

Named Entity Recognition. XLM-RoBERTa

Борисочкин М. И. ИУ5-21М

```
In [1]: from datasets import load_dataset
        from transformers import DataCollatorForTokenClassification
        from transformers import AutoTokenizer, AutoModelForTokenClassification
        from transformers import TrainingArguments, Trainer

        import numpy as np
        import evaluate
```

Загрузка набора данных

Для обучения будем использовать "русскую" часть [WikiNEuRal](#)

```
In [2]: # Загрузка датасета
        dataset = load_dataset("Babelscape/wikineural")
```

```
In [3]: # Пример строки из датасета
        dataset["train_ru"][0]
```

```
Out[3]: {'tokens': ['Детство',
                    'провёл',
                    'в',
                    'Надьсомбате',
                    ',',
                    'с',
                    '1860',
                    'г',
                    '.'],
         'ner_tags': [0, 0, 0, 5, 0, 0, 0, 0, 0],
         'lang': 'ru'}
```

Предобработка данных

```
In [4]: # Загрузка токенизатора
        tokenizer = AutoTokenizer.from_pretrained("FacebookAI/xlm-roberta-base")
```

```
tokenizer_config.json: 0%|          | 0.00/25.0 [00:00<?, ?B/s]
```

```
E:\МГТУ\Магистратура\2 семестр\ММО\ДЗ\HW\env\Lib\site-packages\huggingface_hub\file_download.py:157: UserWarning: `huggingface_hub` cache-system uses symlinks by default to efficiently store duplicated files but your machine does not support them in C:\Users\Михаил\.cache\huggingface\hub\models--FacebookAI--xlm-roberta-base. Caching files will still work but in a degraded version that might require more space on your disk. This warning can be disabled by setting the `HF_HUB_DISABLE_SYMLINKS_WARNING` environment variable. For more details, see https://huggingface.co/docs/huggingface_hub/how-to-cache#limitations.
```

To support symlinks on Windows, you either need to activate Developer Mode or to run Python as an administrator. In order to see activate developer mode, see this article: <https://docs.microsoft.com/en-us/windows/apps/get-started/enable-your-device-for-development>

```
warnings.warn(message)
```

```
E:\МГТУ\Магистратура\2 семестр\ММО\ДЗ\HW\env\Lib\site-packages\huggingface_hub\file_download.py:1132: FutureWarning: `resume_download` is deprecated and will be removed in version 1.0.0. Downloads always resume when possible. If you want to force a new download, use `force_download=True`.
```

```
warnings.warn(
```

```
config.json: 0%|          | 0.00/615 [00:00<?, ?B/s]
```

```
sentencepiece.bpe.model: 0%|          | 0.00/5.07M [00:00<?, ?B/s]
```

```
tokenizer.json: 0%|          | 0.00/9.10M [00:00<?, ?B/s]
```

```
In [5]: # Пример работы токенизатора
example = dataset["train_ru"][0]
tokenized_input = tokenizer(example["tokens"], is_split_into_words=True)
tokens = tokenizer.convert_ids_to_tokens(tokenized_input["input_ids"])
tokens
```

```
Out[5]: ['<s>',
'_Де',
'т',
'ство',
'_пров',
'ёл',
'_в',
'_на',
'дь',
'сом',
'ба',
'те',
'_',
',',
'_с',
'_1860',
'_г',
'_',
',',
'_.',
'</s>']
```

```
In [6]: def tokenize_and_align_labels(examples):
        """Корректировка токенизации
        Parameters
        -----
        examples
            Входное предложение
        Returns
```

```

-----
    tokenized_inputs
    Токенизированный вход
"""
tokenized_inputs = tokenizer(
    examples["tokens"], truncation=True, is_split_into_words=True
)

labels = []
for i, label in enumerate(examples[f"ner_tags"]):
    word_ids = tokenized_inputs.word_ids(batch_index=i)  # Токенизация
    previous_word_idx = None
    label_ids = []
    for word_idx in word_ids:  # Установка значения спец. токенов -100
        if word_idx is None:
            label_ids.append(-100)
        elif (
            word_idx != previous_word_idx
        ):  # Применяем метку только к первому слову в предложении при нескольк
            label_ids.append(label[word_idx])
        else:
            label_ids.append(-100)
            previous_word_idx = word_idx
    labels.append(label_ids)

tokenized_inputs["labels"] = labels
return tokenized_inputs

```

```

In [7]: # Применение токенизатора к датасету
tokenized_dataset = dataset.map(tokenize_and_align_labels, batched=True)

```

```

Map: 0%|          | 0/12372 [00:00<?, ? examples/s]
Map: 0%|          | 0/11597 [00:00<?, ? examples/s]
Map: 0%|          | 0/9618 [00:00<?, ? examples/s]
Map: 0%|          | 0/12678 [00:00<?, ? examples/s]
Map: 0%|          | 0/11069 [00:00<?, ? examples/s]
Map: 0%|          | 0/10547 [00:00<?, ? examples/s]
Map: 0%|          | 0/13585 [00:00<?, ? examples/s]
Map: 0%|          | 0/10160 [00:00<?, ? examples/s]
Map: 0%|          | 0/11580 [00:00<?, ? examples/s]
Map: 0%|          | 0/98640 [00:00<?, ? examples/s]
Map: 0%|          | 0/92720 [00:00<?, ? examples/s]
Map: 0%|          | 0/76320 [00:00<?, ? examples/s]
Map: 0%|          | 0/100800 [00:00<?, ? examples/s]
Map: 0%|          | 0/88400 [00:00<?, ? examples/s]
Map: 0%|          | 0/83680 [00:00<?, ? examples/s]
Map: 0%|          | 0/108160 [00:00<?, ? examples/s]
Map: 0%|          | 0/80560 [00:00<?, ? examples/s]
Map: 0%|          | 0/92320 [00:00<?, ? examples/s]
Map: 0%|          | 0/12330 [00:00<?, ? examples/s]
Map: 0%|          | 0/11590 [00:00<?, ? examples/s]
Map: 0%|          | 0/9540 [00:00<?, ? examples/s]
Map: 0%|          | 0/12600 [00:00<?, ? examples/s]
Map: 0%|          | 0/11050 [00:00<?, ? examples/s]
Map: 0%|          | 0/10460 [00:00<?, ? examples/s]
Map: 0%|          | 0/13520 [00:00<?, ? examples/s]

```

```
Map: 0%|          | 0/10070 [00:00<?, ? examples/s]
Map: 0%|          | 0/11540 [00:00<?, ? examples/s]
```

```
In [8]: # Загрузка DataCollator
data_collator = DataCollatorForTokenClassification(tokenizer=tokenizer)
```

Обучение модели

Метрики качества

```
In [9]: segeval = evaluate.load("segeval")
```

```

In [10]: def compute_metrics(p):
          """Функция для расчёта метрик
          Parameters
          -----
          p
              Предсказание
          Returns
          -----
          metrics
              Метрики качества
          """
          predictions, labels = p
          predictions = np.argmax(predictions, axis=2)
          label_list = [
              "O",
              "B-PER",
              "I-PER",
              "B-ORG",
              "I-ORG",
              "B-LOC",
              "I-LOC",
              "B-MISC",
              "I-MISC",
          ]

          true_predictions = [
              [label_list[p] for (p, l) in zip(prediction, label) if l != -100]
              for prediction, label in zip(predictions, labels)
          ]
          true_labels = [
              [label_list[l] for (p, l) in zip(prediction, label) if l != -100]
              for prediction, label in zip(predictions, labels)
          ]

          results = seqeval.compute(predictions=true_predictions, references=true_labels)
          return {
              "precision": results["overall_precision"],
              "recall": results["overall_recall"],
              "f1": results["overall_f1"],
              "accuracy": results["overall_accuracy"],
          }

```

Загрузка и обучение модели

Для обучения будем использовать [базовую версию](#) XLM-RoBERTa

```

In [11]: id2label = {
          0: "O",
          1: "B-PER",
          2: "I-PER",
          3: "B-ORG",
          4: "I-ORG",
          5: "B-LOC",
          6: "I-LOC",

```

```

    7: "B-MISC",
    8: "I-MISC",
}
label2id = {
    "O": 0,
    "B-PER": 1,
    "I-PER": 2,
    "B-ORG": 3,
    "I-ORG": 4,
    "B-LOC": 5,
    "I-LOC": 6,
    "B-MISC": 7,
    "I-MISC": 8,
}

```

```

In [12]: # Загрузка модели
model = AutoModelForTokenClassification.from_pretrained(
    "FacebookAI/xlm-roberta-base", num_labels=9, id2label=id2label, label2id=label2
)

```

```

model.safetensors:  0%|          | 0.00/1.12G [00:00<?, ?B/s]

```

Some weights of XLMRobertaForTokenClassification were not initialized from the model checkpoint at FacebookAI/xlm-roberta-base and are newly initialized: ['classifier.bias', 'classifier.weight']
 You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

```

In [13]: # Аргументы для обучения
training_args = TrainingArguments(
    output_dir="XLM-RoBERTa",
    learning_rate=2e-5,
    per_device_train_batch_size=16,
    per_device_eval_batch_size=16,
    num_train_epochs=3,
    weight_decay=0.01,
    eval_strategy="epoch",
    save_strategy="epoch",
    load_best_model_at_end=True,
    push_to_hub=False,
)

```

```

In [14]: # Описание тренера
trainer = Trainer(
    model=model,
    args=training_args,
    train_dataset=tokenized_dataset["train_ru"],
    eval_dataset=tokenized_dataset["val_ru"],
    tokenizer=tokenizer,
    data_collator=data_collator,
    compute_metrics=compute_metrics,
)

```

```

In [15]: # Обучение модели
trainer.train()

```

[17310/17310 1:10:23, Epoch 3/3]

Epoch	Training Loss	Validation Loss	Precision	Recall	F1	Accuracy
1	0.052900	0.047823	0.876401	0.888962	0.882637	0.984014
2	0.035900	0.044788	0.886792	0.903350	0.894995	0.985264
3	0.024800	0.045815	0.897186	0.905217	0.901184	0.986112

```
Out[15]: TrainOutput(global_step=17310, training_loss=0.04422337971692275, metrics={'train_runtime': 4223.9292, 'train_samples_per_second': 65.569, 'train_steps_per_second': 4.098, 'total_flos': 1.1444321839124448e+16, 'train_loss': 0.04422337971692275, 'epoch': 3.0})
```

```
In [16]: # Качество лучшей модели
trainer.evaluate()
```

[722/722 00:29]

```
Out[16]: {'eval_loss': 0.044788189232349396,
          'eval_precision': 0.8867924528301887,
          'eval_recall': 0.9033498077979132,
          'eval_f1': 0.8949945593035908,
          'eval_accuracy': 0.9852637987574546,
          'eval_runtime': 32.5625,
          'eval_samples_per_second': 354.395,
          'eval_steps_per_second': 22.173,
          'epoch': 3.0}
```

```
In [17]: # Сохранение модели
trainer.save_model("XLM-RoBERTa/XLM-RoBERTa_best_model/")
```