

# Fine Tuning do BERT no IMDB

Nome: Elton Cardoso do Nascimento

## Instruções:

Treinar e medir a acurácia de um modelo BERT (ou variantes) para classificação binária usando o dataset do IMDB (20k/5k amostras de treino/validação).

Importante:

- [x] Deve-se implementar o próprio laço de treinamento.
- [x] Implementar o acumulo de gradiente.

Dicas:

- BERT geralmente costuma aprender bem uma tarefa com poucas épocas (de 3 a 5 épocas). Se tiver demorando mais de 5 épocas para chegar em 80% de acurácia, ajuste os hiperparametros.
- Solução para erro de memória:
  - [x] Usar bfloat16 permite quase dobrar o batch size

Opcional:

- Pode-se usar a função trainer da biblioteca Transformers/HuggingFace para verificar se seu laço de treinamento está correto. Note que ainda assim é obrigatório implementar o laço próprio.

```
In [ ]: import os # Manipular arquivos
import random # Operações randômicas
import pickle # Serializar/deserializar backups
import time # Medição de tempo
import string # Operações com strings
from concurrent.futures import ThreadPoolExecutor # Parelização
from typing import Tuple, List, Dict, Optional # Type hints

import numpy as np # Operações vetoriais
import matplotlib.pyplot as plt # Plots
import tqdm
import torch # ML
from torch.utils.data import Dataset, DataLoader # Preparação de dados

try:
    import wandb # Logging
except:
    wandb = None
```

## Fixando a seed

```
In [ ]: def reset_seeds(seed:int=123):
        random.seed(seed)
        np.random.seed(seed)
        torch.manual_seed(seed)
```

```
In [ ]: reset_seeds()
```

## Preparando Dados

Primeiro, fazemos download do dataset:

```
In [ ]: if not os.path.isfile("aclImdb.tgz"):
        !curl -LO http://files.fast.ai/data/aclImdb.tgz
        !tar -xzf aclImdb.tgz
```

## Carregando o dataset

Criaremos uma divisão de treino (20k exemplos) e validação (5k exemplos) artificialmente.

```
In [ ]: max_valid = 5000
```

```
In [ ]: def load_texts(folder):
        texts = []
```

```

for path in os.listdir(folder):
    with open(os.path.join(folder, path), encoding="utf8") as f:
        texts.append(f.read())
return texts

```

```

In [ ]:
executor = ThreadPoolExecutor(max_workers=4)

folders = ['aclImdb/train/pos', 'aclImdb/train/neg', 'aclImdb/test/pos', 'aclImdb/test/neg']

futures = []
for folder in folders:
    future = executor.submit(load_texts, folder)

    futures.append(future)

all_texts = []

for future in futures:
    texts = future.result()

    all_texts.append(texts)

executor.shutdown()

x_train_pos = all_texts[0]
x_train_neg = all_texts[1]
x_test_pos = all_texts[2]
x_test_neg = all_texts[3]

```

```

In [ ]:
x_train = x_train_pos + x_train_neg
x_test = x_test_pos + x_test_neg
y_train = [True] * len(x_train_pos) + [False] * len(x_train_neg)
y_test = [True] * len(x_test_pos) + [False] * len(x_test_neg)

```

```

In [ ]:
# Embaralhamos o treino para depois fazermos a divisão treino/valid.
c = list(zip(x_train, y_train))
random.shuffle(c)
x_train, y_train = zip(*c)

x_valid = x_train[-max_valid:]
y_valid = y_train[-max_valid:]
x_train = x_train[:-max_valid]
y_train = y_train[:-max_valid]

```

```

In [ ]:
print(len(x_train), 'amostras de treino.')
print(len(x_valid), 'amostras de desenvolvimento.')
print(len(x_test), 'amostras de teste.')

```

20000 amostras de treino.  
5000 amostras de desenvolvimento.  
25000 amostras de teste.

```

In [ ]:
print('3 primeiras amostras treino:')
for x, y in zip(x_train[:3], y_train[:3]):
    print(y, x[:100])

print('3 últimas amostras treino:')
for x, y in zip(x_train[-3:], y_train[-3:]):
    print(y, x[:100])

print('3 primeiras amostras validação:')
for x, y in zip(x_valid[:3], y_test[:3]):
    print(y, x[:100])

print('3 últimas amostras validação:')
for x, y in zip(x_valid[-3:], y_valid[-3:]):
    print(y, x[:100])

```

3 primeiras amostras treino:  
False POSSIBLE SPOILERS<br /><br />The Spy Who Shagged Me is a muchly overrated and over-hyped sequel. Int  
False The long list of "big" names in this flick (including the ubiquitous John Mills) didn't bowl me over  
True Bette Midler showcases her talents and beauty in "Diva Las Vegas". I am thrilled that I taped it and  
3 últimas amostras treino:  
False I was previously unaware that in the early 1990's Devry University (or was it ITT Tech?) added Film  
True The story and music (George Gershwin!) are wonderful, as are Levant, Guetary, Foch, and, of course,  
True This is my favorite show. I think it is utterly brilliant. Thanks to David Chase for bringing this i  
3 primeiras amostras validação:  
True Why has this not been released? I kind of thought it must be a bit rubbish since it hasn't been. How  
True I was amazingly impressed by this movie. It contained fundamental elements of depression, grief, lon  
True photography was too jumpy to follow. dark scenes hard to see.<br /><br />Had good story line too bad  
3 últimas amostras validação:  
True In the early to mid 1970's, Clifford Irving proposed to write the ultimate biography of Howard Hughe  
True An ultra-modern house in an affluent neighborhood appears to be the cause of each of its inhabitants  
True Some of the best movies that are categorized as "comedies" actually blur between comedy and drama. "

```

In [ ]:
GOOD_MOVIE = 1 #True
BAD_MOVIE = 0 #False

```

## Tokenizador

Preparamos o tokenizador para uso. No caso vamos utilizar o tokenizador preparado para o modelo BERT (é o mesmo do DistilBERT):

```
In [ ]: tokenizer = torch.hub.load('huggingface/pytorch-transformers', 'tokenizer', 'bert-base-cased')
```

Using cache found in C:\Users\eltsu/.cache/torch/hub/huggingface\_pytorch-transformers\_main

Podemos testar o tokenizador imprimindo uma sequência (observe o token inicial "101"=\ e final "102"=\\):

```
In [ ]: tokens = tokenizer(x_train[0], add_special_tokens=True, padding="max_length", max_length=512)

tokens["input_ids"][:10], tokens["input_ids"][-10:]
```

```
Out[ ]: ([101, 153, 9025, 13882, 13360, 2036, 16625, 2346, 17656, 9637],
 [156, 2328, 12165, 2508, 1110, 1141, 10010, 10866, 119, 102])
```

```
In [ ]: BERT_SEP = 102
```

## Dataset e Dataloader

Definimos o dataset para realizar a tokenizador e manipular os dados:

```
In [ ]: class IMDB_Dataset(Dataset):
    """
    Dataset for sentiment analisys

    Input: tokenized review and mask (for padding).
    Output: if is a good (1) or bad (0) review.
    """
    def __init__(self, x_data:List[str], y_data:List[bool], tokenizer) -> None:
        """
        Creates a new dataset.

        Args:
            x_data (List[str]): dataset reviews.
            y_data (List[bool]): dataset targets.
            tokenizer: tokenizer to encode reviews.

        """
        super().__init__()

        self._x_data = tokenizer(x_data,
                                return_tensors="pt", #Return as torch tensor
                                padding=True, #Add padding to small sequences
                                return_token_type_ids=False, #Don't return sequence mask (only one sequence)
                                truncation=True) #Truncate big sentences (max = 512 tokens, with CLS and SEP)

        self._y_data = torch.tensor(y_data, dtype=torch.float32)

        self._size = len(self._y_data)

        self.max_len = 512

    def __len__(self) -> int:
        """
        Gets the size of the dataset.

        Returns:
            int: dataset size.
        """
        return self._size

    def __getitem__(self, idx:int) -> Tuple[torch.Tensor, torch.Tensor, torch.Tensor]:
        """
        Gets a item of the dataset.

        Args:
            idx (int): data index.

        Returns:
            torch.Tensor: dataset input.
            torch.Tensor: dataset attention mask.
            torch.Tensor: dataset target.
        """
        x_data = self._x_data["input_ids"][idx]
        mask = self._x_data["attention_mask"][idx]

        if self.max_len != 512:
            x_data = x_data[:self.max_len]
            x_data[-1] = BERT_SEP

            mask = mask[:self.max_len]
```

```
return x_data, mask, self._y_data[idx]
```

```
In [ ]: datasets = {}

xs = [x_train, x_valid, x_test]
ys = [y_train, y_valid, y_test]
names = ["train", "val", "test"]

for i in range(3):
    dataset = IMDB_Dataset(xs[i], ys[i], tokenizer)
    datasets[names[i]] = dataset
```

Para evitar precisamos realizar várias vezes a tokenização durante o desenvolvimento, podemos serializar o dataset para posteriormente deserializá-lo:

```
In [ ]: file_name = "datasets.bin"
with open(file_name, "wb") as file:
    pickle.dump(datasets, file)
```

```
In [ ]: file_name = "datasets.bin"

with open(file_name, "rb") as file:
    datasets = pickle.load(file)
```

E definimos uma função para criar os dataloaders a partir dos datasets e batch size:

```
In [ ]: def create_dataloaders(datasets:Dict[str, Dataset], batch_size:int) -> Dict[str, DataLoader]:
    ...
    Generate dataloaders from datasets.

    Args:
        datasets (Dict[str, Dataset]): named datasets.
        batch_size (int): batch sizes.

    Returns:
        Dict[str, DataLoader]: dataloaders for the datasets.
    ...

    dataloaders = {}

    for name in datasets:
        dataloaders[name] = DataLoader(datasets[name], batch_size=batch_size, shuffle=True)

    return dataloaders
```

E uma função para alterar o tamanho das sequências que serão geradas:

```
In [ ]: def set_max_len(datasets:Dict[str, IMDB_Dataset], max_len):
    for name in datasets:
        datasets[name].max_len = max_len
```

## Preparação do modelo

Preparamos o modelo a ser utilizado, que é um modelo BERT com uma camada adicional para realizar a classificação, que recebe como entrada o embedding final relacionado ao token CLS:

```
In [ ]: class BinaryClassifierBERT(torch.nn.Module):
    ...
    Classifier model using BERT.
    ...

    def __init__(self, dropout_rate:float=0) -> None:
        ...
        Model constructor.

        Args:
            dropout_rate (float, optional): Dropout before the final layer. Defaults to 0.
        ...
        super().__init__()

        self.bert = torch.hub.load('huggingface/pytorch-transformers', 'model', 'distilbert-base-cased')

        self.dropout = torch.nn.Dropout(dropout_rate)
        self.linear = torch.nn.Linear(768, 1)

    def forward(self, input_ids:torch.Tensor, attention_masks:Optional[torch.Tensor]=None) -> torch.Tensor:
        ...
        Computes the classification for the input.

        Args:
            input_ids (torch.Tensor): tokenized input.
            attention_masks (torch.Tensor, optional): attention mask of the input. Defaults to None.
```

```

Returns:
    torch.Tensor: inference result.
...

bert_output = self.bert(input_ids=input_ids, attention_mask=attention_masks)
c_vector = bert_output.last_hidden_state[:, 0]

y = self.dropout(c_vector)
y = self.linear(y)

return y

```

## Treino

Nesta seção iremos realizar o treino, iniciando pela definição de algumas funções auxiliares.

### Funções auxiliares

Iremos definir três funções auxiliares: uma para calcular a perplexidade a partir da loss, outra para printar informações e uma final para calcular a loss:

```

In [ ]: def ppl(loss:torch.Tensor) -> torch.Tensor:
        """
        Computes the perplexity from the loss.

        Args:
            loss (torch.Tensor): loss to compute the perplexity.

        Returns:
            torch.Tensor: corresponding perplexity.
        """
        return torch.exp(loss)

```

```

In [ ]: def print_info(loss_value:torch.Tensor, epoch:int, total_epochs:int,
                       time:float=0.0, accuracy:Optional[float]=None):
        """
        Prints the information of a epoch.

        Args:
            loss_value (torch.Tensor): epoch loss.
            epoch (int): epoch number.
            total_epochs (int): total number of epochs.
            time (float, optional): time to run the epoch. Don't print if is 0.0. Defaults to 0.0.
            accuracy (float, optional): epoch accuracy.
        """
        ppl_value = ppl(loss_value)

        print(f'Epoch [{epoch+1}/{total_epochs}], \
              Loss: {loss_value.item():.4f}, \
              Perplexity: {ppl_value.item():.4f}', end="")

        if accuracy is not None:
            print(f', Accuracy: {100*accuracy:.4f}%')

        if time != 0:
            print(f", Elapsed Time: {time:.2f} sec")
        else:
            print("")

```

```

In [ ]: MODE_TRAIN = 0
        MODE_EVALUATE = 1

```

```

In [ ]: def compute_loss(model:torch.nn.Module, loader:DataLoader,
                        criterion:torch.nn.Module, mode:int = MODE_EVALUATE,
                        optimizer:Optional[torch.optim.Optimizer]=None) -> Tuple[torch.Tensor, torch.Tensor]:
        """
        Computes the loss from a model across a dataset.

        If in train mode also runs optimizer steps.

        Args:
            model (torch.nn.Module): model to evaluate.
            loader (DataLoader): dataset.
            criterion (torch.nn.Module): loss function to compute.
            mode (int): mode of the computation.
                        If MODE_EVALUATE, computes without gradient, in eval mode and detachs loss.
                        If MODE_TRAIN, computes with gradient and in train mode.
                        Default is MODE_EVALUATE.
            optimizer (torch.optim.Optimizer, optional): optimizer to use in the train mode.

        Returns:
            torch.Tensor: resulting loss.
            torch.Tensor: resulting accuracy
        """
        device = next(iter(model.parameters())).device

```

```

if mode == MODE_EVALUATE:
    model.eval()
    torch.set_grad_enabled(False)
elif mode == MODE_TRAIN:
    model.train()
    torch.set_grad_enabled(True)
else:
    raise ValueError(f"Unknown mode: {mode}.")

total_loss = torch.tensor(0, dtype=torch.float32, device=device)
correct = torch.tensor(0, dtype=torch.float32, device=device)
n = 0
for inputs, masks, targets in tqdm.tqdm(loader):
    inputs = inputs.to(device)
    masks = masks.to(device)

    targets = targets.reshape(-1)
    targets = targets.to(device)

    logits = model(inputs, masks)
    logits = logits.view(-1, logits.shape[-1])

    loss = criterion(logits.squeeze(), targets)
    total_loss += loss*targets.size(0)

    predicted = torch.round(torch.sigmoid(logits.squeeze()))
    correct += (predicted == targets).sum().item()

    n += targets.size(0)

    if mode == MODE_TRAIN:
        optimizer.zero_grad()
        loss.backward()
        optimizer.step()

total_loss /= n
accuracy = correct / n

torch.set_grad_enabled(True)

accuracy = accuracy.detach()
total_loss = total_loss.detach()

return total_loss, accuracy

```

## Inicialização

Começamos o processo de treino inicializando as variáveis.

Definimos se será realizado o logging utilizando o wandb:

```
In [ ]: use_wandb = True
```

Checamos se existe uma GPU disponível:

```
In [ ]: # Verifica se há uma GPU disponível e define o dispositivo para GPU se possível, caso contrário, usa a CPU
device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
device
```

```
Out[ ]: device(type='cuda')
```

Definimos os parâmetros de treino:

```
In [ ]: batch_size = 32
dropout_rate = 0.1
lr = 2.5e-5
n_epoch = 5
optimizer_class = torch.optim.Adam
seq_len = 256
weight_decay = 0

config = {
    "batch_size": batch_size,
    "dropout_rate": dropout_rate,
    "lr": lr,
    "n_epoch": n_epoch,
    "optimizer_class": optimizer_class.__name__,
    "seq_len" : seq_len,
    "weight_decay": weight_decay,
}

if use_wandb:
    run = wandb.init(project="IA024-04-BERT_IMDB", config=config)
    run_name = run.name
    run_id = run.id
    model_name = f"{run_name}-{run_id}.bin"
else:
    #get random name using config
    seed = 0

```

```
for name in config:
    if isinstance(config[name], int) or isinstance(config[name], float):
        seed += config[name]
    else:
        for c in config[name]:
            seed += ord(c)

reset_seeds(seed)
model_name = ''.join(random.choice(string.ascii_uppercase + string.digits) for _ in range(5))
```

Failed to detect the name of this notebook, you can set it manually with the WANDB\_NOTEBOOK\_NAME environment variable to enable code saving.

wandb: Currently logged in as: eltoncn. Use `wandb login --relogin` to force relogin

wandb version 0.16.6 is available! To upgrade, please run: \$ pip install wandb --upgrade

Tracking run with wandb version 0.15.8

Run data is saved locally in d:\Github\IA024\05-BERT\wandb\run-20240408\_205318-dc0slg44

Syncing run **worldly-darkness-8** to [Weights & Biases \(docs\)](#)

View project at [https://wandb.ai/eltoncn/IA024-04-BERT\\_IMDB](https://wandb.ai/eltoncn/IA024-04-BERT_IMDB)

View run at [https://wandb.ai/eltoncn/IA024-04-BERT\\_IMDB/runs/dc0slg44](https://wandb.ai/eltoncn/IA024-04-BERT_IMDB/runs/dc0slg44)

Reiniciamos as sementes:

```
In [ ]: reset_seeds()
```

Criamos o modelo, loss, otimizador e dataloaders:

```
In [ ]: model = BinaryClassifierBERT(dropout_rate)
model.to(device, dtype=torch.bfloat16)

criterion = torch.nn.BCEWithLogitsLoss()
optimizer = optimizer_class(model.parameters(), lr=lr, weight_decay=weight_decay)

set_max_len(datasets, seq_len)
dataloaders = create_dataloaders(datasets, batch_size)
```

Using cache found in C:\Users\eltsu\.cache\torch\hub\huggingface\_pytorch-transformers\_main

### Treino

E finalmente podemos realizar o processo de treino em si:

```
In [ ]: #Informações antes da primeira epoch
prev_loss, prev_accuracy = compute_loss(model, dataloaders["val"], criterion, MODE_EVALUATE)
print_info(prev_loss, -1, n_epoch, 0, prev_accuracy)
```



```
In [ ]: hist = {}
hist["loss_train"] = []
hist["loss_val"] = []
hist["ppl_train"] = []
hist["ppl_val"] = []
hist["accuracy_train"] = []
hist["accuracy_val"] = []

for epoch in range(n_epoch):
    start_time = time.time()

    loss_train, accuracy_train = compute_loss(model, dataloaders["train"], criterion, MODE_TRAIN, optimizer)

    end_time = time.time()

    epoch_duration = end_time - start_time

    ppl_train = ppl(loss_train)

    print_info(loss_train, epoch, n_epoch, epoch_duration, accuracy_train)

    #Validation stats
    loss_val, accuracy_val = compute_loss(model, dataloaders["val"], criterion, MODE_EVALUATE)
    ppl_val = ppl(loss_val)

    print("VAL ", end="")
    print_info(loss_val, epoch, n_epoch, accuracy=accuracy_val)

    #Save history
    hist["loss_train"].append(loss_train.item())
    hist["loss_val"].append(loss_val.item())
    hist["ppl_train"].append(ppl_train.item())
    hist["ppl_val"].append(ppl_val.item())
    hist["accuracy_train"].append(accuracy_train.item())
    hist["accuracy_val"].append(accuracy_val.item())

    log = {
```



```
        "loss_train": loss_train.item(),
        "loss_val": loss_val.item(),
        "ppl_train": ppl_train.item(),
        "ppl_val": ppl_val.item(),
        "accuracy_train": accuracy_train.item(),
        "accuracy_val": accuracy_val.item()
    }

    if use_wandb:
        wandb.log(log)

for key in hist:
    hist[key] = np.array(hist[key])

if use_wandb:
    wandb.finish()
```

```
0%|          | 0/625 [00:00<?, ?it/s]100%|          | 625/625 [08:10<00:00, 1.28it/s]
Epoch [1/5],          Loss: 0.3762,          Perplexity: 1.4568, Accuracy: 82.5050%
, Elapsed Time: 490.46 sec
100%|          | 157/157 [00:41<00:00, 3.78it/s]
VAL Epoch [1/5],          Loss: 0.2958,          Perplexity: 1.3443, Accuracy: 87.7400%

100%|          | 625/625 [08:03<00:00, 1.29it/s]
Epoch [2/5],          Loss: 0.2708,          Perplexity: 1.3109, Accuracy: 88.9000%
, Elapsed Time: 483.77 sec
100%|          | 157/157 [00:41<00:00, 3.80it/s]
VAL Epoch [2/5],          Loss: 0.2811,          Perplexity: 1.3246, Accuracy: 88.4000%

100%|          | 625/625 [07:57<00:00, 1.31it/s]
Epoch [3/5],          Loss: 0.2473,          Perplexity: 1.2805, Accuracy: 89.9900%
, Elapsed Time: 477.57 sec
100%|          | 157/157 [00:41<00:00, 3.79it/s]
VAL Epoch [3/5],          Loss: 0.2718,          Perplexity: 1.3123, Accuracy: 88.8800%

100%|          | 625/625 [08:03<00:00, 1.29it/s]
Epoch [4/5],          Loss: 0.2264,          Perplexity: 1.2540, Accuracy: 91.0150%
, Elapsed Time: 483.53 sec
100%|          | 157/157 [00:42<00:00, 3.66it/s]
VAL Epoch [4/5],          Loss: 0.2757,          Perplexity: 1.3174, Accuracy: 89.1200%

100%|          | 625/625 [07:58<00:00, 1.31it/s]
Epoch [5/5],          Loss: 0.2103,          Perplexity: 1.2340, Accuracy: 91.8200%
, Elapsed Time: 478.51 sec
100%|          | 157/157 [00:42<00:00, 3.68it/s]
VAL Epoch [5/5],          Loss: 0.2694,          Perplexity: 1.3092, Accuracy: 89.3800%
```

Waiting for W&B process to finish... (success).

Run history:

accuracy_train		accuracy_train	0.9182
accuracy_val		accuracy_val	0.8938
loss_train		loss_train	0.21027
loss_val		loss_val	0.26939
ppl_train		ppl_train	1.23401
ppl_val		ppl_val	1.30917

Run summary:

View run **worldly-darkness-8** at: [https://wandb.ai/eltoncn/IA024-04-BERT\\_IMDB/runs/dc0slg44](https://wandb.ai/eltoncn/IA024-04-BERT_IMDB/runs/dc0slg44)

Synced 5 W&B file(s), 0 media file(s), 0 artifact file(s) and 0 other file(s)

Find logs at: `.\wandb\run-20240408_205318-dc0slg44\logs`

Plotamos os gráficos das estatísticas obtidas durante o treinamento, onde podemos observar que o modelo conseguiu treinar corretamente, mas com overfitting:

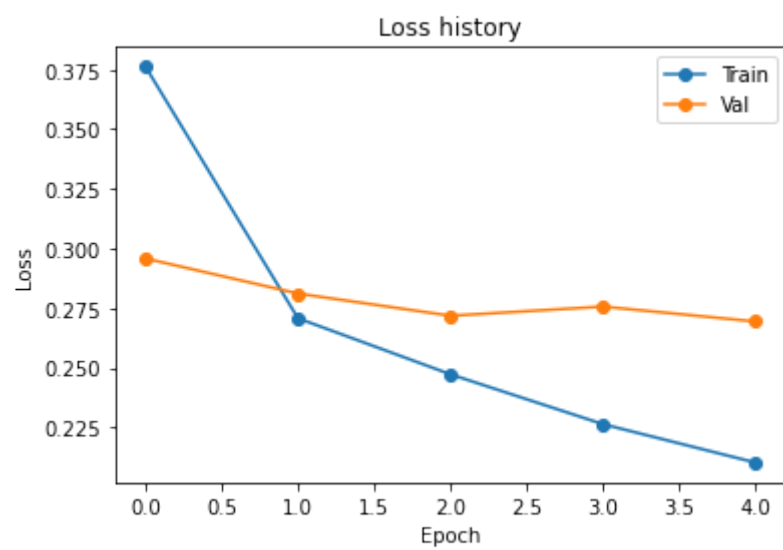
In [ ]:

```
plt.plot(hist["loss_train"], "o-")
plt.plot(hist["loss_val"], "o-")

plt.legend(["Train", "Val"])
plt.xlabel("Epoch")
plt.ylabel("Loss")
plt.title("Loss history")

plt.show()
```

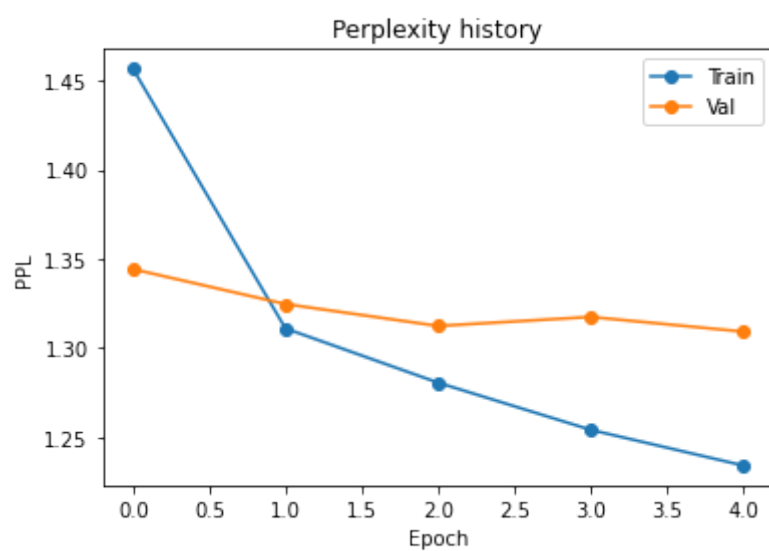




```
In [ ]: plt.plot(hist["ppl_train"], "o-")
plt.plot(hist["ppl_val"], "o-")

plt.legend(["Train", "Val"])
plt.xlabel("Epoch")
plt.ylabel("PPL")
plt.title("Perplexity history")

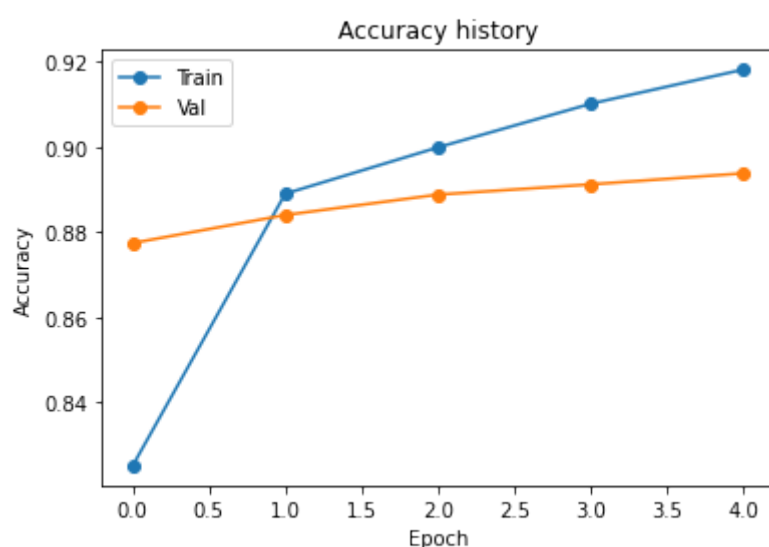
plt.show()
```



```
In [ ]: plt.plot(hist["accuracy_train"], "o-")
plt.plot(hist["accuracy_val"], "o-")

plt.legend(["Train", "Val"])
plt.xlabel("Epoch")
plt.ylabel("Accuracy")
plt.title("Accuracy history")

plt.show()
```



E podemos por fim salvar o modelo para uso posterior:

```
In [ ]: torch.save(model, model_name)
```

## Avaliação

Para avaliação começamos calculando as estatísticas no dataset de teste:

```
In [ ]: test_loss, test_accuracy = compute_loss(model, dataloaders["test"], criterion, mode=MODE_EVALUATE)
test_ppl = ppl(test_loss)

test_loss.item(), test_ppl.item(), test_accuracy.item()
```

100%|██████████| 782/782 [03:19<00:00, 3.93it/s]

Out [ ]: (0.2559235394001007, 1.2916539907455444, 0.8965599536895752)

Calculamos quantos pesos adicionais foram necessários:

```
In [ ]: n_param_bert = sum([p.numel() for p in model.bert.parameters()])

n_param = sum([p.numel() for p in model.parameters()])
n_param-n_param_bert, (n_param-n_param_bert)/n_param_bert
```

Out [ ]: (769, 1.17961222747906947e-05)

E verificamos qualitativamente a saída do modelo:

```
In [ ]: tokens = tokenizer(''This must be gambling debt. Because only when someone threatens to break your legs
                        if you don't pay will you go and agree to make a film like this.'',
                        return_tensors="pt",
                        return_token_type_ids=False,
                        truncation=True,)

with torch.no_grad():
    logits = model(tokens["input_ids"].to(device), tokens["attention_mask"].to(device))

print("Result:", logits.item())
print(f"Is this a good movie? {torch.round(logits).item() == GOOD_MOVIE}")
```

Result: -3.0

Is this a good movie? False

Podemos observar que ele realizou corretamente a classificação