

TÜRKÇE MAKALELERİN SINIFLANDIRILMASI

MURAT DEMİR

Ö19060012

26.12.2019

K-Nearest Neighbor Algoritması Hakkında Kısa Bilgi

KNN algoritması Machine learning'deki en kolay algoritmalarından birisidir ve bir Supervised Learning algoritmasıdır. Çoğunlukla classification amacıyla kullanılır. "K" harfi komşu sayısını ifade eder yani k tane en yakın komşu algoritması olarak Türkçe'ye çevirebiliriz.

TÜRKÇE MAKALELERİN SINIFLANDIRILMASI PROJESİ

Projemiz 4 Makale Türünü Makine Öğrenmesi Yöntemlerinden K-Nearest Neighbor Algoritması ile Tasarlanmıştır. Makale türlerimiz aşağıda verilmiştir.

- Popüler Kültür
- Otomobil
- Müzik
- Ekonomi

İde olarak Eclipse kullanılarak gerçekleştirilmiştir.

Başta Zembereğin Programa Entegrasyonu Aşağıdaki gibi sağlanmıştır.

```
jvmDLLpath = r"C:\Program Files\Java\jdk-13\bin\server\jvm.dll"
jpyype.startJVM(jvmDLLpath, "-Djava.class.path=zemberek-tum-2.0.jar", "-ea")

TR = jpyype.JClass("net.zemberek.tr.yapi.TurkiyeTurkcesi")
tr = TR()

Z = jpyype.JClass("net.zemberek.erisim.Zemberek")
z = Z(tr)
```

Makalelerde geçen kelimelerin köküne inilip stopwords 'lerden kurtarılmıştır.

```
def Zemberek(makale):
    yenido = []
    makale = makale.split()
    for line in makale:
        try:
            kok = z.kelimeCozumle(line)[0].kok().icerik()
            if kok:
                yenido.append(kok)
        except:
            pass
    return yenido

egitim = []
file = open("Makaleler.txt")
for satir in file:
    satir = satir.strip()
    temizle = Zemberek(satir)
    temizle = ' '.join(temizle)
    egitim.append(temizle)
```

Eğitim Makaleler.txt belgesinin okunup ilk 57 satır içerisinde bulunan makalenin 0. Index değerine ,57-119. Satırlar arasında kalan makalenin 1. Index değerine, 119-159. Satırlar arasında kalan makalenin 2. Index değerine ve 159-218. Satırlar arasında kalan makalenin 3. Index değerine atanması sağlanıp sırasıyla

0.index = Popüler Kültür

1.index = Otomobil

2.index = Müzik

3.index = Ekonomi

Alanlarına atanıp KNN algoritması ile eğitilmiştir.

```
vectorizer = TfidfVectorizer(min_df = 0., max_df = 1., use_idf = True)
egitim = vectorizer.fit_transform(egitim)

egitim_y = np.zeros(218)
egitim_y[57:119] = 1
egitim_y[119:159] = 2
egitim_y[159:218] = 3

mkneighbors = KNeighborsClassifier(n_neighbors=15)
mkneighbors.fit(egitim,egitim_y)
DogruEtiketler = ['Popüler Kültür','Otomobil', 'Müzik/Piyano', 'Ekonomi']
```

Ardından TestDokumanlari.txt Dökümanına gireceğimiz Algoritmanın daha önce görmediği yeni makalenin hangi sınıfa dahil edileceği K-Nearest Neighbor algoritması ile hesaplanmış ve programın sağlıklı bir şekilde sonlandırılması sağlanmıştır

```
Dokumani_Temizle = []
for line in testDokumani:
    line = line.strip()
    temizle = Zemberek(line)
    temizle = ' '.join(temizle)
    Dokumani_Temizle.append(temizle)
Test_Dokumani = vectorizer.transform(Dokumani_Temizle)

YuklemEtiketleriKneighbors = mkneighbors.predict(Test_Dokumani)

YuklemEtiketleriKneighbors = YuklemEtiketleriKneighbors.astype(int)
YuklemEtiketleriKneighbors = ''.join(str(v) for v in YuklemEtiketleriKneighbors)
sayac = 0
for i in YuklemEtiketleriKneighbors:
    MevcutSayac = YuklemEtiketleriKneighbors.count(i)
    if(MevcutSayac > sayac):
        sayac = MevcutSayac
        Etiket = i

print("\nDokümanın sinifi: ", DogruEtiketler[np.int(Etiket)], "\n\n\n")

jpye.shutdownJVM()
```