ÜSKÜDAR ÜNİVERSİTESİ FEN BİLİMLERİ ENSTİTÜSÜ BİLGİSAYAR MÜHENDİSLİĞİ YÜKSEK LİSANS PROGRAMI 2022-2023

FİNAL PROJE RAPORU

ÖĞRENCİ NUMARASI: 214327027

ÖĞRENCİ ADI-SOYADI: Hamza Talip Ekşi

PROJE BAŞLIĞI: IMDB Film Yorumlarının Duygu Analizi

1. **GİRİŞ**

Duygu analizi, mesajın duygusal tonunun olumlu, olumsuz veya tarafsız olup olmadığını belirlemek için dijital metnin analiz edilme sürecidir. Fikir madenciliği olarak da bilinen duygu analizi, şirketlerin ürün ve hizmetlerini geliştirmelerine yardımcı olan önemli bir iş zekâsı aracıdır.Bu çalışmada eğlence sektörü üzerinden duygu analizi yapacağız.

2. MATERYALLER VE YÖNTEMLER

Kaggle.com'dan elden edilen IMDB Movie Ratings veri seti ("https://www.kaggle.com/datasets/yasserh/imdb-movie-ratings-sentiment-analysis") kullanılarak bir duygu analizi yapılmıştır.Dataset kullanıcı yorumları ve olumlu-olumsuz(0-1) puandan oluşmaktadır.Ayrıca model oluşturulduktan sonra tahmin yapıması için IMDB sitesinden veri setinde bulunmayan yorumlar alınmıştır.

3. ARAÇLAR

Analizin yapılması için gerekli kod Google Colaboratory üzerinde pyhton dili ile yazılmıştır.

4. KODLAR

```
[ ] !pip install --user tensorflow !pip install --user keras
```

-tensorflow ve keras kütüphaneleri yüklendi.

```
[ ] from tensorflow.python.keras.models import Sequential from tensorflow.python.keras.layers import Dense, GRU, Embedding, CuDNNGRU from tensorflow.python.keras.optimizer_v2 import adam as adam_v2 from keras.preprocessing.text import Tokenizer from keras.utils import pad_sequences
```

- -Yapay sinir ağı oluşturacağımız için Sequential import ediyoruz.
- -Katmanlarını oluşturmak için Dense,bilgi akışını düzenlemek için GRU,kelimeleri vektörleyebilmek için Embedding,GPU üzerinden çalışma için CuDNNGRU import edildi.
- -Optimizer olarak adam_v2 import edildi.
- -Tokenleme işlemi yapılacağından Tokenizer import edildi.

-Kelime vektörlerinin belli bir uzunlukta olmasını istediğimiz için pas_sequences import edildi.

```
[ ] import numpy as np
import pandas as pd
```

-Dizilerle çalışmak için numpy, veri analizi için pandas kullanıldı.

```
from google.colab import drive
drive.mount('/content/gdrive')
```

-Dataseti çekebilmek için drive ,colab'a bağlandı.

```
[ ] %cd /content/gdrive/MyDrive/Colab Notebooks/
/content/gdrive/MyDrive/Colab Notebooks
```

-Datasetin çekleceği klasör girildi.

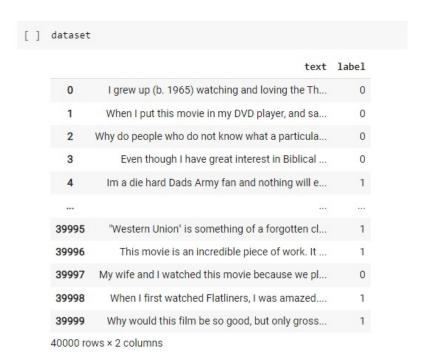
```
!ls

'A TALE OF TWO CITIES.txt'
movie.csv
NlpFinal.ipynb
NlpVize
'Topic models with Gensim.ipynb adlı not defterinin kopyası'
word2vec.ipynb
word2vec.model
```

-Klasördeki dosyalar listelendi.

```
[ ] dataset = pd.read_csv('movie.csv')
```

-Listeden movie.csv dosyası veri seti olarak kullanıldı.



-Dataset kolonları text ve label olacak şekilde 2 tanedir.

```
[ ] target = dataset['label'].values.tolist()
    data = dataset['text'].values.tolist()
```

-Datasette bulunan değerler listeye çevrildi.

```
[ ] cutoff = int(len(data) * 0.80)
    x_train, x_test = data[:cutoff], data[cutoff:]
    y_train, y_test = target[:cutoff], target[cutoff:]
```

-Sahip olan datanın yüzde 80'i eğitim için kullanılacak.Bu yüzden bir cutoff noktası oluşturuldu.

```
[ ] x_train[500]
```

'Just because an event really happened doesn't mean that it will make a good screenplay/ movie. The Cat's Meow, by Peter Bogdanovich claims to be based on a writer paid more attention to creating a bizarre cast of characters than taking time to create a story for the bizarre characters to inhabit. The key moment the producer, believing him to be Chaplin. Basing a key element of a story on someone wearing the wrong hat is trite and contrived. The story attempts to be n, comedy. There is also a lack of empathy for any of the characters. It hardly matters who is shot, who is killed, who is guilty and who is innocent. There ult to care about and the eventual outcome is incidental.'

- x_train 500.indekste bulunan "text" verisi.

```
[ ] x_train[800]
```

'Contrary to most of the comments in this section, I have to say this film just barely escapes the definition 'rubbish!' north-east and know at first hand what the area is like. I am totally sick of films that are supposed to be 'realistic', social background. And funny? I smiled briefly twice and laughed once, but that was at the incongruity of two boy actors as only one likable character in the whole film, and that was the senile grandfather played quietly but very competently ogy department the author studied at. He certainly had no ear for th...'

- x_train 800.indekste bulunan "text" verisi.

```
[ ] y_train[800]
0
```

-y_train 800.indekste bulunan "label" verisi. Verinin "0" olması olumsuz yorum olduğu anlamına gelir. Karşısındaki x_train[800] verisine bakarsak gerçekten olumsuz olduğunu görüyoruz.

```
num_words = 10000
tokenizer = Tokenizer(num_words=num_words)
```

-Yorumlarda en sık geçen ilk 10000 kelime alınıp tokenleştirildi.

```
[ ] tokenizer.fit_on_texts(data)
```

-İlk 10000 data ile devam edileceği belirtildi.

```
tokenizer.word index
'and': 2,
     'a': 3,
     'of': 4,
     'to': 5,
     'is': 6,
     'in': 7,
     'it': 8,
     'i': 9,
     'this': 10,
     'that': 11,
     'was': 12,
     'as': 13,
     'for': 14,
     'movie': 15,
     'with': 16,
     'but': 17,
```

-Bu 10000 kelime görüntülendi.

```
[ ] x_train_tokens = tokenizer.texts_to_sequences(x_train)
```

-Tokenizer tarafından dikkate alınan 10000 kelime bir integer değerine çevrildi.

```
[ ] x_train[800]
```

'Contrary to most of the comments in this section, I have to say this film just barely escapes the definition 'rubbish!'. The only readers who seem to north-east and know at first hand what the area is like. I am totally sick of films that are supposed to be 'realistic', yet portray the working class social background. And funny? I smiled briefly twice and laughed once, but that was at the incongruity of two boy actors with Sunderland accents support as only one likable character in the whole film, and that was the senile grandfather played quietly but very competently by Roy Hudd -- the only non-copy department the author studied at. He certainly had no ear for th...'

```
print(x_train_tokens[800])
[3735, 5, 87, 4, 1, 801, 7, 10, 2460, 9, 24, 5, 130, 10, 18, 38, 1174, 3034,
```

-800.indeksteki textte bulunan kelimelerin token değerleri yazıldı.Nadir kelimeler tokenleştirmede yok sayılır.

```
x_test_tokens = tokenizer.texts_to_sequences(x_test)
```

-Test datasıda aynı işlemden geçirildi. Kelimeler integer değerine çevrilmiş oldu.

```
num_tokens = [len(tokens) for tokens in x_train_tokens + x_test_tokens]
num_tokens = np.array(num_tokens)
```

-Her "text" verisinde kaç token bulunduğu hesaplandı ve bir dizi haline getirildi.

```
np.mean(num_tokens)
217.364725
```

-Filtrelenmiş verinin ortalama token sayısı

```
np.max(num_tokens)
2122
```

-En fazla token sahibi yorum.

```
np.argmax(num_tokens)
34167
```

-En uzun veriye sahip indeksi bulur.

```
x_test[2167]

There\'s a sign on The Lost Highway that says:*MAJOR SPOILERS AHEAD*(but you already knew that, didn\'t you?)Since there\'s a great deal of people pretation of why the plot makes perfect sense. As others have pointed out, one single viewing of this movie is not sufficient. If you have the DVD of tonly upon second or third viewing, please.);)First of all, Mulholland Drive is downright brilliant. A masterpiece. This is the kind of movie that Vogue\'s "It gets inside your head and stays there" really hit the mark.David Lynch deserves praise for creating a movie that not only has a beautification and a very dream-like quality t...."
```

-40000 verinin yüzde sekseni 32000 veri yaptığından 34167 indeksindeki veri cutoff noktasından itibaren test verisidir. En uzun veriye sahip indeks bulundu.

```
max_tokens = np.mean(num_tokens) + 2 * np.std(num_tokens)
max_tokens = int(max_tokens)
max_tokens
```

-Ortalama ve iki standart sapmasını toplayarak yaptığımız hesapla elde ettiğimiz 533 değeri, en fazla 533 token uzunluğuna sahip verileri dikkate alarak verimizin yüzde 95'ini dikkate almış olacağımız anlamına geliyor.

```
[ ] np.sum(num_tokens < max_tokens) / len(num_tokens)
0.945225
```

-533'ten küçük tokenları toplayıp toplam token sayısına bölerek üstteki işlemin sağlaması yapıldı.

```
[ ] x_train_pad = pad_sequences(x_train_tokens, maxlen=max_tokens)
[ ] x_test_pad = pad_sequences(x_test_tokens, maxlen=max_tokens)
```

-533'ten küçük verileri doldurmamız, büyük verileri ise kesmemiz gerekiyor.pad_sequences ile 533'ten küçük verilerde boşluklar yerine "0" yazılırken, büyük verilerde ilk 533 dikkate alınıp geri kalan kısım kesilir.

```
[ ] x_train_pad.shape [ ] x_test_pad.shape (8000, 533)
```

-pad.shape'e baktığımızda train için 32000-533 değeri test için 8000-533 değeri görüldü.

```
[ ] np.array(x_train_tokens[800])
      array([3735,
                           87,
                                       1, 801,
                                                         10, 2460,
                                       38, 1174, 3034,
                                                         1, 4945,
                           10,
                                18,
               60, 6070,
                           34,
                                303,
                                       5,
                                             25, 1854,
                                                          4,
                                                               47,
            1692,
                           6,
                                22,
                                      142,
                                                         67,
                                                              411,
                     8,
                                             34.
                                                   36,
            2187, 2476,
                           2,
                               118,
                                      29,
                                             86,
                                                  523,
                                                         47,
                                                                1, 1635,
               36,
                         236,
                                460, 1157,
                                             4,
                                                  103,
                                                         11,
                                                               22,
                                                                    440,
              25,
                   243, 2056,
                                     765,
                                                                    359, 1831,
                                 1,
                                            716,
                                                               13,
                                                   4, 7187,
             351,
                    17,
                           90, 1378,
                                       89,
                                            83,
                                                    4,
                                                         64,
                                                             1128,
                                                                    955,
                     9, 3031, 1472,
                                       2, 1437,
                                                  280,
                                                         17,
                                                               11,
                                                                     12,
               1,
                      4,
                         104, 419,
                                     149,
                                            16, 2530, 1539,
                                                              263,
                                                                      5, 2485,
              64, 2530,
                                             45,
                                                         60,
                                                               26, 1496,
                                                                          107
                          35, 2030,
                                      454,
                                                   12,
                     1,
                          223,
                                 18,
                                       2,
                                             11,
                                                   12,
                                                         1,
                                                             3573,
                                                                    251, 5312,
                                                        687,
                     51, 9559,
                                 30, 2681,
                                             1,
                                                   60,
                                485,
                                       69,
                                             9,
                                                   38,
                                                        581,
                                                               47, 3675, 2447,
               13.
                    14.
                           1.
               1, 2212, 6879,
                                                   65,
                                       27, 418,
                                                         53, 4591,
                                 29.
                                                                     14.
              693, 1214,
                           2,
                                 9,
                                       95,
                                            24,
                                                  110,
                                                        545,
                                                                1, 2688,
               21,
                    60,
                        236,
                                  9,
                                       3, 2098,
                                                  144,
                                                         3,
                                                              883, 406,
                                                                           19,
                          88, 446, 124,
                                                  2,
                                                         19,
               1, 693,
                                            54.
                                                              10, 1169,
                                                                            8.
             210, 153,
                          8, 210, 789,
                                              8,
                                                  210,
                                                        425,
                                                                2,
[ ] x_train_pad[800]
      array([
                        0.
                                            0.
                                            0,
                        0.
                                            0,
                                                         0,
                        0,
                        0,
                                                         0,
                        0.
                                            0.
                                     0.
                                                         0.
                                            0,
                        0,
                                     0,
                 0.
                        0.
                              0.
                                     0.
                                            0.
                                                         0.
                                                                0.
                                                                       0.
                        0.
                                            0.
                                                         0.
                                                                0,
                                                                                  801,
                                                                       4,
                 0,
                              0,
                                     0,
                                            0, 3735,
                                                               87,
                       10, 2460,
                                           24,
                                                                     18,
                                     9,
                                                   5.
                                                       130.
                                                               10,
                                                                            38, 1174,
                           4945,
                                   864,
                                                      6070,
                                                                    303,
                                                                                   25,
                       1,
                                                  60,
                                                                              5,
              1854
                             47,
                                               1692,
                                                                      22,
                                                                           142,
                      67,
                            411,
                                               2187,
                                                      2476,
                                                                     118,
                36.
                                                                            29,
                                                                                   86,
               523,
                      47,
                                                              236,
                                                                          1157,
                                 1635,
                                                 36,
                              1,
                                            6,
                                                         9,
                                                                     460,
                                                  25,
                                                       243, 2056,
                    7187,
                             13,
                                   359, 1831,
                                                               90,
                                                                   1378,
                           1128,
                                   955,
                                                         9,
                                                            3031,
                                                                             2,
                                                                                1437,
                       64,
                 4.
                                           2,
                                                153,
                                                                    1472
                                           29,
                       17,
                                                                           149,
               280,
                             11,
                                   12,
                                                  1,
                                                              104,
                                                                     419,
                                                                                   16,
                    1539,
                            263,
                                     5, 2485,
                                                 64, 2530,
                                                               35,
                                                                           454,
                                  1496,
                                         107,
                             26,
                        1, 3573,
                                                 17,
                12.
                                  251, 5312,
                                                        51, 9559,
                                                                      30, 2681,
                     687,
                                          176,
                                                        14,
                                                                     485,
                                                                            69,
                60,
                                                 13,
                             47, 3675, 2447,
                                                  1, 2212,
                                                                            27,
                                   14,
                          4591,
                                                      1214,
                                                                      9,
                                                                            95,
                65.
                      53,
                                          1,
                                                693,
                                           2,
                                                                       9,
                                                              236,
                                                                             3, 2098,
               110.
                     545,
                              1, 2688,
                                                 21,
                                                        60,
                            883,
                                                       693,
                                                                     446,
                                                                           124,
               144,
                       3,
                                   406,
                                          19,
                                                               88,
                                                   1,
                                                210,
                      19,
                                                                     210,
                     425,
               210,
                                            6], dtype=int32)
```

-pad işleminden öncesinin ve sonrasının kıyaslaması yapıldı.Boş kısımların "0" ile doldurularak uzunluğunun 533'e tamamlandığı görülüyor.

```
[ ] idx = tokenizer.word_index
  inverse_map = dict(zip(idx.values(), idx.keys()))
```

-Sayıları tekrardan kelimelere çevirebilmek istiyoruz.Bunun için bir sözlük yapısı oluşturuldu.

```
[ ] def tokens_to_string(tokens):
    words = [inverse_map[token] for token in tokens if token!=0]
    text0 = ' '.join(words)
    return text0
```

-Sayılardan oluşan bir listeyi alıp bu sayıları kelimeye dönüştüren bir fonksiyon tanımlandı.

```
[ ] x_train[800]

'Contrary to most of the comments in this section, I have to say this film just barely escapes the definition 'ru north-east and know at first hand what the area is like. I am totally sick of films that are supposed to be 'real social background. And funny? I smiled briefly twice and laughed once, but that was at the incongruity of two boy as only one likable character in the whole film, and that was the senile grandfather played quietly but very comp ogy department the author studied at. He certainly had no ear for th...'

[ ] tokens_to_string(x_train_tokens[800])

'contrary to most of the comments in this section i have to say this film just barely escapes the definition ' the only readers first hand what the area is like i am totally sick of films that are supposed to be yet portray the working class of wherever an arms.
```

d laughed once but that was at the of two boy actors with accents supposedly trying to hide their accents from football fans th competently by roy the only non in the cast as for the writing well i just wonder what university department the author studied

-Ham ve filtrelenmiş halleriyle 800.indeks verisi.

i a native i'm a who's written on the local don't waste your time an..

```
[ ] model = Sequential()
[ ] embedding_size = 50
```

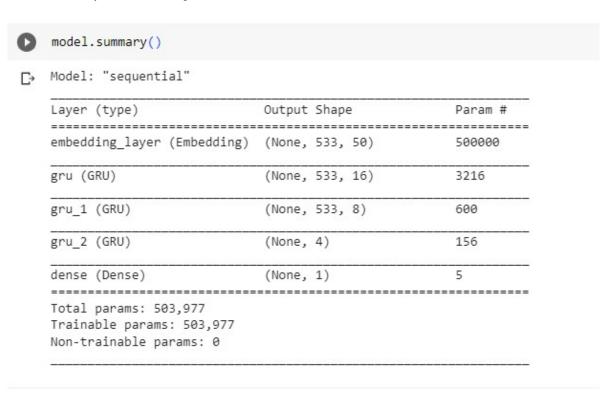
-Model oluşturuldu ve embedding size olarak "50" verildi.Rastgele kelime vektörleri başlatıp model içinde eğiteceğiz.

-10000x50 kelime modele eklendi. Verideki kelime vektörlerini sonraki katmana input olarak verecek.

```
[ ] model.add(GRU(units=16, return_sequences=True))
    model.add(GRU(units=8, return_sequences=True))
    model.add(GRU(units=4))
    model.add(Dense(1, activation='sigmoid'))
```

-GRU yapısı oluşur.Katmanlar arası çıktı alışverişi sağlandı.

-Modeli compile etmek için bazı özellikler verildi.



-Modelin özeti.

-Filtrelenmiş veri bir diziye dönüştürülüp modele yerleştirildi ve eğitime başlandı.

-Modele test edildi ve %88lik bir doğruluk görüldü.

```
[ ] y_pred = model.predict(x=x_test_pad[0:1000])
    y_pred = y_pred.T[0]
```

-İlk 1000 satır için tahminde bulunuldu.

```
[ ] cls_pred = np.array([1.0 if p>0.5 else 0.0 for p in y_pred])
```

-Gelen sonuç 0,5 üzerindeyse 1 yap,altındaysa 0 yap diyerek eşikten geçirildi ve tahmin sınıfı yapısı oluşturuldu.

```
[ ] cls_true = np.array(y_test[0:1000])
```

-Doğruluğu tespit etmek için gerçek sınıflar elde edildi.

```
[ ] incorrect = np.where(cls_pred != cls_true)
  incorrect = incorrect[0]
```

-İki yapı karşılaştırılarak eşleşmeyenlerin indeksleri "incorrect"'e atıldı.

```
[ ] len(incorrect)

127
```

-İlk 1000 satır için kaç yanlış tahmin yapıldığı bulundu.

```
idx = incorrect[0]
idx

0
```

-İlk yanlış tahminin yapıldığı indeks.

```
[ ] text0 = x_test[idx]
    text0

'This movie was a major bait and switch. I rented it because of Rebecca St. James, a popular Christian singer. I have met her
elp out a friend, or a friend of a friend. My first clue that this movie wasn't what it was supposed to be was when I witness
ors were. It was funny how almost everyone in the movie wore solid colors. (There are a few exceptions).Rebecca was verrryyy
se husband has return to the fold. Doesn't she ever leave the house? I had to turn off the movie several times in order to fi
r and go to hell. Jesus spread a message of love and hope. His messag...'
```

-İlk yanlış tahmin "0" indeksinde yapılmış.Bu indeks görüntülendi.Yorumun olumsuz olduğu açık.

```
[ ] y_pred[idx]
0.9440327
```

-Bizim modelimizin bu olumsuz yorumu olumlu olarak yanlış yorumladığı görülüyor.

```
[ ] cls_true[idx]
```

-Gerçek sınıfının "0" yani olumsuz olduğu görülüyor.

```
text1 = "A thrilling masterpiece worth the long run time even better in IMAX! Important to note Rip Lance Reddick he was amazing in his career!"

text2 = "Hands down the best action movie franchise ever."

text3 = "John Wick: Chapter 4 offered me 169 minutes of pure entertainment with an exciting cast of characters and action sequences that'll impress anyone, a perfect action film"

text4 = "Bad acting, bad casting, bad adaptation. Do not watch this if you enjoyed the book. From start to finish it's bad, there's no redeeming qualities to this film! Parker looks text5 = "This movie was nothing like the book, except for the character's names. The acting was atrocious and there was barely any chemistry...I came for the steam, but there is no vext6 = "I'm a huge fan of the book. I was so excited when the movie was announced. Man, did the screenwriter even read the book? This is the biggest jumbled mess of crap! The dialog text7 = "It all leads back to where we once started off as all great trilogies have indicated from the past. But GOTG surpasses expectations with what is nothing short of phenomenal text8 = "This. This is what I've wanted. Yeah some of the jokes are a bit too silly and the tone is a bit confusing at times. But the end result for me is probably the most heartfelt text9 = "I think for what this movie sets out to do, Super Mario Bros. Movie pretty much ticks the boxes of an entertaining, reference fan-service heavy movie. A lot of the jokes are text10 = "The story was simple, but worked wonders with the excellent characters that perfectly reflected their game counterparts. They even had unique aspects given to them that he texts = [text1, text2, text3, text4, text5, text6, text7, text8, text9, text10]
```

-IMDB sitesinden 2023 yılına ait 3 filme yapılan yorumları alarak modelimizde olumlu mu olumsuz mu oldukları test edildi.Bu yorumlar bir texts listesine eklendi.

```
[ ] tokens = tokenizer.texts_to_sequences(texts)
```

-Bu texts listesi tokenize edildi.

-Bu tokenlar daha sonra sınırımız olan 533'e pad edildi. Yani boşluuklar O'lar ile doldurulurken fazlalıklar silindi. 10 yorumumuz var hepsinin 533 token'ı var.

-Modelimizin yaptığı tahminler ortaya çıkmış oldu.

SONUÇLAR

-Tahmin yapılması istenen on yorumun ilk üçünün çok olumlu ,sonraki üç yorumun olumsuz, yedinci yorumun olumlu, sekizinci-dokuzuncu-onuncu yorumların ise nötre yakın olumlu yorumlar olduğu görülüyor.Modelimiz birinci yorumda doğru tahminde bulunsa da sonraki iki yorumu nötre yakın görmüş.Sonraki üç yorumda ise bu yorumların olumsuz olduğunu tahmin edebilmiş.Yedinci yorumu da doğru bir şekilde olumlu olarak tahmin etmiş.Sekiz ve onuncu yorumları çok olumlu bulmuş.Yorumların olumluya kaydığı doğru ama olumsuz ifadelerde var.Bu sebeple oran biraz daha düşük olabilirdi.Dokuzuncu yorumu olumsuza yakın görmüş ama nötr olması gerekirdi.

Sonuç olarak olumsuz yorumlarda gerçekten başarı yüzde yüz iken olumlu yorumların bir kısmında ve nötr yorumlarda tam olarak istenen başarıya sahip değil.

5. ÖNERİ VE TARTIŞMALAR

Modeli daha başarılı hale getirmemiz için eğitim kalitesi artmalı.Bunu sağlamak için ya veri miktarı artırılabilir ya da eğitim yapılırken kullanılan epoch sayısı artırılabilir.

KAYNAKLAR

- -https://www.kaggle.com/datasets/yasserh/imdb-movie-ratings-sentiment-analysis
- -https://www.imdb.com/title/tt10366206/?ref =adv li tt
- -https://m.imdb.com/title/tt2316548/?ref_=m_tt_urv
- -https://www.imdb.com/title/tt6718170/?ref_=adv_li_tt

KONTROL LISTESI

	EVET /HAYIR
Raporunuzu şablonda belirtildiği gibi hazırladınız mı?	Evet
Çalışmanızın sonuçlarını (print screen) rapora eklediniz mi?	Evet
Rapor dosyanızı şablondaki gibi yeniden adlandırdınız mı?	Evet
Raporu sisteme yüklediniz mi?	Evet
Kodları sisteme yüklediniz mi?	Evet