

CSE 484 NATURAL LANGUAGE PROCESSING
HOMEWORK 3: TURKISH CLASSIFIER

ÇAĞRI ÇAYCI
1901042629

1. ABSTRACT

Grammar errors are among the most common mistakes, and they can sometimes lead to the incorrect interpretation of a sentence. One significant example of this occurs with the suffixes '-de' and '-ki' in Turkish. The '-de' suffix in Turkish serves to indicate place or functions as a conjunction. The distinction between these two uses lies in their formatting: in the former case, it is written adjacent to the preceding word, while in the latter case, it is written separately. This project aims to develop a neural network classifier that determines the appropriate placement of '-de' and '-ki' suffixes within sentences. Building a neural network classifier for this grammar rule involves six key steps: data collection, data preprocessing, model definition, training, evaluation, and testing.

2. DATA COLLECTION

To construct a classifier, a substantial dataset must be gathered. After researching available options, the Wikipedia Turkish dump emerged as the most suitable dataset for this purpose. It comprises 4.5 million lines, with each line containing multiple sentences. This vast dataset provides a sufficient number of data samples for this project.

3. DATA PREPROCESSING

a. ADDING SENTENCES

Once all sentences are extracted from the Wikipedia dump, they need to undergo preprocessing before being utilized. Each sentence must meet exactly one of the following conditions to be added to our dataset:

- The sentence contains only one word with the '-de' suffix.
- The sentence contains 'de' solely as a conjunction.
- The sentence contains only one word with the '-ki' suffix.
- The sentence contains 'ki' solely as a conjunction.

If the sentence meets the second or fourth criteria, the conjunction must be merged with the previous word before adding it to the dataset.

```

if re.search(r'\bW*ki\b', sentence) and re.search(r'\bW*de\b', sentence): # If sentence does not contain either '-de' and 'ki', continue.
    continue
elif not re.search(r'\bW*ki\b', sentence) and not re.search(r'\bW*de\b', sentence): # If sentence contains both '-de' and 'ki', continue.
    continue
if re.search(r'\bW*de\b', sentence): # If sentence contains adjacent '-de', append 1 to Y.
    Y_previous.append(1)
elif re.search(r'\sde\s', sentence): # If sentence contains separated '-de', append 0 to Y and merge '-de' with previous word.
    words = sentence.split()
    for i in range(len(words)):
        if(words[i] == 'de'):
            words[i-1] += 'de'
            words.pop(i)
            break
    sentence = ' '.join(words)
    Y_previous.append(0)
elif re.search(r'\bW*ki\b', sentence): # If sentence contains adjacent '-ki', append 1 to Y.
    Y_previous.append(1)
elif re.search(r'\sksi\s', sentence): # If sentence contains separated '-ki', append 0 to Y and merge '-ki' with previous word.
    Y_previous.append(0)
    words = sentence.split()
    for i in range(len(words)):
        if(words[i] == 'ki'):
            words[i-1] += 'ki'
            words.pop(i)
            break
    sentence = ' '.join(words)
else:
    continue
X_previous.append(sentence) # Add sentence to X.

```

(Figure 1)

b. PREPARING DATASET

It is considered that having 50,000 data samples for each class (separated and not separated) is sufficient to build a classifier. Therefore, after the sentences are added to the dataset, only 100,000 data samples were retained.

```

count_p = 0
count_n = 0
Y = []
X = []
for i in range(len(X_previous)): # Make X contains same number (50K) of sample of each class.
    if(count_p == 50000 and count_n == 50000):
        break
    elif(count_p == 50000 and Y_previous[i] == 1):
        continue
    elif(count_n == 50000 and Y_previous[i] == 0):
        continue
    else:
        if(Y_previous[i] == 1):
            count_p += 1
        else:
            count_n += 1
        X.append(X_previous[i])
        Y.append(Y_previous[i])

```

(Figure 2)

```

[{"doc_id": "18" url="https://tr.wikipedia.org/wiki?curid=18" title="Cengiz Han"}
"Yüzyılın başında Orta Asya'daki tüm göçebe bozkır kavimlerini birleştirerek bir ulus haline getirdi ve o ulusu "Moğol" siyasi kimliği çatısı altında topladı."
"Bozkır geleneğinden gelen onlu teşkilatı kullanarak Meritokratik (liyakata bağlı) bir ordu meydana getiren Cengiz Han'ın büyük bir asker olarak ün kazanmasının temelinde, kurduğu posta teşkilatı ve casus ağı ile istihbarat sanartına verdiği b
...
"Dört büyüklerin oynadığı tek TSYD Kupası olması nedeniyle farklılık arzeden 1975 Turnuvasında Galatasaray'ın Beşiktaş'ı 5-1 yendiği gün Fenerbahçe'de kısa süre önce 2-0 yenildiği Trabzonspor'dan rövanşı aldı."
"Son maçta da Beşiktaş'ı üstün bir oyunda 2-0 yenen sarı-lacivertliler 4 büyükünde katıldığı bu turnuvayı 3'te 3 yaparak şampiyon tamamlandı ve tarihi bir başarı kazandı."
"Fenerbahçe bu sonuçla bu tarihte maçta kupayı kazanan takım olurken, son 28 maçta aldığı 18 galibiyet ve 7 beraberliğe bir galibiyet daha ekleyerek ezici üstünlüğünde sürdürdü.")

```

(Figure 3)

c. SENTENCE TOKENIZATION

It's essential to tokenize the sentences for input into the neural network. This process was achieved using the **Tokenizer from keras.preprocessing.text**. As illustrated in the following figure, the sentences were tokenized and sequentially padded. The primary objective of padding is to ensure that inputs are of fixed length. The padding functionality is provided by **pad_sequences from keras.preprocessing.sequence**.

```
Tokenize:
[[list([25, 26, 169, 27, 22, 20, 21, 19, 23, 28, 169, 24, 1292, 753])]
list([787, 623, 192, 16155, 87, 3858, 8134, 86255, 14913, 2, 4018, 338, 1265, 1, 56, 13842, 2302, 350, 5338, 9763, 197, 3240])
list([8134, 28350, 81, 41454, 3186, 1062, 86256, 54845, 62, 2, 580, 315, 1412, 1292, 7163, 15, 2, 939, 4, 2430, 16156, 3484, 1199, 1591, 3186, 1, 8887, 3017, 5, 6914, 13843, 481, 15, 739, 44, 2, 34, 2930])
...
list([244, 22179, 1351, 118, 72947, 875, 134, 129, 2285, 31182, 2289, 20552, 12190, 21783, 101, 46, 20087, 216, 53533, 139, 156, 111, 54, 301, 21982, 203304, 42729, 164])
list([58, 904, 7, 21783, 1815, 2, 24370, 54, 301, 10003, 1155, 54745, 95, 283305, 1426, 3, 17761, 32251, 63, 1266, 1469, 2643, 1, 200, 2, 1024, 560])
list([1013, 3, 26228, 3, 200, 904, 6138, 1032, 408, 2403, 58, 655, 904, 297, 341, 4353, 1, 174, 70879, 2, 4353, 12, 7890, 12041, 53644, 1760])]]
Padding:
[[ [ 25 26 169 ... 0 0 0]
 [ 787 623 192 ... 0 0 0]
 [ 8134 28350 81 ... 0 0 0]
 ...
 [ 244 22179 1351 ... 0 0 0]
 [ 58 904 7 ... 0 0 0]
 [ 1013 3 26228 ... 0 0 0]]
```

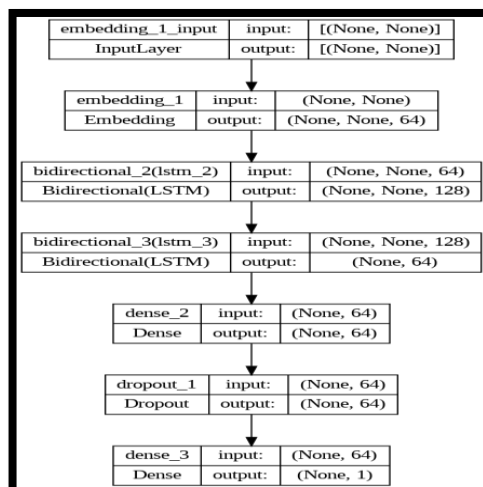
(Figure 4)

d. SHUFFLING DATASET

After the dataset preparation step, it was discovered that the data points belonging to one class were more numerous than the other class among the first 50,000 data points. To mitigate the impact of this situation, the dataset was shuffled by **shuffle from sklearn.utils**.

4. MODEL DEFINITION

To create a model, I conducted a literature review and experimented with various models. After testing, I selected the following model to use in the classifier.



(Figure 5)

Model: "sequential_1"		
Layer (type)	Output Shape	Param #

embedding_1 (Embedding)	(None, None, 64)	13011584
bidirectional_2 (Bidirectional)	(None, None, 128)	66048
bidirectional_3 (Bidirectional)	(None, 64)	41216
dense_2 (Dense)	(None, 64)	4160
dropout_1 (Dropout)	(None, 64)	0
dense_3 (Dense)	(None, 1)	65

Total params: 13123073 (50.06 MB)		
Trainable params: 13123073 (50.06 MB)		
Non-trainable params: 0 (0.00 Byte)		

(Figure 6)

Text data is inherently high-dimensional, sparse, and discrete. An embedding layer helps in transforming this data into dense, continuous vector representations where similar words or phrases are closer together in the embedding space. This transformation enables the model to learn more meaningful representations from the text data, capturing semantic and syntactic similarities.

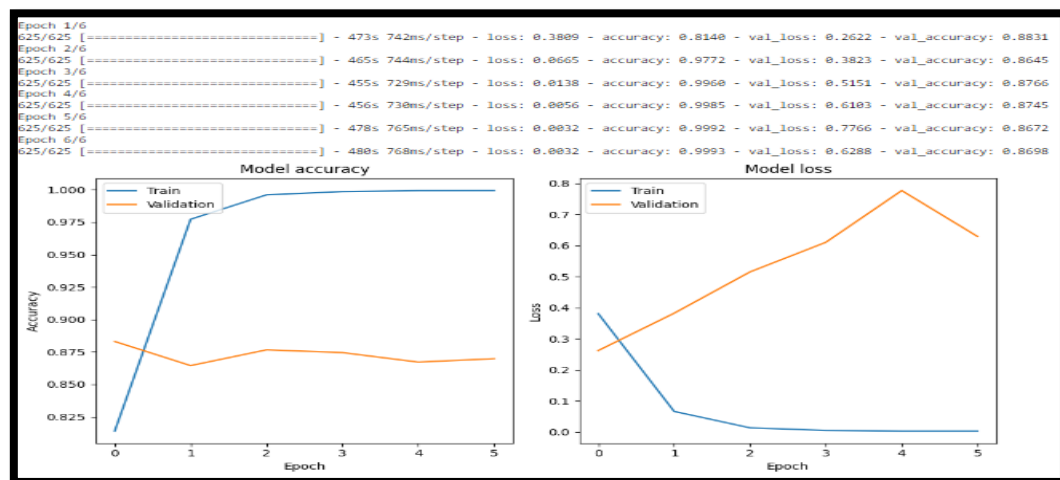
LSTMs are designed to handle sequential data, making them well-suited for processing text, which is inherently sequential in nature. They can capture dependencies and relationships between words across varying distances within a text, allowing them to understand context and meaning more effectively.

Two Dense layers are added to the model. The first Dense layer has 64 units and uses the ReLU (Rectified Linear Activation) activation function, which introduces non-linearity to the model. The second Dense layer has 1 unit and uses the sigmoid activation function. This layer serves as the output layer for binary classification tasks, where the output value (between 0 and 1) represents the probability of the input belonging to the positive class.

A Dropout layer is added after the first Dense layer. Dropout is a regularization technique used to prevent overfitting by randomly setting a fraction of input units to zero during training. In this case, 50% of the input units are randomly dropped out during each training iteration.

Finally, the model is compiled using the Adam optimizer, binary cross-entropy loss function (suitable for binary classification problems), and accuracy as the evaluation metric.

5. TRAINING



(Figure 7)

The model is trained over 6 epochs with batches of size 64. Half of the dataset is allocated for training, and 20% of the training set is used for validation.

During training, the loss decreases significantly from 0.3809 to 0.0032 over the epochs, indicating that the model is effectively learning and improving its predictions on the training data. The training accuracy also increases substantially from 0.8140 to 0.9993, demonstrating that the model's predictions are becoming increasingly accurate on the training data.

However, the validation loss fluctuates and generally increases from 0.2622 to 0.6288, suggesting that the model's performance on unseen data deteriorates over the epochs. This indicates potential overfitting as the model becomes overly specialized to the training data. Despite this, the validation accuracy fluctuates but remains relatively stable around 0.86 to 0.88, indicating that the model's performance on unseen data is consistent but not improving significantly.

Overall, while the model achieves very high accuracy on the training data, its performance on the validation set suggests overfitting, as indicated by the increasing validation loss and relatively stable validation accuracy. This suggests that the model may benefit from regularization techniques or adjustments to improve its generalization performance.

6. TESTING

The model is tested on both unseen 50,000 data samples and seen 50,000 data samples. The results of this test are shown in the following figure.

DATA TYPE	ACCURACY	LOSS
UNSEEN	0.8699	0.629187
SEEN	0.97326	0.128456

(Figure 8)

The system exhibits better performance on seen data samples, as expected. As a next step, additional regularization techniques such as Dropout, Early Stopping, etc., can be applied to the model to mitigate overfitting. These techniques help in improving the model's generalization ability, ensuring that it performs well on unseen data samples as well. These techniques may avoid overfitting and provide better results on unseen data samples.

Some unseen sentences along with their original and predicted labels are shown in the following figures.

Sentence	Ground Truth	Prediction	Predicted Label
Önce eskisi takibini maznet göstererek kadı abdurrehman pasha idaresindeki dostları tüfeklisi askerlerini rumeli'ye gönderilmiş dağılı gösterilmiş dağılı eskisi karşıında bazı beşarlar kazanan yeni ouaın askerlerinden bir kısmı İstanbul'a geri çağırılması corlu'ya ve lüleburgaz'a yerleştirilmiştir	not_separated	0.045349	separated
burada formula 4 x komundaki belirsizliği formula 5 ise x yönündeki momentundaki belirsizliği temsil eder	not_separated	0.999947	not_separated
sonunda batıdaki osmanlı savunma hatlarını kırarak rus ordularının ölü açılması dirençle karşılayınca İstanbul'un eşliğine yesilbölü kadar ilerleyerek osmanlı devleti'nin varlığını tehdit etmiş ve bunun sonucunda osmanlı devleti ayntefakus antlaşmasını imzalamak zorunda kalmıştır	not_separated	0.999972	not_separated
sarıçam bitki	not_separated	0.999987	not_separated
cinsel tercihle göre zaman zaman erkek müşterilerde giderlerdi	separated	0.825658	not_separated
ölüğünde dünyaya en yaşlı politikacısıydı	not_separated	0.999845	not_separated
günde 25'inin 50	not_separated	0.998573	not_separated
hem kendimiz hemde başkaları elbirliğiyle mutlak eğilimimizi ölümlerle aşırıktılara engel oluyor	separated	5.97924e-07	separated
gizli konusunda yöresel kıyafetler daha çok sandıklı yöresine ait karakteristik özellikler gösterirde ayrıntılarla eskiden gelen bazı motifleri yakalanak mümkün olabilmektedir	separated	1.83935e-05	separated
günöğünü 11 merkeze 72 km uzaklıkta yer alan şanta harabeleri merkez ilçe osmanlı köyü sınırları içerisinde bulunmaktadır	not_separated	0.999998	not_separated
söylen çözümlemesinde konuşmacının özelliklerine yönelik açıkça söylenmiş yazılmayan bir anlaşıma vardır	not_separated	0.897779	not_separated
polonya'daki futbol kulüpleri listesi	not_separated	0.999996	not_separated
sütl değişik biçimlerde ürettikleri değerdendirilir krema yapmak taze tereyağı kapattılması tereyağı yağart peynir	not_separated	0.8052108	separated
osmanlı yöneticilerinde barış sağlanmasından yanaydılar	separated	0.545184	not_separated
ancak yine de 1887 Depresine rağmen tarihi ve kültürel dokusunu büyük ölçüde muhafaza etmeyi başarmıştır	not_separated	0.999995	not_separated
çifti erkekle ve çift kadınlar karşılaştıkları tek türdür ve karışık çiftlerde kalifikasyon yarışması yapılmaz	not_separated	0.117886	separated
İbn hattuta kah lüks içinde kah güvencilik içinde bir yaşam sürüyordu	not_separated	0.998407	not_separated
2006 yılında canonical küresel destek ve hizmetler için montreal'de bir ofis açmıştır	separated	0.8001821	separated
karasol iklimin hakim olduğu eldağ dağıt bittinin güderek yok olmasıyla birlikte bozkır step görünümündedir	not_separated	0.999993	not_separated
tümöyle first kavuşu içerisinde kalan 11 doğal sınırlarla kuşatılmış yüksek bir bölgedir	not_separated	0.999998	not_separated

(Figure 9)

Sentence	Ground Truth	Prediction	Predicted Label
1964 1969 yılları arasında trt ankara radyosu tiyatro bölümünde görev aldı aralarında liyada don kisot goriot baba gibi basyapıtların da bulunduğu 18 yapıtı radyo için oynastırdı 180'ü aşkın radyo oyun dizisini yönetti	not_separated	0.999999	not_separated
çinli iso 9000 belgesi müşteriye kaliteyi anlatmak için alınmaktadır bu nedende seçilecek kuruluşun müşteri tarafından kabul görmesi gerekir	separated	4.48556e-06	separated
seferberliğin ilanından sonra kilise ve misyoner okullarının ermeni halkı büyük devletlerin övelliikle rusya'nın yardımı geleceğini ve onları bağımsızlığa kavuşturacaklarına dair yaptıkları telkinler sonucunda birçok olaylar yaşandı	separated	2.5924e-05	separated
burada 2 bölümlü olarak düşümlen tesisin 200 m ² 'lik pide salonu bölümü sol tarafta yer almaktadır	not_separated	0.000143077	separated
tarım sektörü ise ekonomide sadece 4 gibi bir oranı oluşturur	not_separated	0.998734	not_separated
promosyon ise toplulma fiyat indirimi olarak algılanmaktadır asıl amaç tüketiminin güvünde markanın ve ürünün bilinirliğini ve ürünün değeri satışı ve imajı gibi kriterlerin değerini arttırmaktır	not_separated	0.997578	not_separated
dünyadaki en eski tavla 68 parçasıyla beraber güneydoğu iran'da bulunmuştur	not_separated	0.999999	not_separated
siyah beyazlı ekibin orta sahasının değişmez ismi olan ve çok kritik gollerde imza atan kaya köstepen besiktas'ta 4 lig şampiyonluğu yaşadı	separated	0.00188985	separated
bu sözcük ile sunlardan biride kastedilmiş olabilir	separated	8.47518e-06	separated
toplantıya katılan bazı kisiler trump'ın söylemlerini avrupa için tehlikeli bulduklarını ve yıkıcı sonuçlara gödebileceğini belirtirken kimileride abartılı veya bir fırsat olarak gördüklerini açıkladı	separated	0.000276234	separated
burada nowiki formula 2 nowiki planck sabiti ayrıca 1899 yılında zaten tanıttıldı planck'ın kuantum kuramı olarak da bilinir	not_separated	0.999821	not_separated
aztek İmparatorluğu'nun İspanyollar tarafından fethi sırasında 1521'deki tenochtitlan kuşatması sonrasında şehrin neredeyse tamamı yakına uğradı	not_separated	0.943618	not_separated
krışna cocukken tereyağı çalan gençken de fidut çalıp yaramazlık yapan bir tarı olarak bilinse de yetmiş yıllarında daha ciddi tarafının ön plana çıktığı hikemelli bir filozof olarak tasvir edilmiştir	separated	0.0342758	separated
ülkede 1976'lı yıllara kadar bulunmayan demiryolu ağı günümüzde 2008 yılı verilerine göre 649 km'ye ulaşmıştır	not_separated	0.999999	not_separated
dört tanede yeni şarkı buluyordun sweetheart i still believe when you believe ve do you know where you're going to theme from mahogany	separated	0.000225441	separated
eman aynı zamanda bazı berberi grupların da kendilerini arap bir ecada dayamak için aldatıcı ve abartılı bir cabaya girdüklerini ayrıntılı şekilde belgeler	not_separated	0.999846	not_separated
giorgio de chirico	separated	9.43026e-05	separated
belirli bir piyasada belirli bir fiyat düzeyinde tüketicilerin alması hazır oldukları mal miktarının üreticilerin o fiyattan satmaya istekli oldukları miktardan daha fazla olması sonucu ortaya talep fazlası çıkar	not_separated	0.999997	not_separated
sırt yüğecide daha az dik ve üçgendir	separated	0.000125167	separated
oda sıcaklığının üstünde 3 kati halde bulunur	not_separated	0.999999	not_separated

(Figure 10)

7. APPENDIX

[GOOGLE COLAB LINK FOR CODE](#)