

```
1 !pip install -q transformers datasets tokenizers
```

| | |
|------------------------------|-----------------|
| Locate in Drive | .1 MB 5.4 MB/s |
| Open in playground mode | 90 kB 49.4 MB/s |
| | .3 MB 32.8 MB/s |
| | 96 kB 44.8 MB/s |
| New notebook | 95 kB 46.8 MB/s |
| | 9 kB 6.3 MB/s |
| Open notebook | 32 kB 45.2 MB/s |
| Ctrl+O | .1 MB 35.7 MB/s |
| Upload notebook | 43 kB 49.5 MB/s |
| | 92 kB 48.1 MB/s |
| | 71 kB 43.3 MB/s |
| Rename | 60 kB 44.5 MB/s |
| Move | |
| Move to trash | |
| | |
| Save a copy in Drive | |
| | |
| Save a copy as a GitHub Gist | |
| | |
| Save a copy in GitHub | |
| | |
| Save | Ctrl+S |
| | |
| Save and pin revision | Ctrl+M S |
| | |
| Revision history | |
| | |
| Download | ► |
| | |
| Print | Ctrl+P |

```
1 import torch
2 import torch.nn as nn
3 from torch.nn import Linear, LSTMCell
4 from torch.nn import LSTM
5 import torch.nn.functional as F
6 import torch.nn.init as init
7
8 from transformers import GPT2ForSequenceClassification, GPT2TokenizerFast, GPT2Config
9 from datasets import load_dataset
```

```
1 device = "cuda" if torch.cuda.is_available else "cpu"
```

```
1 !nvidia-smi
```

Thu Nov 18 19:09:38 2021

| | | | | | | | | |
|-------------------|-----------|---------------|---------------------------|--------------------|----------|--------------------|-----|----|
| NVIDIA-SMI 495.44 | | | Driver Version: 460.32.03 | | | CUDA Version: 11.2 | | |
| GPU | Name | Persistence-M | Bus-Id | Disp.A | Volatile | Uncorr. | ECC | |
| Fan | Temp | Perf | Pwr:Usage/Cap | Memory-Usage | GPU-Util | Compute | M. | |
| | | | | | | | MIG | M. |
| 0 | Tesla K80 | Off | 00000000:00:04.0 | Off | | | 0 | |
| N/A | 73C | P0 | 73W / 149W | 7879MiB / 11441MiB | 0% | Default | N/A | |
| Processes: | | | | | | | | |

| GPU | GI | CI | PID | Type | Process name | GPU Memory |
|----------------------------|----|----|-----|------|--------------|------------|
| | ID | ID | | | | Usage |
| ===== | | | | | | |
| No running processes found | | | | | | |
| ----- | | | | | | |

Датасет, который мы будем использовать сегодня – тексты из английского твиттера. Они уже почищены от никнеймов, знаков препинания и прочего.

```
1 emotion_dataset = load_dataset("emotion")
```

```

Downloading: 3.62k/? [00:00<00:00,
56.7kB/s]
Downloading: 3.28k/? [00:00<00:00,
89.4kB/s]
Using custom data configuration default
Downloading and preparing dataset emotion/default (download: :
Downloading: 1.66M/1.66M [00:00<00:00,
100% 6.72MB/s]
Downloading: 204k/204k [00:00<00:00,
100% 2.63MB/s]
```

Посмотри, из чего состоит emotion_dataset:

```

1 # emotion_dataset

1 # emotion_dataset["train"]

1 # emotion_dataset["train"]["text"][0]

1 # emotion_dataset["train"]["label"][0]

1 # len(emotion_dataset["train"])
```

Для перевода текста в токены мы будем использовать предобученный BPE-токенайзер.

```

1 tokenizer = GPT2TokenizerFast.from_pretrained("distilgpt2")
2 tokenizer.pad_token = tokenizer.eos_token # У gpt2 нет pad токенов. Вместо них воспольз
```

| | |
|--------------|-------------------------|
| Downloading: | 0.99M/0.99M |
| 100% | [00:00<00:00, 887kB/s] |
| Downloading: | 446k/446k [00:00<00:00, |
| 100% | 887kB/s] |

Подготовь класс, который принимает датасет, токенайзер и имя используемой части (`train`, `validation`, `test`). Используй его для получения данных для обучения.

P.S. Посмотри, как работает токенайзер ([docs](#)) и подумай, как его надо добавить в датасет.

Немного примеров, как работает токенайзер. Это поможет с написанием датасета.

[illegible]

1 # Если надо, попрактикуйся работать с токенайзером здесь

```
1 class TweetDataset(torch.utils.data.Dataset):
2     def __init__(self, part, dataset=emotion_dataset, tokenizer=tokenizer, max_length=1
3         self.part = part
4         self.dataset = dataset
5         self.tokenizer = tokenizer
6         self.max_length = max_length
7
8         self.labels = np.unique(dataset[part]["label"])
9         self.label2num = {l: num for num, l in enumerate(self.labels)}
10
11     def __getitem__(self, idx):
12         """
13         Return dict with tokens, attention_mask and label
14         """
15         text = self.dataset[self.part]["text"][idx]
16         label = self.dataset[self.part]["label"][idx]
17
18         tokenizer_output = tokenizer.encode_plus(
19             text,
20             max_length=self.max_length,
21             padding="max_length",
22             return_tensors="pt")
23
24         target = self.label2num[label]
25         return {
26             "input_ids": tokenizer_output['input_ids'].squeeze(0),
27             "mask": tokenizer_output['attention_mask'].squeeze(0),
28             "target": target
29         }
30
31     def __len__(self):
32         """
33         Return length of dataset
34         """
35         return len(self.dataset[self.part])
```

Создай train, validation и test части датасета. Загрузи их в DataLoaders .

```
1 train_dataset = TweetDataset(part='train')
2 valid_dataset = TweetDataset(part='validation')
3 test_dataset = TweetDataset(part='test')
4
5
6 batch_size = 64 # Задай batch_size
7
8 train_loader = torch.utils.data.DataLoader(
9     dataset=train_dataset, batch_size=batch_size, shuffle=True
10 )
11 valid_loader = torch.utils.data.DataLoader(
12     dataset=valid_dataset, batch_size=batch_size, shuffle=False
13 )
```

```
9 test_loader = torch.utils.data.DataLoader(  
10     dataset=test_dataset, batch_size=batch_size, shuffle=False  
11 )
```

```
1 # batch = next(iter(train_loader))
```

► Начнем с нуля.

Попробуем обучить модель трансформер с нуля решать данную задачу.

```
[ ] ↳ 15 cells hidden
```

► Fine-tuning

Теперь другой подход: загрузим модель, которая обучалась решать задачу Language Modeling. Посмотрим, получим ли мы прирост в качестве.

```
[ ] ↳ 7 cells hidden
```

► Отчет

Покажи здесь, что ты выполнил по этой работе. Ответь на несколько вопросов:

- Какой подход оказался лучше?
Предобученная модель показала качество лучше (0.92 vs 0.89) и получалось это качество быстрее (6 vs 10 эпох)
- На какие слова модель большего всего обращала внимание?
Обе модели смотрели примерно на одинаковые слова, но предобученная модель лучше замечала смысловую связь между словами. Так же у второй модели есть головы, которые смотрят просто на предыдущее слово.
- На каких слоях/головах модель обращала внимание?
На всех слоях видно влияние других слов на слово(тавтология какая-то)

```
[ ] ↳ 1 cell hidden
```

