# Klasik Makine Ogrenmesi Yontemleri ile Mesleklerin Siniflandirilmasi(kariyer.net)

Havvanur Dervişoğlu,drvshavva@gmail.com August 29, 2019

# 1 Çalışmanın Tanımı ve Kullanılan Araçlar

Bu calismadaki adimlar:

- Web scraping ile veri setinin olusturulmasi.
- Olusturulan veri setine preprocessing uygulanmasi.
- On islem asamasından gecmis veri setindeki is tanimi ozelliğine feature extraction uygulanması. Calisma da tf-idf ve count vectorizer oznitelik cikarimi yontemleri farkli parametrelerle denenmis olup siniflandirma basirilari kaydedilmistir.
- Oznitelik cikarimi asamasindan sonra klasik makine ogrenmesi yontemleri ile siniflandirma basarilarinin sonuclari alinmistir. Klasik makine ogrenmesi yontemlerinden; logistic regression, random forest, decision tree classifier ve naive bayes siniflandiricilari kullanilmistir.

### 1.1 Web Scraping icin Kullanilan Araclar

- Calisma da web scraping icin PyCharm calisma ortami kullanilmistir.
- Programa dili olarak python kullanilmistir.
- Kariyer.net'ten veri cekmek icin selenium ve scrapy kutuphaneleri kullanilmistir.

#### 1.2 Siniflandirma icin Kullanilan Araclar

Veri setinde uygulanan preprocessing islemi ich google colab kullanilmistir. Mesleklerin siniflandirilasi icin Pyspark kullanilmistir ve Zeppelin gelistirme ortami kullanilmistir.

# 2 Web Scraping

Kariyer.net web sitesinde veri cekmek icin scrapy ve selenium kutuphaneleri kullanilmistir. Kodda sinif ilan aramalari sinif bazinda yapilmaktadir. Figure 1'de, olusturulacak sinifin adi sinif degiskenine atanmistir, sonuclarin yazilacagi dosya adi ve tipi beklirtilmistir.

Figure 1: Web scraping-1(kariyer.net)

Figure 2'de, selenium kutuphanesinin kullanimi icin gerekli olan web driverin oldugu dizin belirtilmistir ve sonrasinda acilan web sitesinde aramanin yapildigi yere gidilip sinif adi aratilmistir.

```
def parse(self, response):
    #bilgisayarinizda chromedriver'in oldugu dizini burada belirtiyoruz
    webdriver_path = '/home/safir/Desktop/chromedriver'
    self_driver = webdriver.Chrome(webdriver_path)
    self_driver.get('https://www.kariyer.net/is-ilanlari')

#sayfada arama yapilan yere kategori ismini yaziyoruz
    search_item = self_driver.find_element_by_xpath(' // *[ @ id = "txtSearchKeyword"]')
    search_item.clear()
    search_item.send_keys(IlanlarSpider.sinif)
    sleep(1.8)
    #kategori ismini yazdiktan sonra butona tikliyoruz
    next_page = self_driver.find_element_by_xpath('//*[@id="btnSearchKeyword"]')
    self_driver.execute_script("arguments[0].click();", next_page)
    sleep(4.5)

items = KariyerItem()
    #next_page = 'https://www.kariyer.net/is-ilanlari'
    ilan_links = []
    page = 0
    #23 sayfadan ilanlara ait linkleri aliyoruz
```

Figure 2: Web scraping-2(kariyer.net)

Figure 3'de, arama sonucunda cikan ilanlari sayfa sayfa ilerleyerek toplamaktadir.

```
#23 sayfadan ilanlara ait linkleri aliyoruz
while page <= 23:
    scrapy_selector = Selector(text=self_driver.page_source)
    links = scrapy_selector.xpath('//a[@class="link position"]/@href').extract()
    for i in range(0, len(links)):
        ilan_links.append(links[i])

#sonraki sayfaya geciyoruz
    next_page = self_driver.find_element_by_xpath('//a[@id="lnkNextPage"]')
    self_driver.execute_script("arguments[0].click();", next_page)
    page += 1
    sleep(4.8)</pre>
```

Figure 3: Web scraping-3(kariyer.net)

Figure 4'de, toplanan ilan linklerine gidilerek ilana dair istenilen ozellikler alinmaktadir.

```
#topladigimiz ilan linklerine gidip belirtilen ozellikleri aliyoruz
for i in range(0, len(ilan_links)):
    url = 'https://www.kariyer.net' + ilan_links[i]
    self_driver.get(url)
    scrapy_selector = Selector(text=self_driver.page_source)
    genel = scrapy_selector.xpath("""/h3[contains(., 'GENEL NITELİKLER VE İŞ TANIMI')]/following-sibling::node()//text()""").extract()
    if genel:
        items['ilan_baslik'] = scrapy_selector.xpath('//a[@id="jobTitle"]/text()').extract()
        items['sirket_adi'] = scrapy_selector.xpath('//a[@id="jobCompany"]/text()').extract()
        items['genel_nit_is_tanimi'] = genel
    #items['is_tanimi'] = scrapy_selector.xpath("""//div/h3[starts-with(., 'IŞ TANIMI')]/following-sibling::node()/descendant-or-self::t
        items['genel_nit_is_tanimi'] = scrapy_selector.xpath('//div[@class="sub-box aday-kriterleri"]/div[2]/div[1]/div[2]/p/text()').extract()
    items['equitim'] = scrapy_selector.xpath('//div[@class="sub-box aday-kriterleri"]/div[2]/div[a]/div[2]/p/text()').extract()
    items['gitim'] = scrapy_selector.xpath('//div[@class="sub-box aday-kriterleri"]/div[2]/div[a]/div[2]/p/text()').extract()
    items['sektor'] = scrapy_selector.xpath('//div[@class="sub-box aday-kriterleri"]/div[2]/div[a]/div[2]/p/text()').extract()
    items['sektor'] = scrapy_selector.xpath('//div[@class="sub-box pozisyon-bilgileri"]/div[2]/div[2]/p/text()').extract()
    items['claisma_sekli'] = scrapy_selector.xpath('//div[@class="sub-box pozisyon-bilgileri"]/div[2]/div[2]/p/text()').extract()
    items['sehir'] = scrapy_selector.xpath('//div[@class="sub-box pozisyon-bilgileri"]/div[2]/div[2]/p/text()').extract()
    items['sehir'] = scrapy_selector.xpath('//div[@class="sub-box pozisyon-bilgileri"]/div[2]/div[2]/p/text()').extract()
    items['sehir'] = scrapy_selector.xpath('//div[@class="sub-box pozisyon-bilgileri"]/div[2]/div[2]/p/text()').extract()
    items['sehir'] = scrapy_selector.xpath('/div[@class="sub-box pozisyon-bilgileri"]/div[2]/div[a]/p/text()').extract()
    items['sehir
```

Figure 4: Web scraping-4(kariyer.net)

# 3 Veri Seti

Calismada iki farkli veri seti kullanilmistir:

- 1. Kariyer.net web sitesinden genel 5 meslek grubundan olusan turkce veri seti olusturulmustur
- 2. Kariyer.net web sitesinden birbirlerileri ile koreslasyonlari yuksek olan 10 meslekten olusan veri seti olusturulmustur.

#### Amac:

• Turkce veri setinde meslekler arasi korelasyon arttiginda siniflandirma basarisinin nasil degistigini gozlemlemek.

### 3.1 Kariyer.net Veri Seti-1

Bu veri setinde ilanlara ait kullanilan ozellikler; is tanim, label, ilan basligi, dil, egitim seviyesi, sirket, sektor, sehir seklindedir. Siniflardaki ornek sayilari asagida gosterildigi gibidir:

1. banka-sigorta: 954 ilan

2. egitim-ogretim: 822 ilan

3. bilisim-telekom: 768 ilan

4. saglik : 814 ilan

5. yapi-mimar-insaat: 953 ilan

Veri seti olustururken siniflardaki ornek sayilarinin dengeli olmasina dikkat edilmistir.

#### 3.2 Kariyer.net Veri Seti-2

Bu veri setinde ilanlara ait kullanilan ozellikler; is tanim, label, ilan basligi, dil, egitim seviyesi, sirket, sektor, sehir seklindedir. Siniflardaki ornek sayilari asagida gosterildigi gibidir:

1. acente-sigorta-danisman: 81 ilan

2. avukat: 88 ilan

3. hasta-hizmetleri: 93 ilan

4. is-analizi-raporlama : 68 ilan

5. is-guvenlik-saglik: 67 ilan

6. proje-yonetimi: 105 ilan

7. sap: 75 ilan

8. sistem-yonetimi: 76 ilan

9. tibbi-tanitim : 66 ilan

10. web-tasarim: 74 ilan

Veri seti olustururken siniflardaki ornek sayilarinin dengeli olmasina dikkat edilmistir.

# 4 Preprocessing

On islem asamasi siniflandirma yapmadan once verilerden anlamli olan bilgileri cikarmak icin onemli bir adimdir; oncelikle Figure 5'de de goruldugu gibi meslek siniflandirilmasinda kullanilmayacak olan ozellikleri attik.



Figure 5: Preprocess-1(kariyer.net)

Figure 6'da, veri setinde uygulanan islemler gorulmektedir:

- 1. Harfleri kucuk harfe cevirme
- 2. Turkce karakter harfleri ingilizce karaktere cevirme
- 3. Noktalama isaretleri vb isaretleri atma
- 4. Rakamlari silme
- 5. Cumleyi kelimelerine ayirip stop words olup olmadigini kontrol etme ve eger stop word ise almama
- 6. Kelime uzunlugu 3 den kucuk olanlari alma(amac kesme isareti ile ayrilmis kelimelerden kurtulma nin nin gibi)
- 7. Kelimler arasi bir bosluk olacak sekilde birlestirme.

```
def clean(text):
    processed_tweet = []
    text=text.lower()

text=text.replace('$','s')
    text=text.replace('i','i')
    text=text.replace('i','u')
    text=text.replace('i','u')
    text=text.replace('g','g')
    text=text.replace('g','c')

pattern = r"[{}}".format('&+#*.-'"•\'"?!,.():;></-')
    text = re.sub(pattern, " ", text)

text = re.sub(r'[0-9]+', ' ', text)
    text = text.strip()
    #burada eger basta i varsa onu kelimeden ayriyor istanbul => i, stanbul
    #soralim stopwords kullanimini ve usteki
    tokens = WPT.tokenize(text)
    filtered_tokens = [token for token in tokens if token not in stop_word_list
    filtered_tokens = [token for token in tokens if len(token)>3 ]
    text = ''.join(filtered_tokens)
    return text
```

Figure 6: Preprocess-2(kariyer.net)

Figure 7'de, veri setinin son gorunumunun csv dosyasina kaydedilmesi gorulmektedir.

```
clean_df = pd.DataFrame(clean_is_tanim,columns=['genel_nit_is_tanimi'])
clean_df['ilan_baslik'] = clean_ilan_baslik
clean_df['label'] = data.label
clean_df.to_csv('cleaned_web_tasarim.csv',encoding='utf-8')
clean_df.head()
                                                                       ilan_baslik
                               genel_nit_is_tanimi
                                                                                                  label
         universitelerin bilgisayar muhendisligi bilgis...
                                                                      front developer
                                                                                          web tasarim
 1 hedefini dunyanin dort tarafında hayata gecirm...
                                                                       yazilim uzmani
                                                                                           web tasarim
     universitelerin endustri muhendisligi sletme m...
                                                                                           web_tasarim
 3 hedefini dunyanin dort tarafında hayata gecirm... junior yazilim uzmani
                                                                                           web tasarim
       tanimi html taslaklari gereksinimleri gercekle...
                                                                    gelistirme uzmani
                                                                                          web tasarim
```

Figure 7: Preprocess-3(kariyer.net)

Figure 8'de, olusan veri setinde labellara id verilmesi gosterilmistir. Her label icin bir id atamasi yapilmistir.

```
%pyspark
#BURADA LABEL SUTUNUDA YER ALAN SINIFLARA INDEX NUMARASI VERILIYOR
indexer = StringIndexer(inputCol = "label", outputCol = "labelIdx")
df = indexer.fit(df).transform(df)
df.show()
```

Figure 8: Preprocess-4(kariyer.net)

#### 5 Feature Extraction

Veri setlerinde iki farkli oznitelik cikarimi yontemi kullanilmistir. Oznitelik cikarimi uygulanmadan once kariyer.net sitesinden elde edilen veri setindeki ilan tanimi ve basligi ozellikleri tek bir sutunda birlestirilmistir ve bu sutuna oznitelik cikarimi islemi uygulanmistir(Figure 9).

<pre>df = df.withColumn('tum_o df.show()</pre>		ol('genel_ni	t_is_tani	mi'),lit('_'), col
genel_nit_is_tanimi	ilan_baslik	label lab	elIdx	tum_ozellikler
universitelerin b				
hedefini dunyanin	yazilim uzmani web	tasarim	6.0 hede	efini dunyanin
universitelerin e	anali web_	tasarim	6.0 uni	versitelerin e
hedefini dunyanin juni	or yazilim uz web_	tasarim	6.0 hede	efini dunyanin
tanimi html tasla  g	elistirme uzmani web_	tasarim	6.0 tan	imi html tasla
danismanligini ya onyu	z gelistirme web_	tasarim	6.0 dan	ismanligini ya
front developer g onyu	z front yazil web_	tasarim	6.0 from	nt developer g
universitelerin i fron	t developer a web_	_tasarim	6.0 uni	versitelerin i
universitelerin i	senior developer web_	tasarim	6.0 uni	versitelerin i
front development  fr	ontend developer web_	tasarim	6.0 from	nt development
bilsoft yazilim y	yazilim uzmani web_	tasarim	6.0 bil	soft yazilim y
consulting verini	designer web_	tasarim	6.0 con	sulting verini
sisli osmanbey bu	tasarim stajyeri web_	tasarim	6.0 sis	li osmanbey bu
sirketimizde cali	tasarim uzmani web	tasarim	6.0 sir	ketimizde cali
teccuhali html is Ifcon	tend varilim   luch	tacaciml	6 Alter	ruhali html is I

Figure 9: Ozelliklerin birlestirilmesi(kariyer.net)

Figure 10'da, oznitelik islemi uygulanilacak olan sutundaki cumlelerin kelimlerine ayrilmasi islemi gosterilmistir. Bu islem oznitelik cikarimi yontemlerinin o sutunda kullanilmasi icindir.

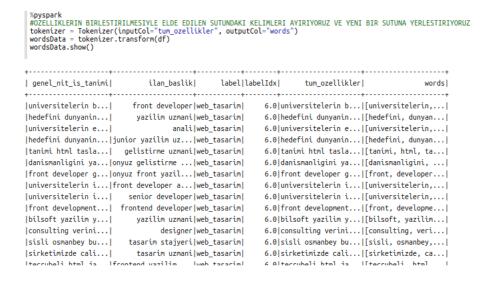


Figure 10: Cumlenin kelimelerine ayrilmasi

Figure 11'de,tf-idf vectorizer oznitelik cikarimi yonteminin kullanim ornegi bulunmaktadir.

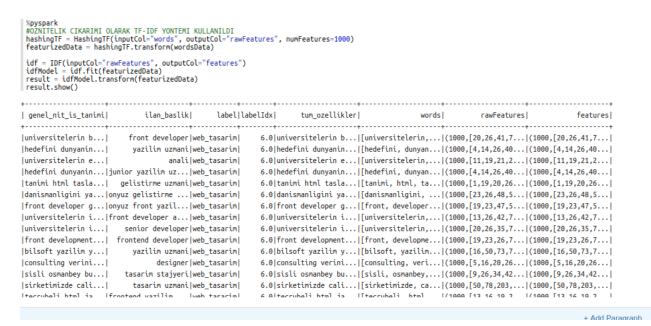


Figure 11: TF-IDF vectorizer

Figure 12'de, count vectorizer oznitelik cikarimi yonteminin kullanim ornegi bulunmaktadir.

```
%pyspark
#0ZNITELIK CIKARIMI OLARAK COUNT_VECTORIZER YONTEMI KULLANILDI
cv = CountVectorizer(inputCol="words", outputCol="features", vocabSize=10000, minDF=5)
model = cv.fit(wordsData)
result = model.transform(wordsData)
result.show()
-----
 genel_nit_is_tanimi|
                        ilan_baslik| label|labelIdx| tum_ozellikler|
                                                                                       words
                                                                                                        features
  ------
universitelerin b...| front developer|web_tasarim|
                                                     6.0|universitelerin b...|[universitelerin....|(2561,[0.1.2.3.4....|
hedefini dunvanin...
                      yazilim uzmani|web_tasarim|
                                                     6.0|hedefini dunvanin...|[hedefini. dunvan...|(2561.[0.1.3.4.5....
universitelerin e...|
                               anali|web_tasarim|
                                                     6.0|universitelerin e...|[universitelerin,...|(2561,[0,1,2,5,8,...|
hedefini dunyanin...|junior yazilim uz...|web_tasarim|
                                                     6.0|hedefini dunyanin...|[hedefini, dunyan...|(2561,[0,1,3,4,6,...
tanimi html tasla...| gelistirme uzmani|web_tasarim|
                                                     6.0|tanimi html tasla...|[tanimi, html, ta...|(2561,[0,1,4,5,7,...|
danismanligini ya...|onyuz gelistirme ...|web_tasarim|
                                                     6.0|danismanligini ya...|[danismanligini, ...|(2561,[0,1,3,4,5,...
front developer g...|onyuz front yazil...|web_tasarim|
                                                     6.0|front developer g...|[front, developer...|(2561,[0,4,5,6,8,...|
universitelerin i...|front developer a...|web_tasarim|
                                                     6.0|universitelerin i...|[universitelerin,...|(2561,[0,2,3,4,5,...|
universitelerin i...|
                     senior developer|web_tasarim|
                                                     6.0|universitelerin i...|[universitelerin,...|(2561,[2,3,4,5,6,...|
front development...| frontend developer|web_tasarim|
                                                     6.0|front development...|[front, developme...|(2561,[1,2,3,4,5,...|
bilsoft yazilim y...|
                      yazilim uzmani|web_tasarim|
                                                     6.0|bilsoft yazilim y...|[bilsoft, yazilim...|(2561,[0,1,4,6,7,...|
                                                     6.0|consulting verini...|[consulting, veri...|(2561,[1,2,9,10,1...|
consulting verini...
                            designer|web tasarim|
                                                     6.0|sisli osmanbey bu...|[sisli, osmanbey....|(2561,[0,1,2,3,4,...|
sisli osmanbey bu...|
                     tasarim stajveri|web tasarim|
sirketimizde cali...|
                      tasarim uzmanilweb tasariml
                                                     6.0|sirketimizde cali...|[sirketimizde, ca...|(2561,[32,91,182....|
A Alterruhali html is | [[terruhali html
```

Figure 12: Count vectorizer

#### 6 Classification

Siniflandirma icin kullanilacak olan ozellikler olusan features sutunu ve labelIdx sutunu(Figure 13). Bu calismada siniflandirma algoritmalarindan logistic regression, decision tree, random forest ve naive bayes kullanilmistir.

```
%pyspark
featured_data = result.select('labelIdx','features')
featured_data.show()
.....
|labelIdx|
               features
     6.0|(1000,[20,26,41,7...|
     6.0|(1000,[4,14,26,40...|
     6.0|(1000,[11,19,21,2...|
     6.0|(1000,[4,14,26,40...|
     6.0|(1000,[1,19,20,26...|
     6.0|(1000,[23,26,48,5...|
     6.0|(1000,[19,23,47,5...|
     6.0|(1000,[13,26,42,7...|
     6.0|(1000,[20,26,35,7...|
     6.0|(1000,[19,23,26,7...|
     6.0|(1000,[16,50,73,7...|
     6.0|(1000,[5,16,20,26...|
     6.0|(1000,[9,26,34,42...|
     6.0|(1000,[50,78,203,...|
```

Figure 13: Siniflandirma icin kullanilan ozellikler

Figure 14'de goruldugu gibi veri seti; train %70 ve test %30 olacak sekilde bolunmustur.

```
%pyspark
#VERI SETINI BOLUYORUZ
splits = featured_data.randomSplit([0.7, 0.3])
train = splits[0]
test = splits[1]
train_rows = train.count()
test_rows = test.count()
print "Training Rows:", train_rows, " Testing Rows:", test_rows
```

Figure 14: Veri setinin bolunmesi

Figure 15'de decision tree classifier algoritmasinin kullanimi gosterilmistir.

```
%pyspark
#SINIFLANDIRMA-1
dt = DecisionTreeClassifier(labelCol="label", featuresCol="features")
dt_model = dt.fit(train)
```

Figure 15: Decision tree classifier kullanimi

Figure 16'da logistic regression algoritmasinin kullanimi gosterilmistir.

```
%pyspark
#SINIFLANDIRMA-2
lr = LogisticRegression(maxIter=20, regParam=0.3, elasticNetParam=0,labelCol="label", featuresCol="features")
lrModel = lr.fit(train)
lr_predictions = lrModel.transform(test)
Took 1 sec Last updated by admin at August 26 2019, 3:16:54 PM.
```

Figure 16: Logistic regression kullanimi

Figure 17'de naive bayes algoritmasinin kullanimi gosterilmistir.

```
%pyspark
#SINIFLANDIRMA-3
from pyspark.ml.classification import NaiveBayes
nb = NaiveBayes(smoothing=1)
model = nb.fit(train)
predictions = model.transform(test)
```

Figure 17: Naive bayes algoritmasinin kullanimi

Figure 18'de decision tree algoritmasinin kullanimi gosterilmistir.

Figure 18: Decision tree algoritmasinin kullanimi

Figure 19'da, siniflandirma algoritmlarina ait dogruluk oraninin nasil alindigi gosterilmistir.

```
%pyspark
from pyspark.ml.evaluation import MulticlassClassificationEvaluator
evaluator = MulticlassClassificationEvaluator(predictionCol="prediction")
evaluator.evaluate(lr_predictions)
0.8841856466100212
```

Figure 19: Siniflandirma algoritmalarinin basari sonucunu alma

#### 7 Sonuclar

Bu bolumde calisma sonucunda elde edilen siniflandirma basarilarinin degerleri verilmistir. Tablo 1'de count vectorizer oznitelik cikarimi yonteminin farkli numWords degerlerindeki dogruluk oranlari verilmistir (sonuclar minDF=5 ve kariyer.net 5 sinifli veri seti icindir).

Table 1: Count vectorizer kariyer.net 5 sinifli veri seti dogruluk oranlari minDF=5 icin

Model	10000	1000	500	100
LR	92 %	90 %	87 %	80 %
RF	85 %	83 %	81 %	79 %
NB	90 %	88%	85%	81 %
DTC	73 %	71 %	71 %	70 %

Table 2'de, count vectorizer oznitelik cikarimi yonteminin iki farkli numWords degerinde ve minDF=2 alinmis sonuclari bulunmaktadir (Sonuclar kariyer.net 5 sinifli veri seti icindir).

Table 2: Count vectorizer kariyer.net 5 sinifli veri seti dogruluk oranlari minDF=2 icin

Model	1000	100
LR	90 %	78 %
RF	81 %	78 %
NB	88%	81 %
DTC	72 %	72 %

Table 3'de, tf-idf oznitelik cikarimi yonteminin farkli numFeatures degerlerindeki sini-flandirma basarilari verilmistir(Sonuclar kariyer.net 5 sinifli veri seti icindir).

Table 3: Tf-idf kariyer.net 5 sinifli veri seti dogruluk oranlari

Model	1000	500	100
DTC	67 %	65 %	52 %
LR	89 %	86 %	74 %
NB	86 %	83%	73 %
RF	79 %	75 %	65 %

Tablo 4'de count vectorizer oznitelik cikarimi yonteminin farkli numWords degerlerindeki dogruluk oranlari verilmistir(sonuclar minDF=5 ve kariyer.net 10 sinifli veri seti icindir).

Table 4: Count vectorizer kariyer.net 10 sinifli veri seti dogruluk oranlari minDF=5 icin

Model	10000	1000	500	100
DTC	52 %	51 %	51 %	48 %
LR	88 %	83 %	82 %	72 %
NB	87 %	82%	85%	76 %
RF	80 %	80 %	84 %	72 %

Table 5'de, tf-idf oznitelik cikarimi yonteminin farkli numFeatures degerlerindeki siniflandirma basarilari verilmistir(Sonuclar kariyer.net 10 sinifli veri seti icindir).

Table 5: Tf-idf kariyer.net 10 sinifli veri seti dogruluk oranlari

Model	1000	500	100
DTC	39 %	36 %	32 %
LR	81 %	75 %	66 %
NB	81 %	74%	61 %
RF	62 %	56 %	53 %