

Project: Text Classification using Neural Network

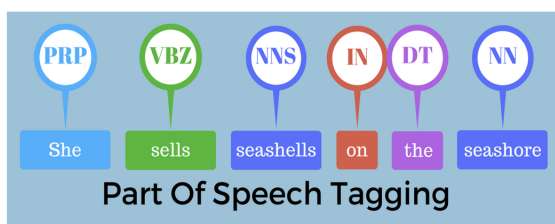
Ngày 29 tháng 10 năm 2023

Phần 1. Giới thiệu

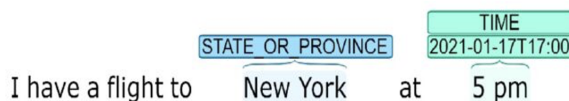
Phân loại văn bản (Text Classification) là một trong những bài toán cơ bản trong xử lý ngôn ngữ tự nhiên. Với yêu cầu của bài toán gán nhãn cho một đơn vị văn bản (văn bản, đoạn văn bản, câu, từ,...) với nhãn được xem xét từ một tập nhãn cho trước.

Phân loại các bài toán phân loại dựa vào đơn vị văn bản cần phân loại có:

- Token-level Classification: gán nhãn cho các đơn vị văn bản là các từ (token). Các bài toán điển hình trong nhóm này là: gán nhãn từ loại (Part-of-Speech Tagging) gán mỗi từ trong câu với các loại từ: danh từ, động từ,... và nhận dạng thực thể (Named Entity Recognition) gán nhãn cho các từ thuộc một thực thể nhất định như người, địa điểm,...

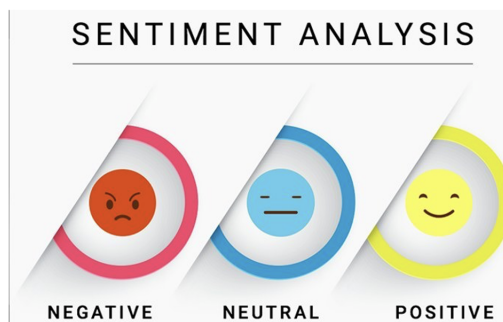
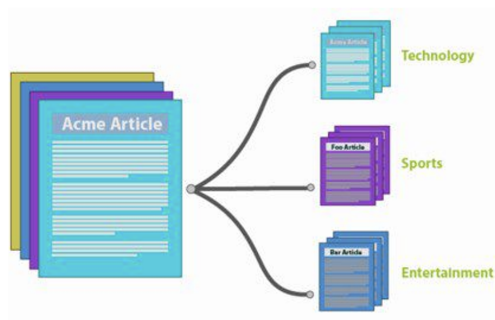


Named Entity Recognition



Hình 1: Bài toán gán nhãn từ loại (POS Tagging) và nhận dạng thực thể (NER).

- Document-level Classification: gán nhãn cho các đơn vị văn bản là các cụm từ, câu, đoạn văn. Các bài toán điển hình gồm có phân loại văn bản vào kinh tế, chính trị,... Phân tích cảm xúc (Sentiment Analysis): phân loại văn bản thành nhãn tích cực, tiêu cực hoặc trung tính.



Hình 2: Phân loại văn bản: kinh tế, chính trị,... và phân tích cảm xúc (Sentiment Analysis).

Phân loại các bài toán phân loại dựa vào số nhãn:

- Binary Classification: phân loại với số nhãn là 2. Ví dụ bài toán phân tích cảm xúc
- Multiclass Classification: phân loại với số nhãn nhiều hơn 2. Ví dụ: kinh tế, chính trị, văn hoá,...
- Multilabel Classification: phân loại một đơn vị văn bản có thể thuộc vào một hoặc nhiều nhãn.

Bài toán phân loại có nhiều ứng dụng quan trọng. Ứng dụng trong phân tích cảm xúc bình luận khách hàng với các sản phẩm. Phân loại các thông tin giả mạo,...

Phần 2. Thực hành

Huấn luyện mô hình phân loại văn bản sử dụng mạng nơ-ron gồm các bước:

- Tải bộ dữ liệu NTC-CSV
- Tiền xử lý dữ liệu văn bản
- Biểu diễn dữ liệu văn bản thành các vectors
- Xây dựng mô hình phân loại
- Huấn luyện mô hình
- Dự đoán, đánh giá mô hình

1. Tải về bộ dữ liệu NTC-SCV

Bộ dữ liệu NTC-SCV gồm 50.000 mẫu trong đó có 30.000 mẫu cho training, 10.000 mẫu cho validation và 10.000 mẫu cho testing. Số nhãn là 2: positive (tích cực) và negative (tiêu cực).

```
1 # Download data from github
2 !git clone https://github.com/congnghia0609/ntc-scv.git
3 !unzip ./ntc-scv/data/data_test.zip -d ./data
4 !unzip ./ntc-scv/data/data_train.zip -d ./data
5 !rm -rf ./ntc-scv
6
7 # Read data from file
8 import os
9 import pandas as pd
10
11 def load_data_from_path(folder_path):
12     examples = []
13     for label in os.listdir(folder_path):
14         full_path = os.path.join(folder_path, label)
15         for file_name in os.listdir(full_path):
16             file_path = os.path.join(full_path, file_name)
17             with open(file_path, "r", encoding="utf-8") as f:
18                 lines = f.readlines()
19                 sentence = " ".join(lines)
20                 if label == "neg":
21                     label = 0
22                 if label == "pos":
23                     label = 1
24                 data = {
25                     'sentence': sentence,
26                     'label': label
27                 }
28                 examples.append(data)
29     return pd.DataFrame(examples)
30
31 folder_paths = {
32     'train': './data/data_train/train',
33     'valid': './data/data_train/test',
34     'test': './data/data_test/test'
35 }
36
37 train_df = load_data_from_path(folder_paths['train'])
38 valid_df = load_data_from_path(folder_paths['valid'])
39 test_df = load_data_from_path(folder_paths['test'])
```

2. Tiền xử lý dữ liệu

Với bộ dữ liệu NTC-SCV có một số bình luận viết bằng tiếng anh hoặc tiếng pháp và chứa các thẻ HTML, đường dẫn URLs. Tiền xử lý dữ liệu gồm 2 bước:

- Xóa bỏ những bình luận không phải tiếng việt

Sử dụng thư viện langid được mô tả như hình sau:



Hình 3: Phát hiện ngôn ngữ sử dụng thư viện langid.

```

1 # Install library
2 !pip install langid
3
4 from langid.langid import LanguageIdentifier, model
5 def identify_vn(df):
6     identifier = LanguageIdentifier.from_modelstring(model, norm_probs=True)
7     not_vi_idx = set()
8     THRESHOLD = 0.9
9     for idx, row in df.iterrows():
10         score = identifier.classify(row["sentence"])
11         if score[0] != "vi" or (score[0] == "vi" and score[1] <= THRESHOLD):
12             not_vi_idx.add(idx)
13     vi_df = df[~df.index.isin(not_vi_idx)]
14     not_vi_df = df[df.index.isin(not_vi_idx)]
15     return vi_df, not_vi_df
16
17 train_df_vi, train_df_other = identify_vn(train_df)
  
```

- Làm sạch dữ liệu

Các bước tiền làm sạch liệu:

- Xóa bỏ thẻ HTML, đường dẫn URL
- Xóa bỏ dấu câu, số
- Xóa bỏ các ký tự đặc biệt, emoticons,...
- Chuẩn hoá khoảng trắng
- Chuyển sang viết thường

```

1 import re
2 import string
3
4 def preprocess_text(text):
5
6     url_pattern = re.compile(r'https?://\s+\www\.\s+')
7     text = url_pattern.sub(r" ", text)
  
```

```

8
9     html_pattern = re.compile(r'<[<>]+>')
10    text = html_pattern.sub(" ", text)
11
12    replace_chars = list(string.punctuation + string.digits)
13    for char in replace_chars:
14        text = text.replace(char, " ")
15
16    emoji_pattern = re.compile("[
17        u"\U0001F600-\U0001F64F" # emoticons
18        u"\U0001F300-\U0001F5FF" # symbols & pictographs
19        u"\U0001F680-\U0001F6FF" # transport & map symbols
20        u"\U0001F1E0-\U0001F1FF" # flags (iOS)
21        u"\U0001F1F2-\U0001F1F4" # Macau flag
22        u"\U0001F1E6-\U0001F1FF" # flags
23        u"\U0001F600-\U0001F64F"
24        u"\U00002702-\U000027B0"
25        u"\U000024C2-\U0001F251"
26        u"\U0001f926-\U0001f937"
27        u"\U0001F1F2"
28        u"\U0001F1F4"
29        u"\U0001F620"
30        u"\u200d"
31        u"\u2640-\u2642"
32    "]" + "", flags=re.UNICODE)
33    text = emoji_pattern.sub(r" ", text)
34
35    text = " ".join(text.split())
36
37    return text.lower()
38
39 train_df_vi['preprocess_sentence'] = [
40     preprocess_text(row['sentence']) for index, row in train_df_vi.iterrows()
41 ]
42 valid_df['preprocess_sentence'] = [
43     preprocess_text(row['sentence']) for index, row in valid_df.iterrows()
44 ]
45 test_df['preprocess_sentence'] = [
46     preprocess_text(row['sentence']) for index, row in test_df.iterrows()
47 ]

```

3. Biểu diễn dữ liệu văn bản thành các vectors

Để biểu diễn dữ liệu văn bản thành các đặc trưng (vectors), chúng ta sử dụng thư viện torchtext.

```

1 !pip install -q torchtext==0.16.0
2
3 # word-based tokenizer
4 from torchtext.data.utils import get_tokenizer
5 tokenizer = get_tokenizer("basic_english")
6
7 # create iter dataset
8 def yield_tokens(sentences, tokenizer):
9     for sentence in sentences:
10         yield tokenizer(sentence)
11
12 # build vocabulary
13 from torchtext.vocab import build_vocab_from_iterator
14
15 vocab_size = 10000
16 vocabulary = build_vocab_from_iterator(

```

```

17     yield_tokens(train_df_vi['preprocess_sentence'], tokenizer),
18     max_tokens=vocab_size,
19     specials=["<unk>"]
20 )
21 vocabulary.set_default_index(vocabulary["<unk>"])
22
23 # convert iter into torchtext dataset
24 from torchtext.data.functional import to_map_style_dataset
25 def prepare_dataset(df):
26     for index, row in df.iterrows():
27         sentence = row['preprocess_sentence']
28         encoded_sentence = vocabulary(tokenizer(sentence))
29         label = row['label']
30         yield encoded_sentence, label
31
32 train_dataset = prepare_dataset(train_df_vi)
33 train_dataset = to_map_style_dataset(train_dataset)
34
35 valid_dataset = prepare_dataset(valid_df)
36 valid_dataset = to_map_style_dataset(valid_dataset)


1 import torch
2
3 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
4
5 def collate_batch(batch):
6     encoded_sentences, labels, offsets = [], [], [0]
7     for encoded_sentence, label in batch:
8         labels.append(label)
9         encoded_sentence = torch.tensor(encoded_sentence, dtype=torch.int64)
10        encoded_sentences.append(encoded_sentence)
11        offsets.append(encoded_sentence.size(0))
12
13    labels = torch.tensor(labels, dtype=torch.int64)
14    offsets = torch.tensor(offsets[:-1]).cumsum(dim=0)
15    encoded_sentences = torch.cat(encoded_sentences)
16    return encoded_sentences.to(device), offsets.to(device), labels.to(device)
17
18 from torch.utils.data import DataLoader
19
20 batch_size = 128
21 train_dataloader = DataLoader(
22     train_dataset,
23     batch_size=batch_size,
24     shuffle=True,
25     collate_fn=collate_batch
26 )
27 valid_dataloader = DataLoader(
28     valid_dataset,
29     batch_size=batch_size,
30     shuffle=False,
31     collate_fn=collate_batch
32 )

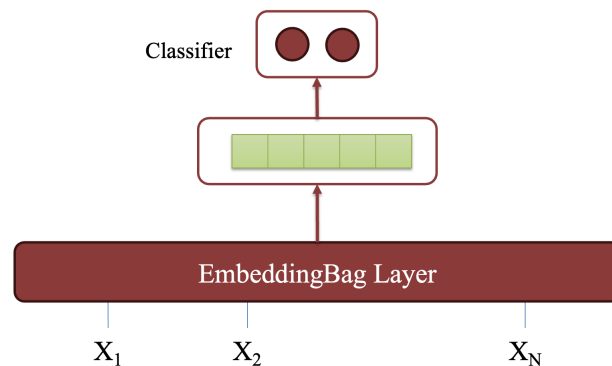
```

4. Xây dựng mô hình phân loại

- Xây dựng mô hình gồm các lớp sau:

1. EmbeddingBag: chuyển các từ với index tương ứng sang biểu diễn vector có chiều embedding_dim, sau đó tính giá trị trung bình các vector trong câu tạo thành vector đại diện cho mỗi văn bản.

2. Linear: chuyển các vector đại diện cho mỗi văn bản thành số node đầu ra là 2 node cho bài toán phân loại 2 lớp.



Hình 4: Mô hình phân loại văn bản sử dụng mạng nơ ron.

```

1 from torch import nn
2
3 class TextClassificationModel(nn.Module):
4     def __init__(self, vocab_size, embed_dim, num_class):
5         super(TextClassificationModel, self).__init__()
6         self.embedding = nn.EmbeddingBag(vocab_size, embed_dim, sparse=False)
7         self.fc = nn.Linear(embed_dim, num_class)
8         self.init_weights()
9
10    def init_weights(self):
11        initrange = 0.5
12        self.embedding.weight.data.uniform_(-initrange, initrange)
13        self.fc.weight.data.uniform_(-initrange, initrange)
14        self.fc.bias.data.zero_()
15
16    def forward(self, inputs, offsets):
17        embedded = self.embedding(inputs, offsets)
18        return self.fc(embedded)
19
20 num_class = len(set(train_df_vi['label']))
21 vocab_size = len(vocabulary)
22 embed_dim = 256
23 model = TextClassificationModel(vocab_size, embed_dim, num_class).to(device)

```

- Định nghĩa hàm loss và optimizer

```

1 learning_rate = 5
2
3 criterion = torch.nn.CrossEntropyLoss()
4 optimizer = torch.optim.SGD(model.parameters(), lr=learning_rate)

```

5. Huấn luyện mô hình

- Định nghĩa hàm huấn luyện mô hình:

```

1 import time
2
3 def train(model, optimizer, criterion, train_dataloader, epoch=0, log_interval=25):
4     model.train()
5     total_acc, total_count = 0, 0
6     losses = []
7     start_time = time.time()

```

```

8
9     for idx, (inputs, offsets, labels) in enumerate(train_dataloader):
10         optimizer.zero_grad()
11         predictions = model(inputs, offsets)
12
13         # compute loss
14         loss = criterion(predictions, labels)
15         losses.append(loss.item())
16
17         # backward
18         loss.backward()
19         torch.nn.utils.clip_grad_norm_(model.parameters(), 0.1)
20         optimizer.step()
21         total_acc += (predictions.argmax(1) == labels).sum().item()
22         total_count += labels.size(0)
23         if idx % log_interval == 0 and idx > 0:
24             elapsed = time.time() - start_time
25             print(
26                 "| epoch {:3d} | {:5d}/{:5d} batches "
27                 "| accuracy {:.8.3f}".format(
28                     epoch, idx, len(train_dataloader), total_acc / total_count
29                 )
30             )
31             total_acc, total_count = 0, 0
32             start_time = time.time()
33
34     epoch_acc = total_acc / total_count
35     epoch_loss = sum(losses) / len(losses)
36     return epoch_acc, epoch_loss

```

- Định nghĩa hàm đánh giá

```

1 def evaluate(model, criterion, valid_dataloader):
2     model.eval()
3     total_acc, total_count = 0, 0
4     losses = []
5
6     with torch.no_grad():
7         for idx, (inputs, offsets, labels) in enumerate(valid_dataloader):
8             predictions = model(inputs, offsets)
9             loss = criterion(predictions, labels)
10            losses.append(loss)
11            total_acc += (predictions.argmax(1) == labels).sum().item()
12            total_count += labels.size(0)
13
14     epoch_acc = total_acc / total_count
15     epoch_loss = sum(losses) / len(losses)
16     return epoch_acc, epoch_loss

```

- Huấn luyện mô hình:

```

1 num_class = len(set(train_df_vi['label']))
2 vocab_size = len(vocabulary)
3 embed_dim = 100
4 model = TextClassificationModel(vocab_size, embed_dim, num_class).to(device)
5
6 learning_rate = 5
7 criterion = torch.nn.CrossEntropyLoss()
8 optimizer = torch.optim.SGD(model.parameters(), lr=learning_rate)
9
10 num_epochs = 5
11 for epoch in range(1, num_epochs+1):

```



```
12     epoch_start_time = time.time()
13     train_acc, train_loss = train(model, optimizer, criterion, train_dataloader, epoch
14 )
15     eval_acc, eval_loss = evaluate(model, criterion, valid_dataloader)
16     print("-" * 59)
17     print(
18         "| End of epoch {:3d} | Time: {:.5.2f}s | Train Accuracy {:.8.3f} | Train Loss
19         {:.8.3f} | Valid Accuracy {:.8.3f} | Valid Loss {:.8.3f} ".format(
20             epoch, time.time() - epoch_start_time, train_acc, train_loss, eval_acc,
21             eval_loss
22         )
23     )
24     print("-" * 59)
```

3. Dự đoán và đánh giá mô hình

- Định nghĩa hàm dự đoán:

```
1 def predict(text):
2     with torch.no_grad():
3         encoded = torch.tensor(vocabulary(tokenizer(text)))
4         output = model(encoded, torch.tensor([0]))
5         return output.argmax(1).item()
```

- Đánh giá độ chính xác trên tập test

```
1 predictions, labels = [], []
2 for index, row in test_df.iterrows():
3     sentence = row['preprocess_sentence']
4     label = row['label']
5     prediction = predict(sentence)
6     predictions.append(prediction)
7     labels.append(label)
8
9 sum(torch.tensor(predictions) == torch.tensor(labels))/len(labels)
```

Kết quả độ chính xác trên tập test là 87

Phần 3. Câu hỏi trắc nghiệm

Câu hỏi 1 Độ đo đánh giá hiệu suất mô hình phân loại văn bản là?

- a) Accuracy
- b) BLEU
- c) ROUGE
- d) MSE

Câu hỏi 2 Bộ dữ liệu được sử dụng cho bài toán phân loại văn bản cho tiếng việt được sử dụng trong project là?

- a) NTC-SCV
- b) IMDB-Review
- c) SQuAD
- d) WikiSum

Câu hỏi 3 Tỷ lệ bộ dữ liệu dùng cho training, validation và testing là?

- a) 0.7 : 0.2 : 0.1
- b) 0.6 : 0.2 : 0.2
- c) 0.7 : 0.1 : 0.2
- d) 0.8 : 0.1 : 0.1

Câu hỏi 4 Bước tiền xử lý dữ liệu không được sử dụng trong phân lập trình là?

- a) Xoá bỏ đường dẫn URL
- b) Xoá bỏ các thẻ HTML
- c) Xoá bỏ các từ dừng
- d) Chuyển sang viết thường

Câu hỏi 5 Thư viện langid được sử dụng để làm gì?

- a) Phát hiện ngôn ngữ cho đoạn văn bản
- b) Dịch máy
- c) Tóm tắt văn bản
- d) Phân loại văn bản

Câu hỏi 6 Phương pháp tách từ được sử dụng trong phân lập trình là?

- a) Word-based Tokenization
- b) Character-based Tokenization
- c) Subword Tokenization
- d) Sentence Tokenization

Câu hỏi 7 Kích thước bộ từ điển được sử dụng trong phân lập trình là?

- a) 20.000

- b) 10.000
- c) 15.000
- d) 25.000

Câu hỏi 8 Phương pháp biểu diễn văn bản được sử dụng trong phần lập trình là?

- a) TF-IDF
- b) Dense Representation
- c) BoW
- d) One-hot Encoding

Câu hỏi 9 Kích thước của biến output trong đoạn code sau là?

```
1 import torch.nn as nn
2
3 vocab_size = 7
4 embedding_dim = 4
5 embedding_sum = nn.EmbeddingBag(vocab_size, embedding_dim, mode='sum')
6 inputs = torch.tensor([1, 2, 4, 5, 4, 3], dtype=torch.long)
7 offsets = torch.tensor([0, 3], dtype=torch.long)
8 outputs = embedding_sum(inputs, offsets)
9 outputs.shape
```

- a) 2 x 4
- b) 7 x 4
- c) 2 x 7
- d) 4 x 7

Câu hỏi 10 Tham số 'mode' trong EmbeddingBag Layer sử dụng trong phần lập trình là?

- a) max
- b) min
- c) mean
- d) sum

Câu hỏi 11 Số node đầu ra của mô hình phân loại được định nghĩa trong phần lập trình là?

- a) 5
- b) 4
- c) 3
- d) 2

Câu hỏi 12 Phương pháp khởi tạo giá trị trọng số bias trong Linear Layer được sử dụng trong phần lập trình là?

- a) zero_()
- b) one_()
- c) random
- d) uniform()

Câu hỏi 13 Hàm sau được sử dụng để làm gì?

```
1 def evaluate(model, criterion, valid_dataloader):
2     model.eval()
3     total_acc, total_count = 0, 0
4     losses = []
5
6     with torch.no_grad():
7         for idx, (inputs, offsets, labels) in enumerate(valid_dataloader):
8             predictions = model(inputs, offsets)
9             loss = criterion(predictions, labels)
10            losses.append(loss)
11            total_acc += (predictions.argmax(1) == labels).sum().item()
12            total_count += labels.size(0)
13
14    epoch_acc = total_acc / total_count
15    epoch_loss = sum(losses) / len(losses)
16    return epoch_acc, epoch_loss
```

- a) Huấn luyện mô hình trên tập training
- b) Đánh giá mô hình trên tập validation
- c) Biểu diễn văn bản thành vectors
- d) Tiền xử lý văn bản

- Hết -