ettiği Kişi Sayısı: 84
Dil: sl
Bölge: IQ

Kullanıcının İlgi Alanları ve İlgili Kelimeler:
İlgili Kategori: Siyasi, İlgili Kelime: savaş,Kelimenin Frekansı: 3

En Çok Kullanılan İlk 5 Kelime:
yer: 5
başladı: 4
yap: 4
savaş: 3
kdy: 3

```

Fig. 2. Tweet Analizi ve Frekans Örneği

VI. PROJENİN UML DİYAGRAMI VE AKIŞ ŞEMASI

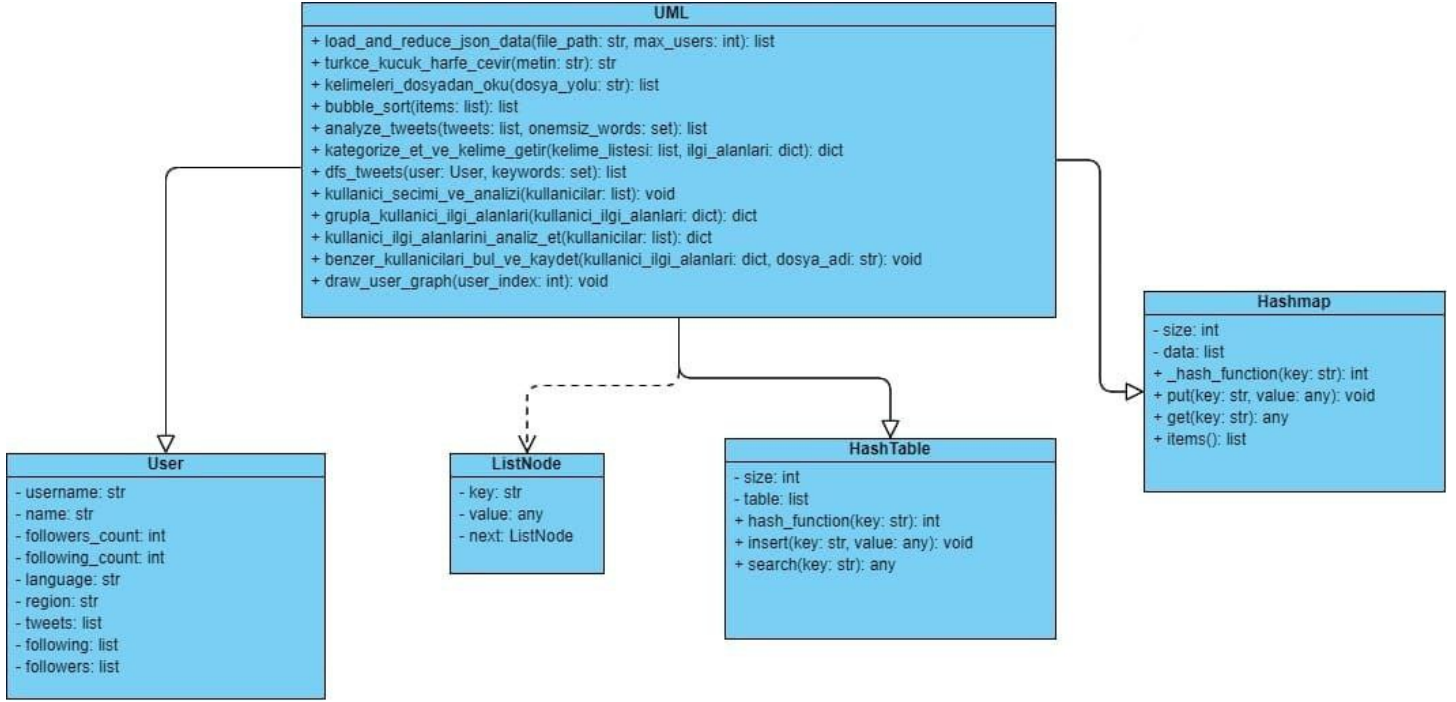


Fig. 7. UML Diyagramı

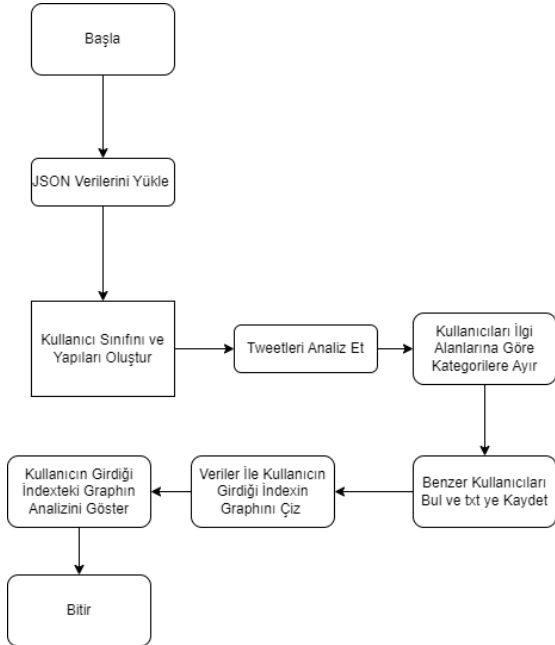


Fig. 8. Akış Diyagramı