

```
[138]: import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
import torch
import torch.nn as nn
import torch.nn.functional as F
import multiprocessing as mp
from sklearn.metrics import roc_curve, auc, roc_auc_score, auc
from tqdm import tqdm, tqdm_notebook as tqdm
from sklearn.preprocessing import StandardScaler
from sklearn.feature_extraction.text import TfidfVectorizer, CountVecorizer
from scipy.sparse import hstack, vstack
import warnings
warnings.filterwarnings('ignore')
```

```
x_train,x_test,y_train,y_test=train_test_split(x,
```

```

In [104]: essay_voc.train('essay')
          fill essay tests (-), joint essays)
          wordtrain_essay_text.split (" ")

          imports pickle
          with open ('C:/glove_vectors', 'wb') as f:
              model = pickle.dump(glove_words = set(model.keys))

          from collections import Counter
          count=Counter(words)
          vocab=sorted(counts,key=count.get,reverse=True)
          vocab=glove_words - count.intersection(vocab)
          int_vocab = []
          for i, word in enumerate(vocab,1):
              int_vocab.append(i)
          train_essay_ints = []
          for essay in essays:
              train_essay_ints.append([vocab.to_int[word] for word in essay.split() if word in vocab])

          essay_voc.test('essay')
          test_essay_ints = []
          for essay in essays:
              test_essay_ints.append([vocab.to_int[word] for word in essay.split() if word in vocab])

In [105]: for fixed_length_input
          def pad_features(essay_ints, seq_lengths):
              features = np.zeros((len(essay_ints), seq_lengths), dtype=int)

              for i, row in enumerate(essay_ints):
                  if len(row)<seq_lengths:
                      continue
                  features[i, :len(row)] = np.array(row)[seq_lengths:]

              return features

In [107]: seq_length = 50

          train_features = pad_features(train_essay_ints, seq_lengths=seq_length)
          test_features=pad_features(test_essay_ints, seq_lengths=seq_length)

          ##sent int(train_features)-len(train_essay_ints)
          ##sent int(train_features[0])-seq_lengths
          print(train_features[1,1:10])

          [[ 285 2482 3269 17056 14253 7456 998 9587 1802 11903]
           17354 2441 26394 2254 12824 23264 17253 31885 31640 12924]
           23720 119 25451 26396 1182 8976 16845 28812 21758 25366]
           2255 2454 25336 1975 14225 24394 994 22697 18312 5925]
           16688 3213 32818 276 18496 18924 9612 28812 16039 4546]]

In [108]: weight_matrix=np.random.normal(size=(len(vocab)+1,500))
          for i in range(1,len(vocab)+1):
              weight_matrix[i]=model.get_vocab[1:]

In [109]: values=w_train['school_state']
          count=Counter(values)
          vocab=sorted(counts,key=count.get,reverse=True)
          vocab.to_int = words = []
          for i, word in enumerate(vocab,1):
              train_school_state_ints=vocab.to_int[word]
          values=w_train['school_state']
          test_school_state_ints=[vocab.to_int[word] if word in vocab else 0 for word in values]
          int_wg_vocab=len(vocab)+1

In [110]: values=w_train['project_grade_category']
          count=Counter(values)
          vocab=sorted(counts,key=count.get,reverse=True)
          vocab.to_int = words = []
          for i, word in enumerate(vocab,1):
              train_grades_ints=vocab.to_int[word]
          values=w_train['project_grade_category']
          test_grades_ints=[vocab.to_int[word] if word in vocab else 0 for word in values]
          int_pg_vocab=len(vocab)+1

In [111]: values=w_train['project_subject_category']
          count=Counter(values)
          vocab=sorted(counts,key=count.get,reverse=True)
          vocab.to_int = words = []
          for i, word in enumerate(vocab,1):
              train_categories_ints=vocab.to_int[word]
          values=w_train['project_subject_category']
          test_categories_ints=[vocab.to_int[word] if word in vocab else 0 for word in values]
          int_sgc_vocab=len(vocab)+1

In [112]: values=w_train['project_subject_subcategories']
          count=Counter(values)
          vocab=sorted(counts,key=count.get,reverse=True)
          vocab.to_int = words = []
          for i, word in enumerate(vocab,1):
              train_sub_categories_ints=vocab.to_int[word]
          values=w_train['project_subject_subcategories']
          test_sub_categories_ints=[vocab.to_int[word] if word in vocab else 0 for word in values]
          int_ssc_vocab=len(vocab)+1

In [113]: values=w_train['teacher_prefix']
          count=Counter(values)
          vocab=sorted(counts,key=count.get,reverse=True)
          vocab.to_int = words = []
          for i, word in enumerate(vocab,1):
              train_teacher_ints=vocab.to_int[word]
          values=w_train['teacher_prefix']
          test_teacher_ints=[vocab.to_int[word] if word in vocab else 0 for word in values]
          int_tv_vocab=len(vocab)+1

In [114]: def check_digits_weight():
          for word in w_val.split():
              if word.isdigit():
                  return 1
              return 0

          standardize=StandardScaler()#standardizing data
          train_prev_projects_array=train['teacher_number_of_previously_posted_projects'].reshape(-1,1)
          train_prev_projects_array=standardize.fit_transform(train_prev_projects_array)
          test_prev_projects_array=test['teacher_number_of_previously_posted_projects'].reshape(-1,1)
          test_prev_projects_array=standardize.transform(test_prev_projects_array)

          train_price_array=train['price'].reshape(-1,1)
          train_price=standardize.fit_transform(train_price_array)
          test_price_array=test['price'].reshape(-1,1)
          test_price=standardize.transform(test_price_array)

          train_quantity_array=train['quantity'].reshape(-1,1)
          train_quantity=standardize.fit_transform(train_quantity_array)
          test_quantity_array=test['quantity'].reshape(-1,1)
          test_quantity=standardize.transform(test_quantity_array)

          train_summary_contains_digit_array=train['project_recommender_name'].isin(check_digits).reshape(-1,1)
          train_summary_contains_digit_array=train['project_recommender_name'].isin(check_digits).reshape(-1,1)
          train_remain_features=np.concatenate((train_prev_projects_train_price_train_quantity_train_summary_contains_digits),axis=1)
          train_remain_features=np.concatenate((test_prev_projects_test_price_test_quantity_test_summary_contains_digits),axis=1)

2.LSTM Model Building (model-1)

In [115]: def create_emb_layer(embedding_matrix,num_embeddings,Padding):
          do not embedding layer with our pre trained weights(in this case word vectors)
          num_embeddings=embedding_dim
          emb_layer=nn.Embedding(num_embeddings, num_embeddings, padding=Padding)
          if num_embeddings!=embedding_dim:
              raise ValueError('num_embeddings must be equal to embedding_dim')
          emb_layer.weight.requires_grad = False

          return emb_layer, num_embeddings, embedding_dim

          class Net(nn.Module):

```

```

super().__init__()
self.embedding = nn.Embeddings(embedding_size, embedding_size)
self.hidden_size = hidden_size

```

```
batch+=1
hidden=net.init_hidden()
a=a.float()
return net(a1,a2,a3,a4)
```