

Natural Language Processing  
CSE 454

Hw#4

Ömer Faruk Akduman  
1801042094

## Content

Homework Concept.....	3
Data Set .....	3
Data Preprocessing.....	3
Cleaning Data and remove unwanted characters from the sentences.....	3
At the end of the Data Preprocessing some samples .....	4
Encode and Tokenize.....	5
Splitting Data .....	5
Model .....	6
Results .....	7
With at most 40 characters sentences.....	7
With at most 30 characters sentences.....	7
Result on the Training Set .....	8
Result on the Test Set.....	8
Blue Score.....	9
Resources .....	10

## Homework Concept

In this homework, I will develop a simple translator from Ottoman Turkish to modern day Turkish. Provided training set at the class Teams page. As an example if the input is “itilaf devletleri , mütareke ahkâmına riayete lüzum görmüyorlar” the translated output will be “itilaf devletleri , ateşkes hükümlerine uymayı gerekli görmüyorlar”.

## Data Set

The name of the dataset file I used is jsonDataset.txt. This file contains information like this

```
{"dataset" : [{"translation":{"ot":"senesi mayısının 19 uncu günü samsuna çıktım  
.", "tr":"senesi mayisinin 19 uncu günü samsuna çıktım ."}},...
```

“ot” means in Ottoman language.

“tr” means in Turkish language.

## Data Preprocessing

Cleaning Data and remove unwanted characters from the sentences.

To make more standardize I remove Turkish characters and make string clean

```
#cleaning data and removes unwanted characters
def str_cleaning(a_str):
    a_str=a_str.lower()
    a_str = a_str.replace("ı","i").replace("ü","u").replace(".", "").replace("ş","s").replace("ç","c").
    result = ''.join(c for c in a_str if c.isalpha() or c == " ")
    return re.sub(' +', ' ', result).rstrip()
```

At the end of the Data Preprocessing some samples

	Turkish	Ottoman
38451	eller agza gitti	eller agza gitti
8813	tut ki uzattim	tut ki uzattim
4513	belge	vesika
18006	oraclakilere soze bakin clemis	huzzara demis gorun kelami
31111	onun dostu evangeliyayi sirkecide bir gazinoda...	onun dostu evangeliyayi sirkecide bir gazinoda...
37647	prenses senetlerin arif beye goturulup goturul...	prenses senetlerin arif beye goturulup goturul...
25265	ruhun bazi anlasilmaz ifade edilmez sikâyetler...	ruhun bazi anlasilmaz ifade edilmez sikâyetler...
21236	ayagindan ipi gevsetmeyi akil etmez o da	ayagindan ipi gevsetmeyi akletmez o da
11833	elbette boyle bir iliskiye girismek kesin olar...	bittabi boyle bir munasebete girismek katiyyen...
16706	salih pasa bu durumu biliyor ve buna bilerek n...	salih pasa bu vaziyeti bilerek ve bu hale bile...

## Encode and Tokenize

I tokenize sentences and use one hot encode to encode.

```
[ ] def create_tokenizer(lines):
    # fit a tokenizer
    tokenizer = Tokenizer()
    tokenizer.fit_on_texts(lines)
    return tokenizer

def max_len(lines):
    # max sentence length
    return max(len(line.split()) for line in lines)

def encode_sequences(tokenizer, length, lines):
    # encode and pad sequences
    X = tokenizer.texts_to_sequences(lines) # integer encode sequences
    X = pad_sequences(X, maxlen=length, padding='post') # pad sequences with 0 values
    return X

def encode_output(sequences, vocab_size):
    # one hot encode target sequence
    ylist = list()
    for sequence in sequences:
        encoded = to_categorical(sequence, num_classes=vocab_size)
        ylist.append(encoded)
    y = np.array(ylist)
    y = y.reshape(sequences.shape[0], sequences.shape[1], vocab_size)
    return y
```

## Splitting Data

I divided the dataset into 2 parts, train and test.

```
# Prepare training data
trainX = encode_sequences(src_tokenizer, src_length, train[:, idx_src])
trainY = encode_sequences(tar_tokenizer, tar_length, train[:, idx_tar])
trainY = encode_output(trainY, tar_vocab_size)

# Prepare test data
testX = encode_sequences(src_tokenizer, src_length, test[:, idx_src])
testY = encode_sequences(tar_tokenizer, tar_length, test[:, idx_tar])
testY = encode_output(testY, tar_vocab_size)
```

## Model

I used LSTM model also used some features like Early stopping, due to *model training time* “approximately it takes 8 hours for one model in colab” I only train 3 different model and put here best one with these features.

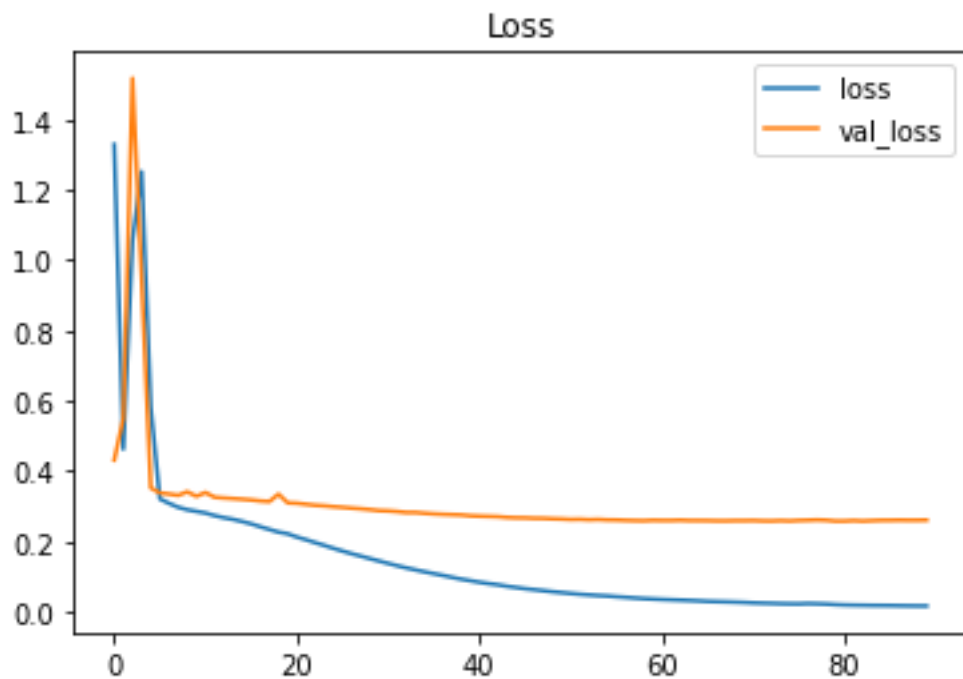
```
[ ] def create_model(src_vocab, tar_vocab, src_timesteps, tar_timesteps, n_units):
    # Create the model
    model = Sequential()
    model.add(Embedding(src_vocab, n_units, input_length=src_timesteps, mask_zero=True))
    model.add(LSTM(n_units))
    model.add(RepeatVector(tar_timesteps))
    model.add(LSTM(n_units, return_sequences=True))
    model.add(TimeDistributed(Dense(tar_vocab, activation='softmax'))))
    return model

# Create model
model = create_model(src_vocab_size, tar_vocab_size, src_length, tar_length, 256)
model.compile(optimizer='adam', loss='categorical_crossentropy')

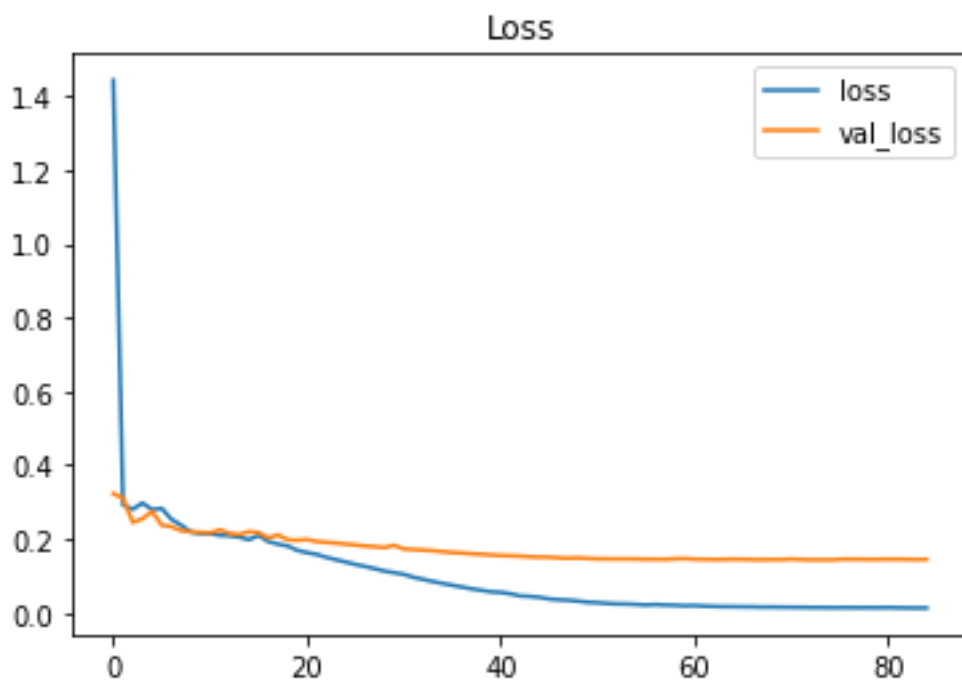
history = model.fit(trainX,
                    trainY,
                    epochs=200,
                    batch_size=64,
                    validation_split=0.1,
                    verbose=1,
                    callbacks=[
                        EarlyStopping(
                            monitor='val_loss',
                            patience=10,
                            restore_best_weights=True
                        )
                    ])
])
```

## Results

With at most 40 characters sentences



With at most 30 characters sentences



## Result on the Training Set

### Result on the Training Set ### OTTOMAN (SOURCE)	TURKISH (TARGET)	AUTOMATIC TRANSLATION IN TURKISH
birkac bez parcası cikardilar at basından belayı nerede birini sevmez mi gece evime gittim fayda cikmadi kalbi kanaya kanaya yandi sana para lazim yazihanede bir daha caldi taniyamadi ama suphelendi nicin gitmemeli ne nezahet bu hocam hayranim beye bir sey soylar diye gelsin a be cok oturacak misiniz ona guzel diyor geceleri galiba bekci geziyor bana sadece margarita deyin y rab o ne encumen ne alem hangi arkadasindaydin nihayet ben nazim beyi kabul etmedim daha neler yerden goge kadar hak etmistı gunler geciyor	birkac bez parcası cikardilar at basından belayı nerede birini sevmez mi gece evime gittim fayda cikmadi kalbi kanaya kanaya yandi sana para gerek yazihanede bir daha caldi taniyamadi ama suphelendi nicin gitmemeli ne incelik bu hocam hayranim beye bir sey soylar diye gelsin a be cok oturacak misiniz ona guzel diyor geceleri galiba bekci geziyor bana sadece margarita deyin ya rab o ne meclistir ne cm hangi arkadasindaydin nihayet ben nazim beyi kabul etmedim daha neler yerden goge kadar hak etmistı gunler geciyor	birkac bez parcası cikardilar at basından belayı nerede birini sevmez mi gece evime gittim fayda cikmadi kalbi kanaya kanaya yandi sana para gerek yazihanede bir daha caldi taniyamadi ama suphelendi nicin gitmemeli ne incelik bu hocam hayranim kapici bir kadar soylar bey bey mu cok oturacak misiniz ona bey ki geceleri galiba bekci geziyor bana sadece margarita deyin ya rab o ne meclistir ne cm hangi arkadasindaydin nihayet ben nazim beyi kabul etmedim daha neler yerden goge kadar hak etmistı gunler geciyor

## Result on the Test Set

### Result on the Test Set ### OTTOMAN (SOURCE)	TURKISH (TARGET)	AUTOMATIC TRANSLATION IN TURKISH
yalnizlikten canim sikiliyor ihtiyarladim hanimefendi alisverisin ne daha neler daha neler hukumeti de ayni halde siz nereden biliyorsunuz hangi siir bir degil haklisin teyze silahimiz var midir erzurumdan k o kalirsiniz oyle mi kesfettigim gibi bu sabah annesi geldi yazan biraderimdir evet acaba hata mi ettim aman sor bana merak oldu hafta icinde biter salincakci sirri sezai korkak korkak telasla haykirdilar tatmin edilmege muhtactir melahat burhanettini sevmiyor kaynatmis dervis ahmed	yalnizlikten canim sikiliyor ihtiyarladim hanimefendi alisverisin ne daha neler daha neler hukumeti de ayni durumda siz nereden biliyorsunuz hangi siir bir degil haklisin teyze silahimiz var midir erzurumdan k o kalirsiniz oyle mi kesfettigim gibi bu sabah annesi geldi yazan agabeyimdir evet acaba yanlis mi yaptim aman sor bana merak oldu hafta icinde biter salincakci sirri sezai korkak korkak telasla haykirdilar yatistirilmesi gerekir melahat burhanettini sevmiyor kaynatmis dervis ahmed	canim canim hanimefendi ne daha daha zavalli neler hukumeti de durumda durumda pasa miyiz biliyorsunuz hangi siir bir degil haklisin teyze olur olur olur erzurumdan k o oyle mi saika gibi bu bu ankaraya yaptim is evet acaba yanlis mi yaptim onu onu zavalli oldu oldu kadin icinde icinde salincakci sirri sezai gogus gecirdi saika ve ikinci tanidiklardan seviyor siz da seviyor bakistilar dervis ahmed



## Blue Score

### Blue Score function

```
def bleu_score(model, tokenizer, sources, raw_dataset):
    actual, predicted = [], []
    for i, source in enumerate(sources):
        source = source.reshape((1, source.shape[0]))
        translation = predict_seq(model, tar_tokenizer, source)
        raw_target, raw_src = raw_dataset[i]
        actual.append([raw_target.split()])
        predicted.append(translation.split())

    bleu_dic = {}
    bleu_dic['1-grams'] = corpus_bleu(actual, predicted, weights=(1.0, 0, 0, 0))
    bleu_dic['1-2-grams'] = corpus_bleu(actual, predicted, weights=(0.5, 0.5, 0, 0))
    bleu_dic['1-3-grams'] = corpus_bleu(actual, predicted, weights=(0.3, 0.3, 0.3, 0))
    bleu_dic['1-4-grams'] = corpus_bleu(actual, predicted, weights=(0.25, 0.25, 0.25, 0.25))

    return bleu_dic

bleu_train = bleu_score(model, tar_tokenizer, trainX, train)
bleu_test = bleu_score(model, tar_tokenizer, testX, test)
```

### Blue Score Results

```
[27] bleu_train

{'1-grams': 0.9063629236970987,
 '1-2-grams': 0.8653306887308325,
 '1-3-grams': 0.8023825817067682,
 '1-4-grams': 0.6401645192328522}

bleu_test

{'1-grams': 0.4469939190143961,
 '1-2-grams': 0.35591626789313996,
 '1-3-grams': 0.3279145847813834,
 '1-4-grams': 0.22312765673187918}
```

## Resources

- 1- GitHub
- 2- [www.spicework.com](http://www.spicework.com)
- 3- Stackoverflow
- 4- Wikipedia
- 5- [kaggle.com](http://kaggle.com)
- 6- chatgpt
- 7- <https://towardsdatascience.com/>