



Doğal Dil İşlemeye Giriş Ödevi

Dersin Yürütücüsü: *Banu Diri*

Öğrencinin;

Adı: Muhammet Kayra

Soyadı: Bulut

Numarası: 20011901

Giriş

Verilen ödevde REGEX ile adres verisi içerisinde adres bloklarının ayrıştırılması işlemi yapılmıştır.

Veri

600 (599) satırdan oluşan veri her satırda farklı şekillerde yazılmış adreslerden oluşmaktadır.

Veride en çok tekrar eden özellik her satırın sonunda il ve ilçelerin “/” karakteriyle birbirlerinden ayrılmış olmasıdır. Bunun satırdan satıra değişmekle birlikte cadde, sokak ve no, mahalle gibi bilgileri içermektedir. Nadir olsa da bazı satırlarda iş hanı, avm, bulvar gibi kısımlar da bulunmaktadır.

Problem

Temel sorunumuz bu verilerin hangilerinde hangi özelliklerin(Adres bloklarının) bulunup bulunmadığını bulmaktır. Bunu çözdükten sonra da verinin hangi kısmının bu özelliği gösterdiğini bulmaktır. Örneğin eğer bir satır içinde “mah.” bloğu geçiyorsa bu satırda mahalle bilgisi bulunuyor diyebiliriz. Bunun ardından da mahalleden önceki öbeği alarak bu öbeğin mahalle bilgisini içerdiğini söyleyebiliriz. Örneğin “23 Nisan Mah.” bu öbekte “Mah.”tan önce yer alan “23 Nisan” buradaki mahallenin adını belirtmektedir.

Yöntem

Regular expression, “<https://regex101.com/>” yardımıyla yazılmıştır. Elimizdeki veri çok fazla sayıda kombinasyon içerdiğinden verilerin bir kısmından feragat etmek durumunda kalıyoruz.

Bu kaybı minimuma indirmek için veri seti incelenip en çok tekrar eden veri gruplarını kapsayacak regular expression yazılmaya çalışılmıştır. Her satırda en sonda il ve ilçe blokları bulunmaktadır. Yazılan regex her satır için il ve ilçe bilgisini ayırtmaktadır. Satır başında ise en çok rastlanan üç seçenek bulunmaktadır bunlar:

1. Mahalle
2. Cadde
3. Sokak

Satırlar okunurken eğer bloklar bulunuyorsa bu sırada bulunacağı varsayılmıştır. Farklı sıralarda girilen adresler hata vermektedir. Bu bilgiler ayrıştırılırken satır,

((([A-z0-9ŞÜÂĞİÖşçğïöÇü\.\-]*[Mm][Aa]?[Hh][Aa]?[Ll]?[Ll]?[Ee]?[Ss]?[Ii]?[Ii]?[?].?)?)([A-z0-9ŞÜÂĞİÖşçğïöÇü\.\-]*[Bb]U?[Ll][Vv]A?R?I?[\.\.?])(["A-zÜÂçğïö-öÇĞİÖŞşü,*CA?DD?E?S?[Ii]?[\.\.?])(["A-zĞİÖŞÜÂçğö-9V,iöÇşü]*SO?[Kk]A?K?Ğ?I?[\.\.?])(.*BLOK)?([A-zÜÂçğïö-9-öÇĞİÖŞşü)*APT\.\.?)?(?.*NO?:?|[0-9]+[A-z\-\?]?(V?[0-9]*V?[A-z]?)(.*)?(?([A-zÇĞİÖçğïöŞÜâşü)*)?(V?)([0-9]{5})?)([A-zÇŞÜÂçğïĞİÖöşü]*)

	Mah	Cad	Bulvar	Sok	Blok	Apt	No	Daire	PK	Semt	Şehir
Kaç adet	513	434	22	187	4	6	494	323	5	598	598
Bulunan	502	425	19	151	3	5	456	307	1	589	597
Başarı(%)	97,9	97,9	86,37	80,07	75	83,33	92,31	95,05	20	98,5	99,83

Bulunamayan adresler;

1-)Fenerbahçe, İğrıp Sk. No:13, 34726 İstanbul

2-)KAPTAN PAŞA MAH. YAY GEÇİDİ SOK. NO:6/C BEYOĞLU

3-)Ziya Gökalp, Mahallesi, Eski Turgut Özal Cd. No:29 D:Kat.2, 34000 İkitelli
Osib/Başakşehir/İstanbul

4-)Mehmet Nesih Özmen, Fatih Cad. & Karadal Sokakı , Merter Keresteciler Sitesi
No:20, 34173 Güngören/İstanbul