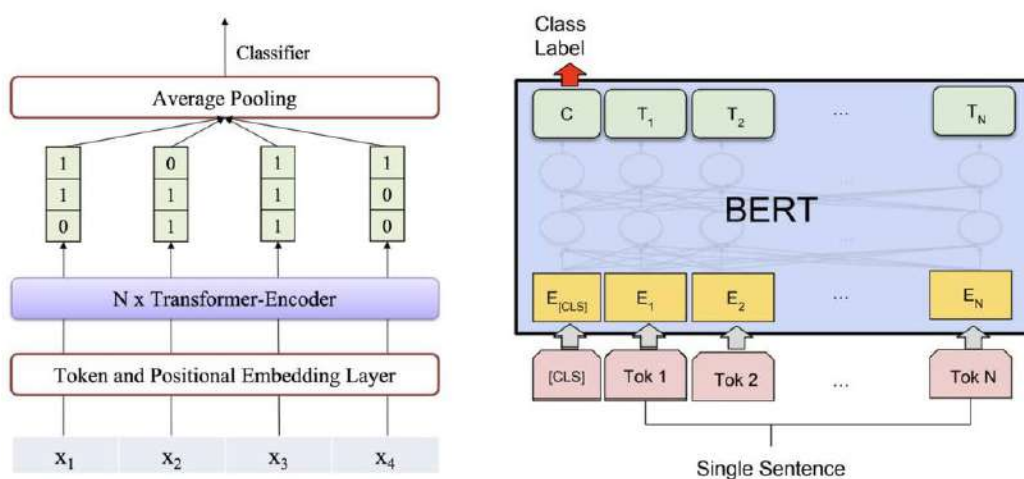


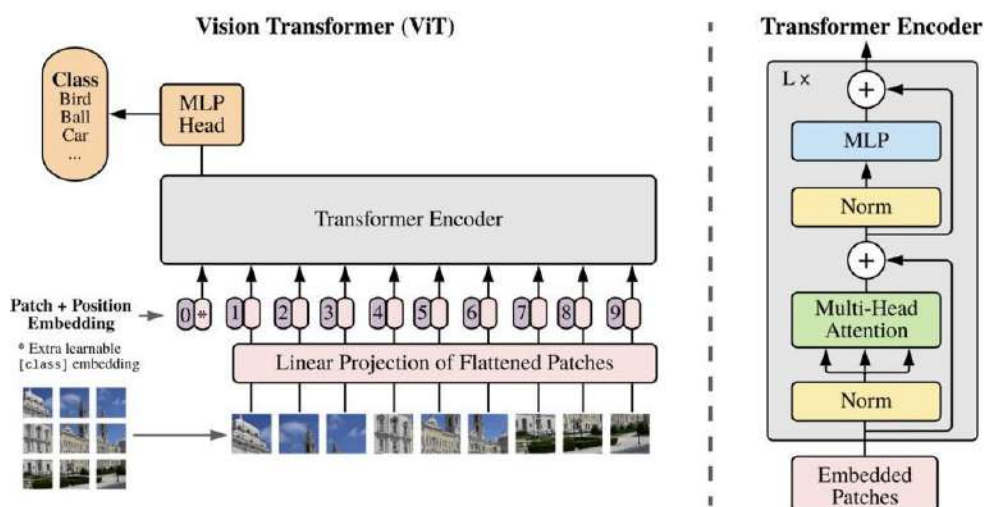
## Exercise: Transformer Applications

Ngày 4 tháng 12 năm 2023

### Phần 1. Transformer



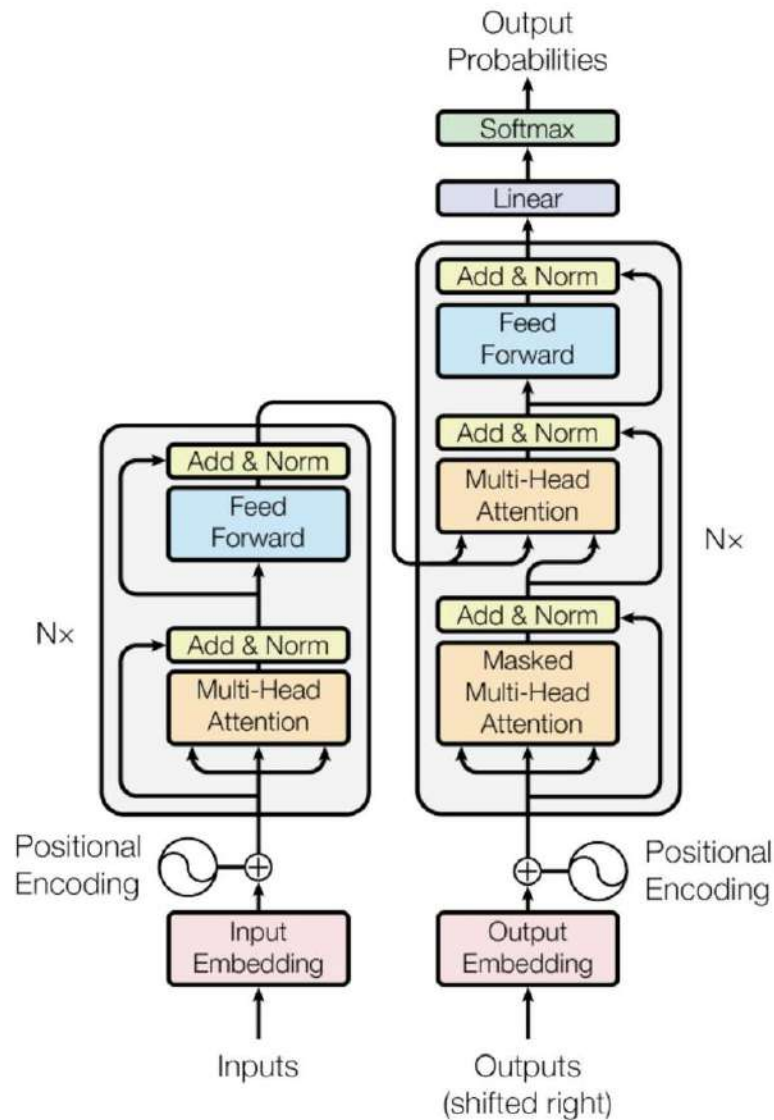
Hình 1: Mô hình Transformer-Encoder và BERT cho bài toán phân loại văn bản.



Hình 2: Mô hình Vision Transformer cho bài toán phân loại hình ảnh.

Mô hình Transformer ra đời với kỹ thuật cốt lõi và cơ chế Attention, đã đạt được những kết quả ấn tượng cho các bài toán về dữ liệu văn bản rồi từ đó mở rộng cho các kiểu dữ liệu như hình ảnh và âm thanh,... Trong phần này, chúng ta sẽ tìm hiểu các thành phần trong mô hình Transformer và ứng dụng cho bài toán phân loại văn bản.

## 1.1. Kiến trúc Transformer



Hình 3: Mô hình kiến trúc Transformer.

Kiến trúc Transformer gồm các thành phần:

- **Input Embedding:** Biểu diễn các token đầu vào (Thường được tách bởi Subword-based Tokenization) thành các Dense Vector.
- **Positional Encoding:** Biểu diễn vị trí (thứ tự) của các token trong câu. Thường được tính dựa vào hàm sinusoid hoặc được học trong quá trình huấn luyện mô hình.

- Các khối encoder: Để mã hoá các tokens đầu vào thành các contextual embedding. Bao gồm: Multi-Head Attention, Add - Normalization, Feed Forward
- Các khối decoder: nhận input là các token lịch sử và trạng thái mã hoá từ encoder, giải mã dự đoán token tiếp theo. Gồm: Masked Multi-Head Attention (dựa vào token lịch sử của decoder), Multi-Head Attention (dựa vào encoder và trạng thái hiện tại decoder), Add - Normalization, Feed Forward
- Language Model Head: Projection và Softmax dự đoán token tiếp theo vs xác suất lớn nhất.

## 1. Input Embedding, Positional Encoding

```

1 class TokenAndPositionEmbedding(nn.Module):
2     def __init__(self, vocab_size, embed_dim, max_length, device='cpu'):
3         super().__init__()
4         self.device = device
5         self.word_emb = nn.Embedding(
6             num_embeddings=vocab_size,
7             embedding_dim=embed_dim
8         )
9         self.pos_emb = nn.Embedding(
10            num_embeddings=max_length,
11            embedding_dim=embed_dim
12        )
13
14    def forward(self, x):
15        N, seq_len = x.size()
16        positions = torch.arange(0, seq_len).expand(N, seq_len).to(self.device)
17        output1 = self.word_emb(x)
18        output2 = self.pos_emb(positions)
19        output = output1 + output2
20        return output

```

## 2. Encoder

```

1 class TransformerEncoderBlock(nn.Module):
2     def __init__(self, embed_dim, num_heads, ff_dim, dropout=0.1):
3         super().__init__()
4         self.attn = nn.MultiheadAttention(
5             embed_dim=embed_dim,
6             num_heads=num_heads,
7             batch_first=True
8         )
9         self.ffn = nn.Sequential(
10            nn.Linear(in_features=embed_dim, out_features=ff_dim, bias=True),
11            nn.ReLU(),
12            nn.Linear(in_features=ff_dim, out_features=embed_dim, bias=True)
13        )
14        self.layernorm_1 = nn.LayerNorm(normalized_shape=embed_dim, eps=1e-6)
15        self.layernorm_2 = nn.LayerNorm(normalized_shape=embed_dim, eps=1e-6)
16        self.dropout_1 = nn.Dropout(p=dropout)
17        self.dropout_2 = nn.Dropout(p=dropout)
18
19    def forward(self, query, key, value):
20        attn_output, _ = self.attn(query, key, value)
21        attn_output = self.dropout_1(attn_output)
22        out_1 = self.layernorm_1(query + attn_output)
23        ffn_output = self.ffn(out_1)
24        ffn_output = self.dropout_2(ffn_output)
25        out_2 = self.layernorm_2(out_1 + ffn_output)
26        return out_2

```



```

27
28 class TransformerEncoder(nn.Module):
29     def __init__(self,
30                 src_vocab_size, embed_dim, max_length, num_layers, num_heads, ff_dim,
31                 dropout=0.1, device='cpu'
32             ):
33         super().__init__()
34         self.embedding = TokenAndPositionEmbedding(
35             src_vocab_size, embed_dim, max_length, device
36         )
37         self.layers = nn.ModuleList(
38             [
39                 TransformerEncoderBlock(
40                     embed_dim, num_heads, ff_dim, dropout
41                 ) for i in range(num_layers)
42             ]
43         )
44
45     def forward(self, x):
46         output = self.embedding(x)
47         for layer in self.layers:
48             output = layer(output, output, output)
49         return output

```

### 3. Decoder

```

1 class TransformerDecoderBlock(nn.Module):
2     def __init__(self, embed_dim, num_heads, ff_dim, dropout=0.1):
3         super().__init__()
4         self.attn = nn.MultiheadAttention(
5             embed_dim=embed_dim,
6             num_heads=num_heads,
7             batch_first=True
8         )
9         self.cross_attn = nn.MultiheadAttention(
10             embed_dim=embed_dim,
11             num_heads=num_heads,
12             batch_first=True
13         )
14         self.ffn = nn.Sequential(
15             nn.Linear(in_features=embed_dim, out_features=ff_dim, bias=True),
16             nn.ReLU(),
17             nn.Linear(in_features=ff_dim, out_features=embed_dim, bias=True)
18         )
19         self.layernorm_1 = nn.LayerNorm(normalized_shape=embed_dim, eps=1e-6)
20         self.layernorm_2 = nn.LayerNorm(normalized_shape=embed_dim, eps=1e-6)
21         self.layernorm_3 = nn.LayerNorm(normalized_shape=embed_dim, eps=1e-6)
22         self.dropout_1 = nn.Dropout(p=dropout)
23         self.dropout_2 = nn.Dropout(p=dropout)
24         self.dropout_3 = nn.Dropout(p=dropout)
25
26     def forward(self, x, enc_output, src_mask, tgt_mask):
27         attn_output, _ = self.attn(x, x, x, attn_mask=tgt_mask)
28         attn_output = self.dropout_1(attn_output)
29         out_1 = self.layernorm_1(x + attn_output)
30
31         attn_output, _ = self.cross_attn(
32             out_1, enc_output, enc_output, attn_mask=src_mask
33         )
34         attn_output = self.dropout_2(attn_output)
35         out_2 = self.layernorm_2(out_1 + attn_output)

```

```

36         ffn_output = self.ffn(out_2)
37         ffn_output = self.dropout_2(ffn_output)
38         out_3 = self.layer_norm_2(out_2 + ffn_output)
39         return out_3
40
41
42 class TransformerDecoder(nn.Module):
43     def __init__(self,
44                 tgt_vocab_size, embed_dim, max_length, num_layers, num_heads, ff_dim,
45                 dropout=0.1, device='cpu',
46                 ):
47         super().__init__()
48         self.embedding = TokenAndPositionEmbedding(
49             tgt_vocab_size, embed_dim, max_length, device
50         )
51         self.layers = nn.ModuleList(
52             [
53                 TransformerDecoderBlock(
54                     embed_dim, num_heads, ff_dim, dropout
55                 ) for i in range(num_layers)
56             ]
57         )
58
59     def forward(self, x, enc_output, src_mask, tgt_mask):
60         output = self.embedding(x)
61         for layer in self.layers:
62             output = layer(output, enc_output, src_mask, tgt_mask)
63         return output

```

#### 4. Transformer

```

1 class Transformer(nn.Module):
2     def __init__(self,
3                 src_vocab_size, tgt_vocab_size,
4                 embed_dim, max_length, num_layers, num_heads, ff_dim,
5                 dropout=0.1, device='cpu',
6                 ):
7         super().__init__()
8         self.device = device
9         self.encoder = TransformerEncoder(
10             src_vocab_size, embed_dim, max_length, num_layers, num_heads, ff_dim
11         )
12         self.decoder = TransformerDecoder(
13             tgt_vocab_size, embed_dim, max_length, num_layers, num_heads, ff_dim
14         )
15         self.fc = nn.Linear(embed_dim, tgt_vocab_size)
16
17     def generate_mask(self, src, tgt):
18         src_seq_len = src.shape[1]
19         tgt_seq_len = tgt.shape[1]
20
21         src_mask = torch.zeros(
22             (src_seq_len, src_seq_len),
23             device=self.device).type(torch.bool)
24
25         tgt_mask = (torch.triu(torch.ones(
26             (tgt_seq_len, tgt_seq_len),
27             device=self.device)
28         ) == 1).transpose(0, 1)
29         tgt_mask = tgt_mask.float().masked_fill(
30             tgt_mask == 0, float('-inf')).masked_fill(tgt_mask == 1, float(0.0))

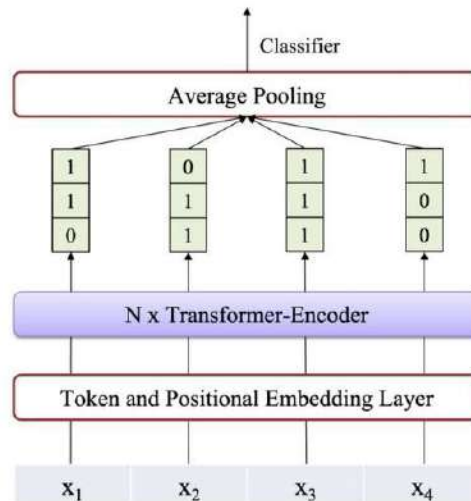
```

```
31         return src_mask, tgt_mask
32
33     def forward(self, src, tgt):
34         src_mask, tgt_mask = self.generate_mask(src, tgt)
35         enc_output = self.encoder(src)
36         dec_output = self.decoder(tgt, enc_output, src_mask, tgt_mask)
37         output = self.fc(dec_output)
38         return output
```

## 5. Thử nghiệm

```
1 batch_size = 128
2 src_vocab_size = 1000
3 tgt_vocab_size = 2000
4 embed_dim = 200
5 max_length = 100
6 num_layers = 2
7 num_heads = 4
8 ff_dim = 256
9
10 model = Transformer(
11     src_vocab_size, tgt_vocab_size,
12     embed_dim, max_length, num_layers, num_heads, ff_dim
13 )
14
15 src = torch.randint(
16     high=2,
17     size=(batch_size, max_length),
18     dtype=torch.int64
19 )
20
21 tgt = torch.randint(
22     high=2,
23     size=(batch_size, max_length),
24     dtype=torch.int64
25 )
26
27 prediction = model(src, tgt)
28 prediction.shape # batch_size x max_length x tgt_vocab_size
```

## 1.2. Text Classification



Hình 4: Phân loại văn bản sử dụng Transformer-Encoder

Ở phần này, chúng ta sử dụng mô hình Transformer-Encoder cho bài toán phân loại văn bản. Kiến trúc mô hình gồm các phần:

- Input Embedding và Positional Encoding
- Các lớp Encoder
- Lớp Average Pooling: lấy trung bình biểu diễn các từ thành biểu diễn cho câu
- Classifier: Linear layer

### 1. Load Dataset

```

1 # download
2 !git clone https://github.com/congngghia0609/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 # load data
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. Preprocessing

```

1 import re
2 import string
3
4 def preprocess_text(text):
5     # remove URLs https://www.
6     url_pattern = re.compile(r'https?://\s+\www\.\s+')
7     text = url_pattern.sub(" ", text)
8
9     # remove HTML Tags: <>
10    html_pattern = re.compile(r'<[<>]+>')
11    text = html_pattern.sub(" ", text)
12
13    # remove puncs and digits
14    replace_chars = list(string.punctuation + string.digits)
15    for char in replace_chars:
16        text = text.replace(char, " ")
17
18    # remove emoji
19    emoji_pattern = re.compile("[
20        u"\U0001F600-\U0001F64F" # emoticons
21        u"\U0001F300-\U0001F5FF" # symbols & pictographs
22        u"\U0001F680-\U0001F6FF" # transport & map symbols
23        u"\U0001F1E0-\U0001F1FF" # flags (iOS)
24        u"\U0001F1F2-\U0001F1F4" # Macau flag
25        u"\U0001F1E6-\U0001F1FF" # flags
26        u"\U0001F600-\U0001F64F"
27        u"\U00002702-\U000027B0"
28        u"\U000024C2-\U0001F251"
29        u"\U0001f926-\U0001f937"
30        u"\U0001F1F2"
31        u"\U0001F1F4"
32        u"\U0001F620"
33        u"\u200d"
34        u"\u2640-\u2642"
35    ]+", flags=re.UNICODE)
36    text = emoji_pattern.sub(" ", text)
37
38    # normalize whitespace
39    text = " ".join(text.split())
40
41    # lowercasing
42    text = text.lower()

```



```

43     return text
44
45 train_df['sentence'] = [preprocess_text(row['sentence']) for index, row in train_df.
    iterrows()]
46 valid_df['sentence'] = [preprocess_text(row['sentence']) for index, row in valid_df.
    iterrows()]
47 test_df['sentence'] = [preprocess_text(row['sentence']) for index, row in test_df.
    iterrows()]

```

### 3. Representation

```

1  # !pip install -q torchtext==0.16.0
2
3  def yield_tokens(sentences, tokenizer):
4      for sentence in sentences:
5          yield tokenizer(sentence)
6
7
8  # word-based tokenizer
9  from torchtext.data.utils import get_tokenizer
10 tokenizer = get_tokenizer("basic_english")
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['preprocess_sentence'], tokenizer),
18     max_tokens=vocab_size,
19     specials=["<pad>", "<unk>"]
20 )
21 vocabulary.set_default_index(vocabulary["<unk>"])
22
23 # convert torchtext dataset
24 from torchtext.data.functional import to_map_style_dataset
25
26 def prepare_dataset(df):
27     # create iterator for dataset: (sentence, label)
28     for index, row in df.iterrows():
29         sentence = row['preprocess_sentence']
30         encoded_sentence = vocabulary(tokenizer(sentence))
31         label = row['label']
32         yield encoded_sentence, label
33
34 train_dataset = prepare_dataset(train_df)
35 train_dataset = to_map_style_dataset(train_dataset)
36
37 valid_dataset = prepare_dataset(valid_df)
38 valid_dataset = to_map_style_dataset(valid_dataset)
39
40 test_dataset = prepare_dataset(test_df)
41 test_dataset = to_map_style_dataset(test_dataset)

```

### 4. Dataloader

```

1  import torch
2
3  seq_length = 100
4
5  def collate_batch(batch):
6      # create inputs, offsets, labels for batch
7      sentences, labels = list(zip(*batch))

```

```

8     encoded_sentences = [
9         sentence+([0]* (seq_length-len(sentence))) if len(sentence) < seq_length else
10        sentence[:seq_length]
11        for sentence in sentences
12    ]
13
14    encoded_sentences = torch.tensor(encoded_sentences, dtype=torch.int64)
15    labels = torch.tensor(labels)
16
17    return encoded_sentences, labels
18
19 from torch.utils.data import DataLoader
20
21 batch_size = 128
22
23 train_dataloader = DataLoader(
24     train_dataset,
25     batch_size=batch_size,
26     shuffle=True,
27     collate_fn=collate_batch
28 )
29
30 valid_dataloader = DataLoader(
31     valid_dataset,
32     batch_size=batch_size,
33     shuffle=False,
34     collate_fn=collate_batch
35 )
36
37 test_dataloader = DataLoader(
38     test_dataset,
39     batch_size=batch_size,
40     shuffle=False,
41     collate_fn=collate_batch
42 )

```

## 5. Trainer

```

1 # train epoch
2 import time
3
4 def train_epoch(model, optimizer, criterion, train_dataloader, device, epoch=0,
5                 log_interval=50):
6     model.train()
7     total_acc, total_count = 0, 0
8     losses = []
9     start_time = time.time()
10
11     for idx, (inputs, labels) in enumerate(train_dataloader):
12         inputs = inputs.to(device)
13         labels = labels.to(device)
14
15         optimizer.zero_grad()
16
17         predictions = model(inputs)
18
19         # compute loss
20         loss = criterion(predictions, labels)
21         losses.append(loss.item())
22
23         # backward
24         loss.backward()

```

```

24     optimizer.step()
25     total_acc += (predictions.argmax(1) == labels).sum().item()
26     total_count += labels.size(0)
27     if idx % log_interval == 0 and idx > 0:
28         elapsed = time.time() - start_time
29         print(
30             "| epoch {:3d} | {:5d}/{:5d} batches |
31             "| accuracy {:.3f}".format(
32                 epoch, idx, len(train_dataloader), total_acc / total_count
33             )
34         )
35         total_acc, total_count = 0, 0
36         start_time = time.time()
37
38     epoch_acc = total_acc / total_count
39     epoch_loss = sum(losses) / len(losses)
40     return epoch_acc, epoch_loss
41
42 # evaluate
43 def evaluate_epoch(model, criterion, valid_dataloader, device):
44     model.eval()
45     total_acc, total_count = 0, 0
46     losses = []
47
48     with torch.no_grad():
49         for idx, (inputs, labels) in enumerate(valid_dataloader):
50             inputs = inputs.to(device)
51             labels = labels.to(device)
52
53             predictions = model(inputs)
54
55             loss = criterion(predictions, labels)
56             losses.append(loss.item())
57
58             total_acc += (predictions.argmax(1) == labels).sum().item()
59             total_count += labels.size(0)
60
61     epoch_acc = total_acc / total_count
62     epoch_loss = sum(losses) / len(losses)
63     return epoch_acc, epoch_loss
64
65 # train
66 def train(model, model_name, save_model, optimizer, criterion, train_dataloader,
67           valid_dataloader, num_epochs, device):
68     train_accs, train_losses = [], []
69     eval_accs, eval_losses = [], []
70     best_loss_eval = 100
71     times = []
72     for epoch in range(1, num_epochs+1):
73         epoch_start_time = time.time()
74         # Training
75         train_acc, train_loss = train_epoch(model, optimizer, criterion,
76         train_dataloader, device, epoch)
77         train_accs.append(train_acc)
78         train_losses.append(train_loss)
79
80         # Evaluation
81         eval_acc, eval_loss = evaluate_epoch(model, criterion, valid_dataloader,
82         device)

```



```

81     eval_accs.append(eval_acc)
82     eval_losses.append(eval_loss)
83
84     # Save best model
85     if eval_loss < best_loss_eval:
86         torch.save(model.state_dict(), save_model + f'/{model_name}.pt')
87
88     times.append(time.time() - epoch_start_time)
89     # Print loss, acc end epoch
90     print("-" * 59)
91     print(
92         "| End of epoch {:3d} | Time: {:.5.2f}s | Train Accuracy {:.8.3f} | Train
Loss {:.8.3f} "
93         "| Valid Accuracy {:.8.3f} | Valid Loss {:.8.3f} ".format(
94             epoch, time.time() - epoch_start_time, train_acc, train_loss, eval_acc
, eval_loss
95         )
96     )
97     print("-" * 59)
98
99     # Load best model
100    model.load_state_dict(torch.load(save_model + f'/{model_name}.pt'))
101    model.eval()
102    metrics = {
103        'train_accuracy': train_accs,
104        'train_loss': train_losses,
105        'valid_accuracy': eval_accs,
106        'valid_loss': eval_losses,
107        'time': times
108    }
109    return model, metrics
110
111 # report
112 import matplotlib.pyplot as plt
113
114 def plot_result(num_epochs, train_accs, eval_accs, train_losses, eval_losses):
115     epochs = list(range(num_epochs))
116     fig, axs = plt.subplots(nrows = 1, ncols = 2 , figsize = (12,6))
117     axs[0].plot(epochs, train_accs, label = "Training")
118     axs[0].plot(epochs, eval_accs, label = "Evaluation")
119     axs[1].plot(epochs, train_losses, label = "Training")
120     axs[1].plot(epochs, eval_losses, label = "Evaluation")
121     axs[0].set_xlabel("Epochs")
122     axs[1].set_xlabel("Epochs")
123     axs[0].set_ylabel("Accuracy")
124     axs[1].set_ylabel("Loss")
125     plt.legend()

```

## 6. Modeling

```

1 class TransformerEncoderCls(nn.Module):
2     def __init__(self,
3         vocab_size, max_length, num_layers, embed_dim, num_heads, ff_dim,
4         dropout=0.1, device='cpu'
5     ):
6         super().__init__()
7         self.encoder = TransformerEncoder(
8             vocab_size, embed_dim, max_length, num_layers, num_heads, ff_dim, dropout,
9             device
10        )
11        self.pooling = nn.AvgPool1d(kernel_size=max_length)

```

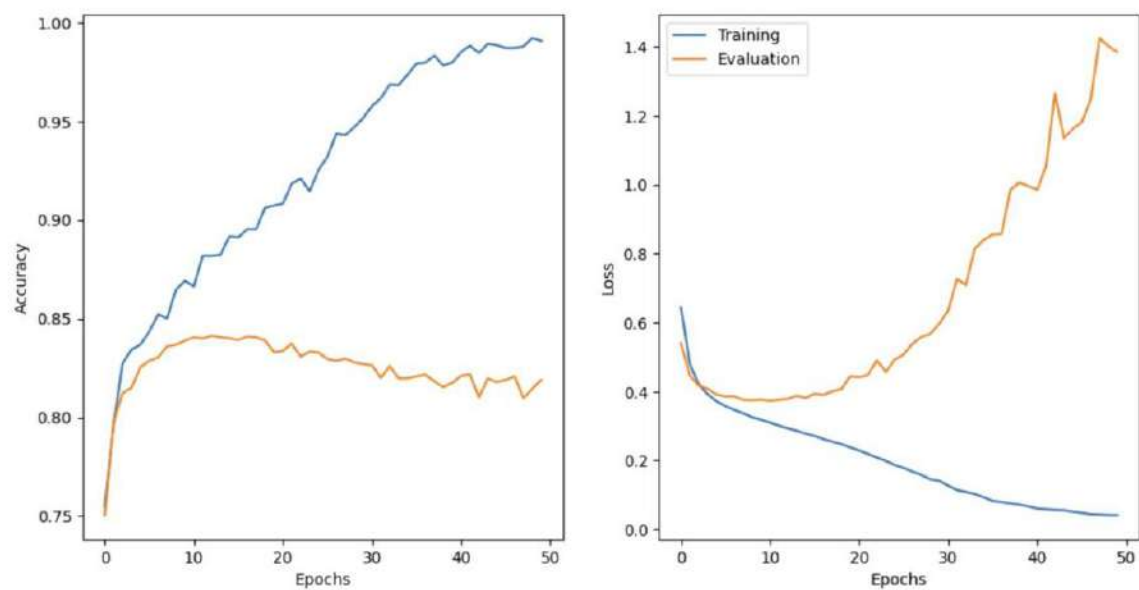


```
11     self.fc1 = nn.Linear(in_features=embed_dim, out_features=20)
12     self.fc2 = nn.Linear(in_features=20, out_features=2)
13     self.dropout = nn.Dropout(p=dropout)
14     self.relu = nn.ReLU()
15     def forward(self, x):
16         output = self.encoder(x)
17         output = self.pooling(output.permute(0,2,1)).squeeze()
18         output = self.dropout(output)
19         output = self.fc1(output)
20         output = self.dropout(output)
21         output = self.fc2(output)
22         return output
```

## 7. Training

```
1 import torch.optim as optim
2
3 vocab_size = 10000
4 max_length = 100
5 embed_dim = 200
6 num_layers = 2
7 num_heads = 4
8 ff_dim = 128
9 dropout=0.1
10
11 model = TransformerEncoderCls(
12     vocab_size, max_length, num_layers, embed_dim, num_heads, ff_dim, dropout
13 )
14
15 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
16
17 model = TransformerEncoderCls(
18     vocab_size, max_length, num_layers, embed_dim, num_heads, ff_dim, dropout, device
19 )
20 model.to(device)
21
22 criterion = torch.nn.CrossEntropyLoss()
23 optimizer = optim.Adam(model.parameters(), lr=0.00005)
24
25 num_epochs = 50
26 save_model = './model'
27 os.makedirs(save_model, exist_ok = True)
28 model_name = 'model'
29
30 model, metrics = train(
31     model, model_name, save_model, optimizer, criterion, train_dataloader,
32     valid_dataloader, num_epochs, device
33 )
```

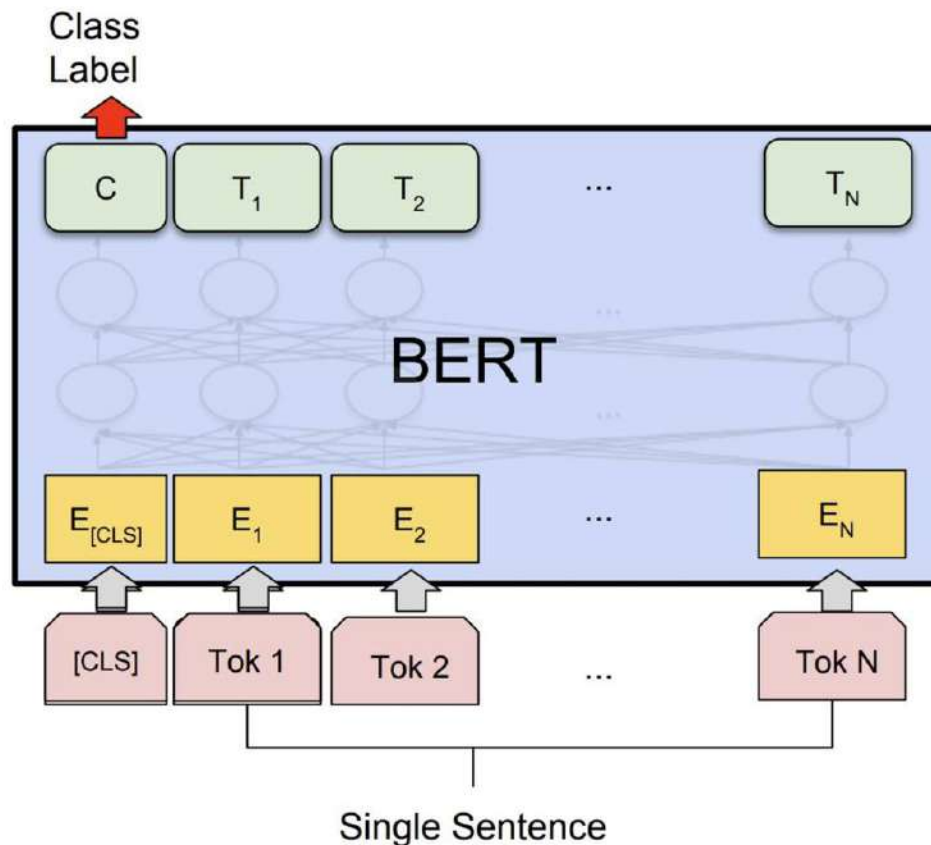
Kết quả training và accuracy trên tập test là 82.55



Hình 5: Quá trình huấn luyện bộ dữ liệu NTC-SCV với mô hình Transformer-Encoder.

## Phần 2. Text Classification using BERT

Một trong những mô hình pretrained đầu tiên cho dữ liệu văn bản dựa vào kiến trúc mô hình Transformer được ứng dụng cho các downstream task khác nhau đó là BERT. Trong phần này chúng ta sẽ fine tuning BERT cho bài toán phân loại trên bộ dữ liệu NTC-SCV dựa vào thư viện transformers của huggingface.



Hình 6: Fine-Tuning BERT cho bài toán phân loại văn bản.

### 1. Load Dataset

```

1 # install libs
2 !pip install -q -U transformers datasets accelerate evaluate
3
4 # download
5 !git clone https://github.com/congnghia0609/ntc-scv.git
6 !unzip ./ntc-scv/data/data_test.zip -d ./data
7 !unzip ./ntc-scv/data/data_train.zip -d ./data
8 !rm -rf ./ntc-scv
9
10 # load data
11 import os
12 import pandas as pd
13
14 def load_data_from_path(folder_path):
15     examples = []

```

```

16 for label in os.listdir(folder_path):
17     full_path = os.path.join(folder_path, label)
18     for file_name in os.listdir(full_path):
19         file_path = os.path.join(full_path, file_name)
20         with open(file_path, "r", encoding="utf-8") as f:
21             lines = f.readlines()
22             sentence = " ".join(lines)
23             if label == "neg":
24                 label = 0
25             if label == "pos":
26                 label = 1
27             data = {
28                 'sentence': sentence,
29                 'label': label
30             }
31             examples.append(data)
32 return pd.DataFrame(examples)
33
34 folder_paths = {
35     'train': './data/data_train/train',
36     'valid': './data/data_train/test',
37     'test': './data/data_test/test'
38 }
39
40 train_df = load_data_from_path(folder_paths['train'])
41 valid_df = load_data_from_path(folder_paths['valid'])
42 test_df = load_data_from_path(folder_paths['test'])
43
44 # convert to Dataset Object
45 from datasets import Dataset, DatasetDict
46
47 raw_dataset = DatasetDict({
48     'train': Dataset.from_pandas(train_df),
49     'valid': Dataset.from_pandas(valid_df),
50     'test': Dataset.from_pandas(test_df)
51 })

```

## 2. Preprocessing

```

1 import re
2 import string
3
4 def preprocess_text(text):
5     # remove URLs https://www.
6     url_pattern = re.compile(r'https?://\s+\www\.\s+')
7     text = url_pattern.sub(" ", text)
8
9     # remove HTML Tags: <>
10    html_pattern = re.compile(r'<[<>]+>')
11    text = html_pattern.sub(" ", text)
12
13    # remove puncs and digits
14    replace_chars = list(string.punctuation + string.digits)
15    for char in replace_chars:
16        text = text.replace(char, " ")
17
18    # remove emoji
19    emoji_pattern = re.compile("[
20        u"\U0001F600-\U0001F64F" # emoticons
21        u"\U0001F300-\U0001F5FF" # symbols & pictographs

```



```

22         u"\U0001F680-\U0001F6FF" # transport & map symbols
23         u"\U0001F1E0-\U0001F1FF" # flags (iOS)
24         u"\U0001F1F2-\U0001F1F4" # Macau flag
25         u"\U0001F1E6-\U0001F1FF" # flags
26         u"\U0001F600-\U0001F64F"
27         u"\U00002702-\U000027B0"
28         u"\U000024C2-\U0001F251"
29         u"\U0001f926-\U0001f937"
30         u"\U0001F1F2"
31         u"\U0001F1F4"
32         u"\U0001F620"
33         u"\u200d"
34         u"\u2640-\u2642"
35         "]" + ", flags=re.UNICODE)
36     text = emoji_pattern.sub(r" ", text)
37
38     # normalize whitespace
39     text = " ".join(text.split())
40
41     # lowercasing
42     text = text.lower()
43     return text
44
45 train_df['sentence'] = [preprocess_text(row['sentence']) for index, row in train_df.
46                          iterrows()]
47 valid_df['sentence'] = [preprocess_text(row['sentence']) for index, row in valid_df.
48                          iterrows()]
49 test_df['sentence'] = [preprocess_text(row['sentence']) for index, row in test_df.
50                         iterrows()]
51
52 # tokenization
53 from transformers import AutoTokenizer
54
55 model_name = "bert-base-uncased"
56
57 tokenizer = AutoTokenizer.from_pretrained(
58     model_name,
59     use_fast=True
60 )
61 max_seq_length = 100
62 max_seq_length = min(max_seq_length, tokenizer.model_max_length)
63
64 def preprocess_function(examples):
65     # Tokenize the texts
66
67     result = tokenizer(
68         examples["sentence"],
69         padding="max_length",
70         max_length=max_seq_length,
71         truncation=True
72     )
73     result["label"] = examples['label']
74
75     return result
76
77 # Running the preprocessing pipeline on all the datasets
78 processed_dataset = raw_dataset.map(
79     preprocess_function,
80     batched=True,
81     desc="Running tokenizer on dataset",
82 )

```

```
32
33 # collator with padding max length
34 from transformers import DataCollatorWithPadding
35
36 data_collator = DataCollatorWithPadding(tokenizer=tokenizer)
```

### 3. Modeling

```
1 from transformers import AutoConfig, AutoModelForSequenceClassification
2
3 num_labels = 2
4
5 config = AutoConfig.from_pretrained(
6     model_name,
7     num_labels=num_labels,
8     finetuning_task="text-classification"
9 )
10 model = AutoModelForSequenceClassification.from_pretrained(
11     model_name,
12     config=config
13 )
```

### 4. Metric

```
1 import numpy as np
2 import evaluate
3
4 metric = evaluate.load("accuracy")
5 def compute_metrics(eval_pred):
6     predictions, labels = eval_pred
7     predictions = np.argmax(predictions, axis=1)
8     result = metric.compute(predictions=predictions, references=labels)
9     return result
```

### 5. Trainer

```
1 from transformers import TrainingArguments, Trainer
2
3 training_args = TrainingArguments(
4     output_dir="save_model",
5     learning_rate=2e-5,
6     per_device_train_batch_size=128,
7     per_device_eval_batch_size=128,
8     num_train_epochs=10,
9     evaluation_strategy="epoch",
10    save_strategy="epoch",
11    load_best_model_at_end=True
12 )
13
14 trainer = Trainer(
15     model=model,
16     args=training_args,
17     train_dataset=processed_dataset["train"],
18     eval_dataset=processed_dataset["valid"],
19     compute_metrics=compute_metrics,
20     tokenizer=tokenizer,
21     data_collator=data_collator,
22 )
23
```

```
24 trainer.train()
```

## 6. Training

Kết quả training và accuracy trên tập test là 85.52

Epoch	Training Loss	Validation Loss	Accuracy
1	No log	0.423262	0.809400
2	No log	0.401464	0.824700
3	0.447100	0.392191	0.836700
4	0.447100	0.373473	0.842200
5	0.316700	0.396957	0.840800
6	0.316700	0.416141	0.842200
7	0.245100	0.454328	0.843700
8	0.245100	0.468748	0.845100
9	0.187800	0.486744	0.841400
10	0.187800	0.489360	0.842600

Hình 7: Quá trình huấn luyện bộ dữ liệu NTC-SCV dựa vào fine tuning BERT.

# Phần 3. Vision Transformer

## 1. Load Dataset

```
1 import torch
2 import torchvision.transforms as transforms
3 from torch.utils.data import DataLoader, random_split
4 import torch.optim as optim
5 from torchvision.datasets import ImageFolder
6 from torch import nn
7 import math
8 import os

1 # download
2 !gdown 11Buzytn4vIh4x_0qz8MY29JMMdIqSzj-
3 !unzip ./flower_photos.zip
4
5 # load data
6 data_patch = "./flower_photos"
7 dataset = ImageFolder(root=data_patch)
8 num_samples = len(dataset)
9 classes = dataset.classes
10 num_classes = len(dataset.classes)
11
12 # split
13 TRAIN_RATIO, VALID_RATIO = 0.8, 0.1
14 n_train_examples = int(num_samples * TRAIN_RATIO)
15 n_valid_examples = int(num_samples * VALID_RATIO)
16 n_test_examples = num_samples - n_train_examples - n_valid_examples
17 train_dataset, valid_dataset, test_dataset = random_split(
18     dataset,
19     [n_train_examples, n_valid_examples, n_test_examples]
20 )
```

## 2. Preprocessing

```
1 # resize + convert to tensor
2 IMG_SIZE = 224
3
4 train_transforms = transforms.Compose([
5     transforms.Resize((IMG_SIZE, IMG_SIZE)),
6     transforms.RandomHorizontalFlip(),
7     transforms.RandomRotation(0.2),
8     transforms.ToTensor(),
9     transforms.Normalize([0.5, 0.5, 0.5], [0.5, 0.5, 0.5])
10 ])
11
12 test_transforms = transforms.Compose([
13     transforms.Resize((IMG_SIZE, IMG_SIZE)),
14     transforms.ToTensor(),
15     transforms.Normalize([0.5, 0.5, 0.5], [0.5, 0.5, 0.5])
16 ])
17
18 # apply
19 train_dataset.dataset.transform = train_transforms
20 valid_dataset.dataset.transform = test_transforms
21 test_dataset.dataset.transform = test_transforms
```