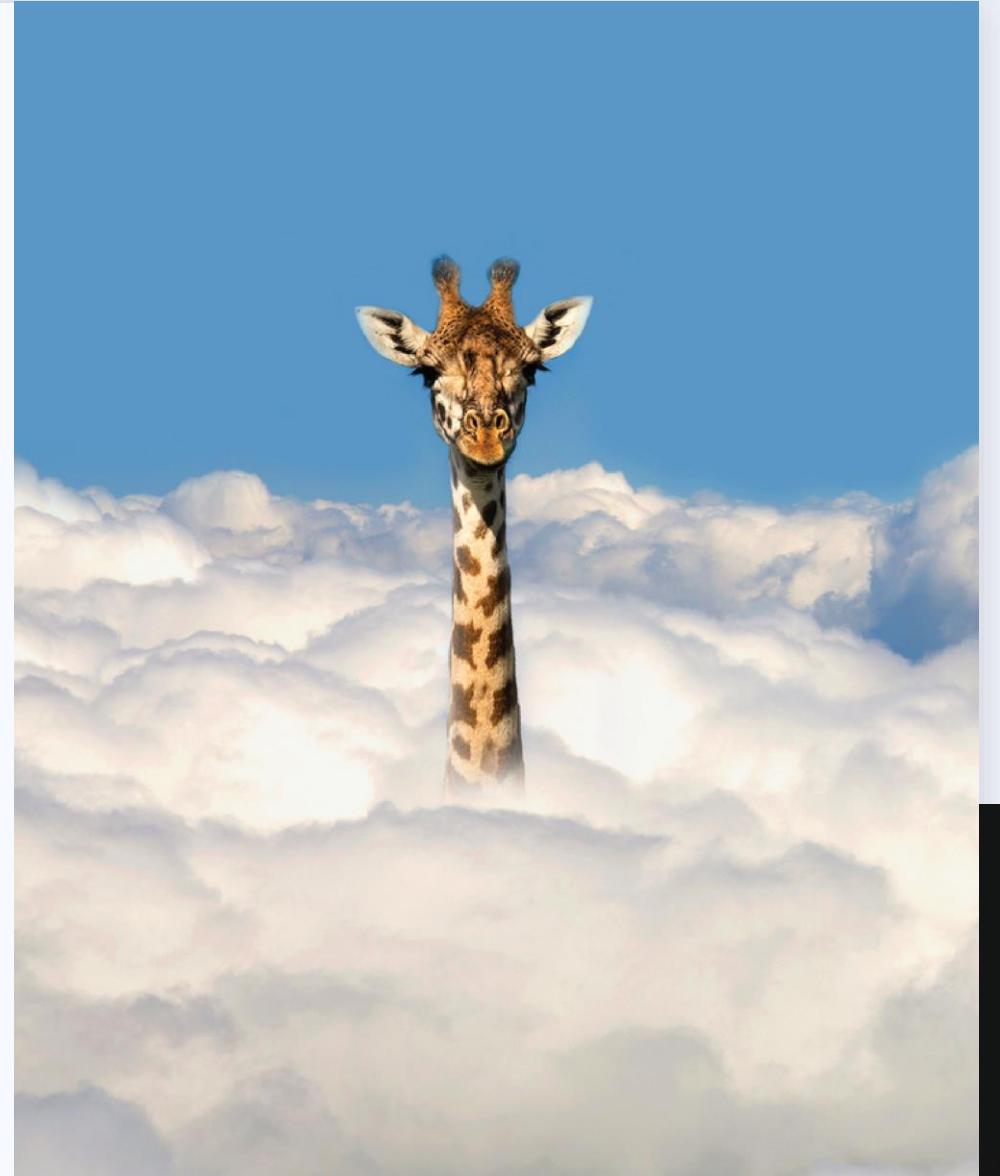


MULTILINGUAL NER

ASH GROUP



PROJECT STRUCTURE

| XLM-Roberta Named Entity Recognition

How to start

Paths and Parameters

Install Libraries

Download Datasets

Import Libraries

Load Data

EN Dataframes

FA Dataframes

Create tags table

Create Proper Input Data

Create Individual Test File

XLM-Roberta

TokenFromSubtoken

Conditional Random Field

CRFLayer

NERModel

NERTokenizer

NERDataset

NERWrapper

NER

Train, and Test Model

Train

Save Model

Test with Interface

MULTI LINGUAL NER

“Fine-tune XLM Roberta model to support English and Persian language named entity recognition.

MULTI LINGUAL NER

- *Given a sequence of tokens (words, and possibly punctuation marks), provide a tag from a predefined tag set for each token in the sequence.*

Tags :

```
[ 'I-CW', 'I-CORP', 'B-PROD', 'I-PROD',  
  'I-GRP', 'O', 'B-LOC', 'B-GRP',  
  'B-CORP', 'B-CW', 'I-PER', 'B-PER',  
  'I-LOC' ]
```

PATH AND PARAMETERS

- If Debug == True then fine-tuning will be executed with limited partition of data specified with MAX_EXAMPLES.
- Other global parameters are well commented.

```
1 DEBUG = True # Debug mode: If True then training set will be limited to MAX_EXAMPLES
2 MAX_EXAMPLES = 1000 # Size of Debug Set
3 MAX_EPOCHS = 15 # Max Number of epochs
4 CHECKPOINT_PATH = "./drive/MyDrive" # Root directory of checkpoints
5 MODEL_PATH = "./model" # Root directory of pytorch model
6 MODEL_NAME = "model.pt" # name of final pytorch model
```

DOWNLOAD DATASET

- Create data folder and download datasets

```
1 # Create data folder and download required datasets
2 ! mkdir data
3 ! wget https://github.com/language-ml/4-token-classification/raw/main/Multilingual-NER/en_test.csv -P ./data
4 ! wget https://github.com/language-ml/4-token-classification/raw/main/Multilingual-NER/en_train.csv -P ./data
5 ! wget https://github.com/language-ml/4-token-classification/raw/main/Multilingual-NER/fa_test.csv -P ./data
6 ! wget https://github.com/language-ml/4-token-classification/raw/main/Multilingual-NER/fa_train.csv -P ./data
```

LIBRARY

List of necessary libraries.

```
1 # import primitive libraries
2 import os
3 import pandas as pd
4 from tqdm import tqdm
5 import numpy as np
6 import json
7
8 # import seqeval to report classifier performance metrics
9 from seqeval.metrics import accuracy_score, precision_score, recall_score, f1_score
10 from seqeval.scheme import IOB2
11
12 # import torch related modules
13 import torch
14 from torch.utils.data import DataLoader
15 from torch.utils.data import Dataset
16 from torch.nn.utils.rnn import pad_sequence
17 import torch.nn as nn
18
19 # import pytorch lightning library
20 import pytorch_lightning as pl
21 from torchcrf import CRF as SUPERCRF
22
23 # import NLTK to create better tokenizer
24 import nltk
25 from nltk.tokenize import RegexpTokenizer
26
27 # Transformers : Roberta Model
28 from transformers import XLMRobertaTokenizerFast
29 from transformers import XLMRobertaModel, XLMRobertaConfig
30
31 # import sklearn inorder to split data into train-evaluation-test
32 from sklearn.model_selection import GroupShuffleSplit
33
34 # import Typings
35 from typing import Union, Dict, List, Tuple, Any, Optional
```

Load Data

```
1 # A function to load dataset into dataframe
2 def load_data(name, path='./data', test=False):
3     print(f'Processing {name}')
4     # Read CSV
5     df = pd.read_csv(path + '/' + name)
6     # Define columns
7     df.columns = ['index', 'Token', 'Tag']
8     df.set_index('index', drop=True, inplace=True)
9     # Remove rows which starts with # id
10    mask = df.Token.str.startswith('# id')
11    indices = list(map(lambda x:x+1, df.index[mask].tolist())) + [df.index[-1]]
12    # Assign Sent number to each row
13    sentence = 0
14    df['Sent'] = None
15    for index in tqdm(range(len(indices)-1)):
16        df.loc[indices[index]:indices[index+1], 'Sent'] = sentence
17        sentence += 1
18    df.drop(df.index[mask], inplace=True)
19    # Drop Tag column if Dataframe is for test
20    if test:
21        df.drop(columns='Tag', inplace=True)
22    return df
23
```

XLM-Roberta Named Entity Recognition

How to start

Paths and Parameters

Install Libraries

Download Datasets

Import Libraries

Load Data

EN Dataframes

FA Dataframes

Create tags table

Create Proper Input Data

Create Individual Test File

XLM-Roberta

TokenFromSubtoken

Conditional Random Field

CRFLayer

NERModel

NERTokenizer

NERDataset

NERWrapper

NER

Train, and Test Model

Train

Save Model

Test with Interface

Load Data

```
1 # EN Data
2 en_data = load_data('en_train.csv')
3 # clean tag column
4 en_data['Tag'] = en_data['Tag'].apply(lambda tag: tag.strip() if tag != ' _ 0' else '0')
5 # clean Token column
6 en_data['Token'] = en_data['Token'].apply(lambda token: token.strip())
7
8 # EN Test (Deploy)
9 en_deploy_test = load_data('en_test.csv', test=True)
10
11 # FA Data
12 fa_data = load_data('fa_train.csv')
13 # clean tag column
14 fa_data['Tag'] = fa_data['Tag'].apply(lambda tag: tag.strip() if tag != ' _ 0' else '0')
15 # clean Token column
16 fa_data['Token'] = fa_data['Token'].apply(lambda token: token.strip())
17
18 # FA TEST (Deploy)
19 fa_deploy_test = load_data('fa_test.csv', test=True)
```

XLM-Roberta Named Entity Recognition

How to start

Paths and Parameters

Install Libraries

Download Datasets

Import Libraries

Load Data

EN Dataframes

FA Dataframes

Create tags table

Create Proper Input Data

Create Individual Test File

XLM-Roberta

TokenFromSubtoken

Conditional Random Field

CRFLayer

NERModel

NERTokenizer

NERDataset

NERWrapper

NER

Train, and Test Model

Train

Save Model

Test with Interface

Preparing Data

- Preparing data
 - Train
 - Test
 - Evaluate

GroupShuffleSplit would proportionally split data based on **VAL_TEST_SIZE** and **TEST_SIZE** which are ratio of Val set + Test set to training set and Test set to Val set.

```
1 # Split Dataframe into Train and Val+Test
2 splitter = GroupShuffleSplit(test_size=VAL_TEST_SIZE, n_splits=1, random_state = 42)
3 split = splitter.split(en_data, groups=en_data['Sent'])
4 en_train_inds, en_val_test_inds = next(split)
5
6 # En Train and Val+Test
7 en_train = en_data.iloc[en_train_inds]
8 en_val_test = en_data.iloc[en_val_test_inds]
9
10 # Split Dataframe into Val and Test
11 splitter = GroupShuffleSplit(test_size=TEST_SIZE, n_splits=1, random_state = 42)
12 split = splitter.split(en_val_test, groups=en_val_test['Sent'])
13 en_val_inds, en_test_inds = next(split)
14
15 # En Val and Test
16 en_val = en_val_test.iloc[en_val_inds]
17 en_test = en_val_test.iloc[en_test_inds]
```

Preparing Data

- We need to create suitable training and test input files in order to feed into model.

```
1 # Create proper input files
2 def write_file(en_df, fa_df, name, path='./data'):
3     # Get the first sentence. The parameter will be updated during execution of function and show the last processed sentence number.
4     last_sent = en_df.iloc[0].Sent
5
6     with open(f'{path}/{name}.txt', 'a') as writer:
7         # Write English (Token, Tag) pairs
8         for index, row in tqdm(en_df.iterrows(), desc = f'Writing English {name} file'):
9             token, tag, current_sent = row
10
11             if (current_sent == last_sent):
12                 writer.write(f"{token}\t{tag}\n")
13             else:
14                 last_sent = current_sent
15                 writer.write(f"\n")
16                 writer.write(f"{token}\t{tag}\n")
17
18         # Write Farsi (Token, Tag) pairs
19         last_row_df = fa_df.iloc[-1].Sent
20
21         for index, row in tqdm(fa_df.iterrows(), desc = f'Writing Farsi {name} file'):
22             token, tag, current_sent = row
23
24             if current_sent!=last_row_df:
25                 if (current_sent == last_sent):
26                     writer.write(f"{token}\t{tag}\n")
27
28                 else:
29                     last_sent = current_sent
30                     writer.write(f"\n")
31                     writer.write(f"{token}\t{tag}\n")
32             else:
33                 # Last Sentence of file without newline char
34                 writer.write(f"\n")
35                 writer.write(f"{token}\t{tag}")
```

Activate

Create Test File

- Write Individual Persian and English test files to test final model.

```
1 def write_individual_test(df, name, path='./data'):
2     # Get the first sentence. The parameter will be updated during execution c
3     last_sent = df.iloc[0].Sent
4
5     with open(f'{path}/{name}.txt', 'a') as writer:
6         last_row_df = df.iloc[-1].Sent
7
8         for index, row in tqdm(df.iterrows(), desc = f'Writing {name} file'):
9             token, tag, current_sent = row
10
11             if current_sent!=last_row_df:
12                 if (current_sent == last_sent):
13                     writer.write(f"{token}\t{tag}\n")
14
15                 else:
16                     last_sent = current_sent
17                     writer.write(f"\n")
18                     writer.write(f"{token}\t{tag}\n")
19             else:
20                 # Last Sentence of file without newline char
21                 writer.write(f"\n")
22                 writer.write(f"{token}\t{tag}")
```

XLM Roberta

Following subsections show different parts of out XLM-Roberta model which is implemented with help of Pytorch and Pytorch Lightning.

XLM-Roberta Named Entity Recognition

How to start

Paths and Parameters

Install Libraries

Download Datasets

Import Libraries

Load Data

EN Dataframes

FA Dataframes

Create tags table

Create Proper Input Data

XLM-Roberta

TokenFromSubtoken

Conditional Random Field

CRFLayer

NERModel

NERTokenizer

NERDataset

NERWrapper

NER

Train, Save and Test Model

Train

Save Model

Test

Test with Interface

Token From Subtoken

Pad, Reshape and finally convert to tensor each units and masks.

```
1 class TokenFromSubtoken(torch.nn.Module):
2
3     def forward(self, units: torch.Tensor, mask: torch.Tensor) -> torch.Tensor:
4
5         device = units.device
6         nf_int = units.size()[-1]
7         batch_size = units.size()[0]
8
9         # number of TOKENS in each sentence
10        token_seq_lengths = torch.sum(mask, 1).to(torch.int64)
11        # number of words
12        n_words = torch.sum(token_seq_lengths)
13        # max token seq len
14        max_token_seq_len = torch.max(token_seq_lengths)
15
16        idxs = torch.stack(torch.nonzero(mask, as_tuple=True), dim=1)
17        # padding is for computing change from one sample to another in the batch
18        sample_ids_in_batch = torch.nn.functional.pad(input=idxs[:, 0], pad=[1, 0])
19
20        a = (~torch.eq(sample_ids_in_batch[1:], sample_ids_in_batch[:-1])).to(torch.int64)
21
22        # transforming sample start masks to the sample starts themselves
23        q = a * torch.arange(n_words, device=device).to(torch.int64)
24        count_to_substract = torch.nn.functional.pad(torch.masked_select(q, q.to(torch.bool)), [1, 0])
25
26        new_word_indices = torch.arange(n_words, device=device).to(torch.int64) - count_to_substract[torch.cumsum(a, 0)]
27
28        n_total_word_elements = max_token_seq_len*torch.ones_like(token_seq_lengths, device=device).sum()
29        word_indices_flat = (idxs[:, 0] * max_token_seq_len + new_word_indices).to(torch.int64)
30        #x_mask = torch.sum(torch.nn.functional.one_hot(word_indices_flat, n_total_word_elements), 0)
31        #x_mask = x_mask.to(torch.bool)
32        x_mask = torch.zeros(n_total_word_elements, dtype=torch.bool, device=device)
33        x_mask[word_indices_flat] = torch.ones_like(word_indices_flat, device=device, dtype=torch.bool)
34        # to get absolute indices we add max_token_seq_len:
35
36        nonword_indices_flat = (~x_mask).nonzero().squeeze(-1)
37
38        elements = units[mask.bool()]
39
40        paddings = torch.zeros_like(nonword_indices_flat, dtype=elements.dtype).unsqueeze(-1).repeat(1,nf_int).to(device)
41
42        # tensor_flat -> [x, x, x, x, x, 0, x, 0, 0]
43        tensor_flat_unordered = torch.cat([elements, paddings])
44        _, order_idx = torch.sort(torch.cat([word_indices_flat, nonword_indices_flat]))
45        tensor_flat = tensor_flat_unordered[order_idx]
46
47        tensor = torch.reshape(tensor_flat, (-1, max_token_seq_len, nf_int))
48
49        return tensor
```

CRF

- CRF Layer implemented with help of **TorchCRF** library, but we need to override **_Viterbi_decode** to make it more compatible with our model.

```
1 class CRF(SUPERCRF):
2
3     # override viterbi decoder in order to make it compatible with our code
4     def _viterbi_decode(self, emissions: torch.FloatTensor,
5                         mask: torch.ByteTensor) -> List[List[int]]:
6         # emissions: (seq_length, batch_size, num_tags)
7         # mask: (seq_length, batch_size)
8         assert emissions.dim() == 3 and mask.dim() == 2
9         assert emissions.shape[:2] == mask.shape
10        assert emissions.size(2) == self.num_tags
11        assert mask[0].all()
12
13        seq_length, batch_size = mask.shape
14
15        # Start transition and first emission
16        # shape: (batch_size, num_tags)
17        score = self.start_transitions + emissions[0]
18        history = []
19
20
21        for i in range(1, seq_length):
22            # Broadcast viterbi score for every possible next tag
23            # shape: (batch_size, num_tags, 1)
24            broadcast_score = score.unsqueeze(2)
25
26            # Broadcast emission score for every possible current tag
27            # shape: (batch_size, 1, num_tags)
28            broadcast_emission = emissions[i].unsqueeze(1)
29
30            # shape: (batch_size, num_tags, num_tags)
31            next_score = broadcast_score + self.transitions + broadcast_emission
32
33            # shape: (batch_size, num_tags)
34            next_score, indices = next_score.max(dim=1)
35
36
37            # shape: (batch_size, num_tags)
38            score = torch.where(mask[i].unsqueeze(1), next_score, score)
39            history.append(indices)
40
```

CRF Layer

- Forward
 - Decide output logits based on backbone network
- Backward
 - decode based on CRF weights

```
1 class CRFLayer(nn.Module):
2     def __init__(self, embedding_size, n_labels):
3
4         super(CRFLayer, self).__init__()
5         self.dropout = nn.Dropout(0.1)
6         self.output_dense = nn.Linear(embedding_size, n_labels)
7         self.crf = CRF(n_labels, batch_first=True)
8         self.token_from_subtoken = TokenFromSubtoken()
9
10        # Forward: decide output logits based on backbone network
11        def forward(self, embedding, mask):
12            logits = self.output_dense(self.dropout(embedding))
13            logits = self.token_from_subtoken(logits, mask)
14            pad_mask = self.token_from_subtoken(mask.unsqueeze(-1), mask).squeeze(-1).bool()
15            return logits, pad_mask
16
17        # Decode: decode based on CRF weights
18        def decode(self, logits, pad_mask):
19            return self.crf.decode(logits, pad_mask)
20
21        # Evaluation Loss: calculate mean log likelihood of CRF layer
22        def eval_loss(self, logits, targets, pad_mask):
23            mean_log_likelihood = self.crf(logits, targets, pad_mask, reduction='sum').mean()
24            return -mean_log_likelihood
25
```


NER Model

Following class include pretrained ***XLMLRobertaModel*** and implements basic necessary methods like predict, decode, loss calculation and freezing specific number of layers.

```
1 class NERModel(nn.Module):
2
3     def __init__(self, n_labels:int, roberta_path:str):
4         super(NERModel,self).__init__()
5         self.roberta = XLMLRobertaModel.from_pretrained(roberta_path)
6         self.crf = CRFLayer(self.roberta.config.hidden_size, n_labels)
7
8     # Forward: pass embeddings to CRF layer in order to evaluate logits from subword sequence
9     def forward(self,
10                 input_ids:torch.Tensor,
11                 attention_mask:torch.Tensor,
12                 token_type_ids:torch.Tensor,
13                 mask:torch.Tensor) -> torch.Tensor:
14
15         embedding = self.roberta(input_ids=input_ids,
16                                   attention_mask=attention_mask,
17                                   token_type_ids=token_type_ids)[0]
18         logits, pad_mask = self.crf(embedding, mask)
19         return logits, pad_mask
20
21     # Disable Gradient and Predict with model
22     @torch.no_grad()
23     def predict(self, inputs:Tuple[torch.Tensor]) -> torch.Tensor:
24         input_ids, attention_mask, token_type_ids, mask = inputs
25         logits, pad_mask = self(input_ids, attention_mask, token_type_ids, mask)
26         decoded = self.crf.decode(logits, pad_mask)
27         return decoded, pad_mask
28
29     # Decode: pass to crf decoder and decode based on CRF weights
30     def decode(self, logits, pad_mask):
31         """Decode logits using CRF weights
32         """
33         return self.crf.decode(logits, pad_mask)
34
35     # Evaluation Loss: pass to crf eval_loss and calculate mean log likelihood of CRF layer
36     def eval_loss(self, logits, targets, pad_mask):
37         return self.crf.eval_loss(logits, targets, pad_mask)
38
39     # Determine number of layers to be fine-tuned (!freeze)
40     def freeze_roberta(self, n_freeze:int=6):
41         for param in self.roberta.parameters():
42             param.requires_grad = False
43
44         for param in self.roberta.encoder.layer[n_freeze:].parameters():
45             param.requires_grad = True
```

NER Tokenizer

Following class inherits benefits of **NLTK** tokenizer and **XLNetRobertaTokenizer** in order to create better tokenizer which is responsible to tokenize input batches.

```
1 class NERTokenizer(object):
2
3     MAX_LEN=512
4     BATCH_LENGTH_LIMT = 380 # Max number of roberta tokens in one sentence.
5
6     # Modified version of http://stackoverflow.com/questions/36353125/nltk-regular-expression-tokenizer
7     PATTERN = r'''(?x)          # set flag to allow verbose regexps
8         (?:[A-Z]\.)+          # abbreviations, e.g. U.S.A. or U.S.A #
9         | (?:\d+\.)          # numbers
10        | \w+(?:[-.]\w+)*      # words with optional internal hyphens
11        | \$?\d+(?:\.\d+)?%?   # currency and percentages, e.g. $12.40, 82%
12        | \.\.\.              # ellipsis, and special chars below, includes ], [
13        | [-\]\[\.,;'"?():_`"/\^_{}|~!@&$]
14    '''
15
16    def __init__(self, base_model:str, to_device:str='cpu'):
17        super(NERTokenizer,self).__init__()
18        self.roberta_tokenizer = XLNetRobertaTokenizerFast.from_pretrained(base_model, do_lower_case=False)
19        self.to_device = to_device
20
21        self.word_tokenizer = RegexpTokenizer(self.PATTERN)
22        self.sent_tokenizer = nltk.data.load('tokenizers/punkt/english.pickle')
23
24    # tokenize batch of tokens
25    def tokenize_batch(self, inputs, pad_to = None) -> torch.Tensor:
26        batch = [inputs] if isinstance(inputs[0], str) else inputs
27
28        input_ids, attention_mask, token_type_ids, mask = [], [], [], []
29        for tokens in batch:
30            input_ids_tmp, attention_mask_tmp, token_type_ids_tmp, mask_tmp = self._tokenize_words(tokens)
31            input_ids.append(input_ids_tmp)
32            attention_mask.append(attention_mask_tmp)
33            token_type_ids.append(token_type_ids_tmp)
34            mask.append(mask_tmp)
35
36        input_ids = pad_sequence(input_ids, batch_first=True, padding_value=self.roberta_tokenizer.pad_token_id)
37        attention_mask = pad_sequence(attention_mask, batch_first=True, padding_value=0)
38        token_type_ids = pad_sequence(token_type_ids, batch_first=True, padding_value=0)
39        mask = pad_sequence(mask, batch_first=True, padding_value=0)
```

NER Dataset

Any Pytorch Model Fine-tuning needs a good Pytorch dataset with `__len__` and `__getitem__` conventional methods.

```
1 # Pytorch Dataset
2 class NERDataset(Dataset):
3
4     def __init__(self, data_path:str, label_tags: List[str], base_model:str="xlm-roberta-base",
5                 default_label=0, max_length:int=512, to_device="cpu"):
6         self.tokenizer = NERTokenizer(base_model=base_model, to_device=to_device)
7         self.label_tags = label_tags
8         self.name_to_label = {x: i for i, x in enumerate(self.label_tags)}
9         self.default_label = default_label
10        self.max_length = max_length
11
12
13        # open file (train, test or val)
14        with open(data_path, 'r') as f:
15            data_text = f.read()
16        self.data = []
17
18        # the loop notices the change of sentence with double newline
19        for sentence in filter(lambda x: len(x)>2, data_text.split('\n\n')):
20            sample = []
21            # each word laid in sepertaed lines
22            for wordline in sentence.split('\n'):
23                if wordline=='':
24                    continue
25                # the word and label are seperated from each other with tab
26                word, label = wordline.split('\t')
27                sample.append((word, label))
28            self.data.append(sample)
29
30        # len of dataset
31        def __len__(self):
32            return len(self.data)
33
34
35        def __getitem__(self, idx):
36            item = self.data[idx]
37            words, labels = list(zip(*item))
38
39            labels_idx = [self.name_to_label.get(x, self.default_label) for x in labels]
40            y = torch.tensor(labels_idx, dtype=torch.long)
41            diff = self.max_length - y.shape[-1]
42            y = torch.nn.functional.pad(y, (0, diff), value=self.default_label)
43            X = self.tokenizer.tokenize_batch(list(words), pad_to=self.max_length)
44
45            return X, y
```

NER Wrapper

Following Class inherits from Pytorch Lightning module class which will help us to pass this class to the **trainer** object and train model with predefined methods in this class.

```
1 # Lightning Wrapper for NER Model
2
3 class NERWrapper(pl.LightningModule):
4     def __init__(self,
5                  learning_rate = 2e-5,
6                  weight_decay = 0.0,
7                  batch_size = 16,
8                  freeze_layers = 8,
9                  tags = tags_table,
10                 train_path = "./data/train.txt",
11                 val_path = "./data/val.txt",
12                 test_path = "./data/test.txt",
13                 pretrained_path = None,
14                 *args, **kwargs
15             ):
16
17         super(NERWrapper, self).__init__()
18         self.save_hyperparameters('learning_rate', 'weight_decay', 'batch_size')
19         self.tags, self.train_path, self.val_path, self.test_path = tags, train_path, val_path, test_path
20         self.model = NERModel(n_labels=len(self.tags), roberta_path="xlm-roberta-base")
21
22         if pretrained_path is not None:
23             self.model.load_state_dict(torch.load(pretrained_path))
24             self.model.freeze_roberta(freeze_layers)
25
26     def forward(self, *args, **kwargs):
27         return self.model.forward(*args, **kwargs)
28
29     def _step(self, batch, batch_idx):
30         (input_ids, attention_mask, token_type_ids, mask), labels = batch
31         logits, pad_mask = self.model(input_ids, attention_mask, token_type_ids, mask)
32         labels = labels[:, :logits.shape[1]]
33         loss = self.model.eval_loss(logits, labels, pad_mask)
34         preds_tag_idx = self.model.decode(logits, pad_mask)
35         preds_tag = [[self.tags[start.item()] for m, start in zip(mask, sample) if m] for mask, sample in zip(pad_mask, preds_tag_idx)]
36         labels_tag = [[self.tags[start.item()] for m, start in zip(mask, sample) if m] for mask, sample in zip(pad_mask, labels)]
37         tensorboard_logs = {'batch_loss': loss}
38         for metric, value in tensorboard_logs.items():
39             self.log(metric, value, prog_bar=True)
40         return {'loss': loss, 'preds': preds_tag, 'labels': labels_tag}
```

NER Interface

You can use this interface and pass a raw string and receive NER tags in response.

```
1 class NER(object):
2
3     def __init__(self, model_path, model_name = MODEL_NAME, tags = tags_table):
4
5         self.tags = tags
6         self.device = "cuda" if torch.cuda.is_available() else "cpu"
7         roberta_path = "xlm-roberta-base"
8         self.model = NERModel(n_labels=len(self.tags), roberta_path=roberta_path).to(self.device)
9         state_dict = torch.load(os.path.join(model_path, model_name))
10        self.model.load_state_dict(state_dict, strict=False)
11        self.model.eval()
12        self.tokenizer = NERTokenizer(base_model=roberta_path, to_device=self.device)
13
14    @torch.no_grad()
15    def __call__(self, raw_text):
16
17        outputs_flat, spans_flat, entities = [], [], []
18        for batch, words, spans in self.tokenizer.prepare_row_text(raw_text):
19            output, pad_mask = self.model.predict(batch)
20            outputs_flat.extend(output[pad_mask.bool()].reshape(-1).tolist())
21            spans_flat += sum(spans, [])
22
23        for tag_idx, (start, end) in zip(outputs_flat, spans_flat):
24            tag = self.tags[tag_idx]
25            if tag != 'O':
26                entities.append({'Text': raw_text[start:end], 'Tag': tag})
27
28        return entities
```

Train and Test Model

We train model with training dataset and we use validation test on predefined steps and finally output model will test with test datasets and result will be displayed

XLM-Roberta Named Entity Recognition

How to start

Paths and Parameters

Install Libraries

Download Datasets

Import Libraries

Load Data

EN Dataframes

FA Dataframes

Create tags table

Create Proper Input Data

Create Individual Test File

XLM-Roberta

TokenFromSubtoken

Conditional Random Field

CRFLayer

NERModel

NERTokenizer

NERDataset

NERWrapper

NER

Train, and Test Model

Train

Save Model

Test with Interface

Train + Test

Following method will train and test over whole data.

```
1 # train function
2 def train_test(test_path):
3     plmodel = NERWrapper(test_path = test_path)
4     trainer = pl.Trainer(default_root_dir=CHECKPOINT_PATH, gpus=1, max_epochs=MAX_EPOCHS)
5     trainer.fit(plmodel)
6     result = trainer.test(plmodel)
7     print(result)
```

	Name	Type	Params
0	model	NERModel	278 M
28.4 M	Trainable params		
249 M	Non-trainable params		
278 M	Total params		

NER Interface

Here is example of how you can use NER interface which has been explained before.

```
1 # del ner
2 ner = NER(model_path = MODEL_PATH, model_name = MODEL_NAME, tags = tags_table) # Load pretrained model.

1 text = """"his playlist includes sonny sharrock , gza , country teasers and the notorious b.i.g.""
2 result = ner(text)
3 result
```


EVALUATION

Here are final evaluation result of our model after hours and hours of sitting behind the monitor proving google colab that we are not robots.

```
DATALOADER:0 TEST RESULTS
{'batch_loss': 24.60813919067383,
 'test_F1': 0.7578543,
 'test_accuracy': 0.9484375,
 'test_loss': 60.608143005371094,
 'test_precision': 0.772465,
 'test_recall': 0.7423492}
```