



Yıldız Teknik Üniversitesi
Bilgisayar Mühendisliği Bölümü

Doğal Dil İşleme Dersi Ödev-1

Mert TÜRKMENOĞLU
17011005

KULLANILAN REGEX:

$$\begin{aligned} & /([\text{Ss}][\text{Oo}][\text{Kk}](\backslash.\backslash\text{s})^*)|([\text{Ss}][\text{Oo}][\text{Kk}][\text{Aa}][\text{Kk}](\backslash.\backslash\text{s})|([\text{Ss}][\text{Kk}](\backslash.\backslash\text{s})|([\text{Bb}][\text{Uu}][\text{Ll}][\text{Vv}][\text{Aa}][\text{Rr}][\text{Ii}]|([\text{Bb}][\text{Uu}][\text{Ll}][\text{Vv}](\backslash.\backslash\text{s})^*)|(\backslash\text{s})+[\text{Bb}][\text{Ll}][\text{Vv}](\backslash.\backslash\text{s})+)|([\text{Mm}][\text{Aa}][\text{Hh}](\backslash.\backslash\text{s})+)|([\text{Mm}][\text{Hh}](\backslash.\backslash\text{s})+)|(\backslash\text{s})+([\text{Mm}](\backslash.\backslash\text{s})|([\text{Cc}][\text{Aa}][\text{Dd}][\text{Dd}][\text{Ee}][\text{Ss}][\text{IiIi}](\backslash.\backslash\text{s})^*)|([\text{Cc}][\text{Aa}][\text{Dd}](\backslash.\backslash\text{s})^*)|([\text{Cc}][\text{Dd}](\backslash.\backslash\text{s})|([\text{Cc}](\backslash.\backslash\text{s})|([\text{Aa}][\text{Pp}][\text{Aa}][\text{Rr}][\text{Tt}][\text{Mm}][\text{Aa}][\text{Nn}][\text{Li}](\backslash.\backslash\text{s})|([\text{Aa}][\text{Pp}][\text{Tt}](\backslash.\backslash\text{s})|[\text{Nn}][\text{Oo}](\backslash.\backslash\text{s})|(\backslash\text{s})/\text{gm} \end{aligned}$$

Yukarıda verilen RegEx'i kullanarak input dosyasındaki her bir satırı ayırıp JSON formatında kaydeden Python kodu:

```

def grouper(L, n):
    args = [iter(L)] * n
    return ([e for e in t if e != None] for t in itertools.zip_longest(*args))

def get_lines(f_name):
    res = []

    with open(f_name, 'r') as f:
        for line in f:
            res.append([e.strip() for e in re.split(my_regex, line[:-1]) if not e in invalid])

    return res

def make_dict_from_line(r):
    d = {}
    splitted = re.split(r' ', r[:-2])

    for e in grouper(r[:-2], 2):
        if e[0] != '' and e[0] != ' ':
            label = ([ l for l, regex in label_regex_dict.items() if re.search(regex, " ".join(e)) != None] + ['other'])[0]
            d[label] = e[0]

    no_or_desc = " ".join(splitted[:-1]).strip()
    if no_or_desc != '':
        d['no' if re.search(r'[0-9]', no_or_desc) != None else 'desc'] = no_or_desc

    d['county'] = splitted[-1].strip()
    d['province'] = r[-1]

    return d

def write_to_file(f_name, data):
    with open(f_name, 'w') as f:
        f.write(json.dumps({'data': data}, indent=2, ensure_ascii=False))

lines = get_lines(INPUT_FILE_NAME)
result = [make_dict_from_line(line) for line in lines]
write_to_file('output.json', result)

```

output.json dosyasının bir kısmı:

```
{
  "data": [
    {
      "desc": "YENİBOSNA METRO İSTASYONU",
      "county": "BAKIRKÖY",
      "province": "İSTANBUL"
    },
    {
      "avenue": "KENNEDY",
      "desc": "SİRKEÇİ ARABALI VAPUR İSKELESİ",
      "county": "FATİH",
      "province": "İSTANBUL"
    },
    {
      "district": "YAVUZTÜRK",
      "avenue": "KARADENİZ",
      "no": "2",
      "county": "ÜSKÜDAR",
      "province": "İSTANBUL"
    },
    {
      "district": "HAMİDİYE",
      "street": "ALPEREN",
      "no": "15/2",
      "county": "ÇEKMEKÖY",
      "province": "İSTANBUL"
    },
    {
      "district": "UĞUR MUMCU",
      "avenue": "YUNUS EMRE",
      "no": "25",
      "county": "KARTAL",
      "province": "İSTANBUL"
    },
    {
      "other": "BAĞLARBAŞI",
      "avenue": "İNÖNÜ",
      "no": "3",
      "county": "MALTEPE",
      "province": "İSTANBUL"
    },
  ]
}
```

Hata oluşan durumlara örnekler:

- MIGROS ALIŞVERİŞ MERKEZİ E5 KARAYOLU ÜZERİ GIRIS KAT
BEYLİKDÜZÜ/ İSTANBUL

- CUMHURİYET MAH. ŞEHİTLER CAD. BEYLİKDÜZÜ BULVAR EVLERİ
ÇARŞISI NO:134 ESENYURT/ İSTANBUL

- MEHTERÇEŞME MAH. CUMH.CAD 1810 SOK. NO:1 ESENYURT/ İSTANBUL

Bu örneklerin JSON dosyasından alınıp birleştirilmiş hali:

```
{
  "data": [
    {
      "no": "MIGROS ALIŞVERİŞ MERKEZİ E5 KARAYOLU ÜZERİ GIRIS KAT", "county": "BEYLİKDÜZÜ",
      "province": "İSTANBUL"
    },
    {
      "district": "CUMHURİYET",
      "avenue": "ŞEHİTLER",
      "boulevard": "BEYLİKDÜZÜ",
      "other": "AR EVLERİ ÇARŞISI",
      "no": "134",
      "county": "ESENYURT",
      "province": "İSTANBUL"
    },
    {
      "district": "CU",
      "street": "1810",
      "no": "1",
      "county": "ESENYURT",
      "province": "İSTANBUL"
    }
  ]
}
```

Sonuç:

- RegEx, düzgün formatta verilen adreslerin tamamını tanımıştır.
- Yazımında farklılık / yanlışlık bulunan adreslerde bunları “other” kategorisinde değerlendirmiştir.
- Uç noktalara ilişkin hatalı örnekler yukarıda belirtilmiştir.
- Tamamiyle başarısız olduğu hiçbir veri bulunmamaktadır.