



Öğrenci Adı ve Soyadı	
Öğrenci Numarası	
Ders Adı	Doğal Dil İşlemeye Giriş
Sınav Tarihi ve saati	6.12.2024 11.00
Sınav Süresi	50 dakika
Sınav Salonu	E1-31, E1-41, E1-51
Öğretim Elemanı Adı	Dr.Öğr.Üyesi Zeynep Altan

Soru Numarası	1	2	3	4	5	6	Toplam
Alınan Not							

1) Aşağıdaki düzgün ifadelerin her birinin hangi dizgilerle eşleşebileceğini { } içerisinde yazın. Her bir dizgiyi nasıl elde ettiğinizi açık olarak yazın. Sadece cevap değerlendirilmez (20 puan)

i) [Aa] (n (in)? | yem?)

{Anın, An, anın, an, Aye, aye, Ayem, ...}

'A' yada 'a' den sonra 'n' gelir. Sonrasında 'in' - ayem  
'A' ya da 'a' den sonra 'y' 'e' 'i' mi gelir. Sonrasında 'm' olur

ii) [Aa] kil \\*1\*

{Akıl\*\*1, akıl\*\*1}

'k' anlamı 'den sonra ne varsa o zembelen gelmesidir. Burada \* sembolü 2 kere gelecektir

iii) A(ra)+ba

{Araba, Araraba, Arara...raba}

'ra' 'ini' bir ya da daha fazla sayıda  
tekrar eder. Öncesinde 'A' - Sonrasında 'ba' vardı

iv) \d+(\.ld+)?(lw+)

{15.49 kod 15, 25 kod 25}

'ld' herhengi bir desimaldir. '+' bir veya daha fazla payda desimaldir  
'lw+' ifadesi 'ya' vardır 'ya da' anlamı. 'ini' bir veya daha fazla desimaldir

2) Sari | sari düzgün ifadesi aynı dizgi ya da dizgilerle eşleşecek şekilde verilen farklı iki şekilde daha ifade edilebilir. Bu düzgün ifadeler nelerdir? Hangi dizgiler eşleşir? Cevabınızı mutlaka açıklayarak vermelisiniz (10 puan).

[Ss]ari

(S|s)ari

her ibisi de S ya da s anlamına eladır S ya da s

Düzgün ifade, Sari veya sari dizgileri ile eşleşir

3) i) Satır başında (dizginin başında) alfanümerik karakterlerden oluşan bir sıralanışı herhangi bir tamsayı izlemektedir ve satır (dizgi) sonlanmaktadır. Bu düzgün ifadeyi 2 farklı şekilde nasıl yazarsınız? Yaptığınız işlemleri açıklamanız gerekir. Bu düzgün ifade ile eşleşen 2 farklı dizgi yazın (10 puan).

$\wedge [a-zA-Z0-9-] + [0-9] + \$$  ya da  $\wedge \wedge [a-zA-Z0-9] * (0-9) + \$$   
 1. dizgi 15 kod / 7 ve -All 45 eşleşen dizgilerdir

ii) i) deki düzgün ifadeyi olumsuz olarak yazdığımızda eşleşen dizgileriniz neler olacaktır? Cevap için düzgün ifadenin de yazılması zorunludur. Yapılan işlemler açıklmalıdır (10 puan)

$\wedge [ \wedge \wedge [a-zA-Z0-9] * (0-9) + ] \$$  ya da  $\wedge [ \wedge (a-zA-Z0-9) * (0-9) + ] \$$

Klavyedeki özel karakterlerin bir dizilişini (örneğin !+), klavyedeki rakamların bir dizilişini (örneğin 1+2+3) dışındaki karakterler, örneğin za'lar  $\Rightarrow$  !+2+3

4) Çalışmaya ait derlem örneği "01" "12" "123" "234" olarak verilmektedir. Bu derleme byte pair encoding (BPE) belirteç(token) öğrenici algoritmasını 2 kere uyguladığımızda sonucunuz ne olur? Bu algoritma ne amaçla kullanılmaktadır? (Belirteç sonu imini hesaplamalarınıza eklemeyiniz) (15 puan).

01	0	1 defa
12	1	3 defa
123	2	3 defa
234	3	2 defa
	4	1 defa
	12	2 defa
	23	2 defa

sözleşme = {0, 1, 2, 3, 4}  
 sözleşme {0, 1, 2, 3, 4, 12, 23, 111  
 12 birleştirilir ve sözleşme eklenir  
 23 birleştirilir ve sözleşme eklenir  
 Algoritmanın 2 kere uygulanması isteniyor

$$\hat{P}(w_i | c_j) = \frac{\text{count}(w_i, c_j)}{\sum_{w \in V} \text{count}(w, c_j)}$$

Sözleşme listeli derlemelerden derlemeyi önlemek için sözleşme orijinalinden farklı bir yeri elemanlar (öğeler) eklenir.

5) Çok terimli Naive Bayes algoritması ile ne ifade edilmektedir? Açıklayınız. (10 puan).

Maximum Likelihood Estimation  $\rightarrow$  simülasyon  $\rightarrow$  Çok terimli bir Naive Bayes algoritmasıdır.

$$\hat{P}(w_i | c_j) = \frac{\text{count}(w_i, c_j)}{\sum_{w \in V} \text{count}(w, c_j)}$$

Orijinali Markov varsayımıdır. Sınırdaki sözleşmenin gelme olasılığıdır n-gram modeller olarak gerçekleştirilir. Metin sınıflandırıcıdır.

Bir  $c_j$  sınıflandırıcısında  $w_i$  sözleşmenin olma olasılığı hesaplanır. Sonuçta sayma ile ulaşılmıştır. Tüm dokümanlarda elde edilen megal doküman elde edilir.

$w_i$  sözleşmenin  $c_j$  sınıfında (konusunda) kaç kez görüldüğü pay ifadesidir. Payda ise, sözleşmeli (v) yani tüm sözleşme topluğunun birleşiminde ki sözleşmelerin  $c_j$  sınıfında (konusunda)