

Проект ООП: GlossingTool

прототип интерфейса для глоссинга на локальном хосте

python: FastAPI

ML: CNN

+ небольшой фронт

Все, что касается классов

Классы

- **Lemmalnserter** основной класс для работы с подстановкой лемм и морфем
- **DataEncoder** базовый класс для энкодеров
- **EntryDataset** датасет для работы с данными
- **GlossingPreprocessor** препроцессор данных
- **CharTokenizer** токенизатор
- **Trainer** класс для обучения

критерий 1

Наследование от собственных классов

- **MorphBIOEncoder** от DataEncoder
- **SegmDataModule** от DataModule
- **MorphSegmentationCNN** от BasicNeuralClassifier

критерий 3

Наследование от стандартных классов

- **EntryDataset** от Dataset
- **TokenType** от Enum
- **BasicNeuralClassifier** от nn.Module
- Для “Duck Typing”: **DataclassLike** от Protocol

критерий 2

Наследование от абстрактных классов

наследование от ABC:

- **DataEncoder**
- **SystemPipeline**
- **BaseInference**
- **DataModule**

критерий 4

Используемые библиотеки

ML-составляющая

torch

numpy

scikit-learn

Дополнительно

jaxtyping для типизации тензоров

tqdm

python-dotenv

pyyaml

pytest

критерий 5

Декораторы

@abstractmethod в абстрактных классах

@classmethod ExperimentConfig.from_yaml()
InferenceConfig.from_yaml()

@staticmethod EntryDataset.collate_fn()

критерий 6

Собственные декораторы

регистрация модели/датамодуля/пайплайна
для трейна с помощью декоратора, пример:

@register_model

критерий 7

Перегрузка операторов

перегрузка операторов для
истории обучения модели

Пример для TrainingHistory

```
9  class TrainingHistory:  
10 |  
11 |     def __getitem__(self, idx: int) -> Dict[str, Any]:  
12 |         """Доступ к данным конкретной эпохи по индексу."""  
13 |         if idx < 0 or idx >= len(self.epochs):  
14 |             raise IndexError(f"Индекс {idx} вне диапазона: всего {len(self.epochs)} эпох")  
15 |  
16 |         return {  
17 |             'epoch': self.epochs[idx],  
18 |             'train_loss': self.train_losses[idx],  
19 |             'val_loss': self.val_losses[idx],  
20 |             'val_accuracy': self.val_accuracies[idx],  
21 |             'val_word_accuracy': self.val_word_accuracies[idx]  
22 |         }  
23 |  
24 |     def __iadd__(self, epoch_data: Dict[str, Any]) -> 'TrainingHistory':  
25 |         """Добавление данных новой эпохи через +=."""  
26 |         self.epochs.append(epoch_data['epoch'])  
27 |         self.train_losses.append(epoch_data['train_loss'])  
28 |         self.val_losses.append(epoch_data['val_loss'])  
29 |         self.val_accuracies.append(epoch_data['val_accuracy'])  
30 |         self.val_word_accuracies.append(epoch_data['val_word_accuracy'])  
31 |  
32 |         return self
```

критерий 8

Исключения

NotImplementedError

RuntimeError IndexError ValueError

критерий 9

Кастомные исключения

LengthMismatch DataMissing
ComponentInitializationError

```
4  class DataError(Exception):  
5  |     """Базовое для ошибок с данными""""  
6  |     pass  
7  |  
8  class LengthMismatch(DataError):  
9  |     """Ошибка соответствия по количеству в GlossingPipeline""""  
10 |     def __init__(self, text: str, sent: str, value_one: Any, value_two: Any):  
11 |         super().__init__(f"Количество слов в тексте {text} в предложении {sent} не совпадает с данными: {len(value_one)} vs {len(value_two)}")  
12 |  
13 class DataIsMissing(DataError):  
14 |     """Ошибка отсутствия данных""""  
15 |     def __init__(self):  
16 |         super().__init__(f"Данные не были предоставлены")  
17 |  
18 class ComponentError(Exception):  
19 |     """Базовое для ошибок компонентов""""  
20 |  
21 class ComponentInitializationError(ComponentError):  
22 |     """Ошибка инициализации компонента""""  
23 |     def __init__(self, component_name: str, reason: str):  
24 |         self.component_name = component_name  
25 |         self.reason = reason  
26 |         super().__init__(f"Не получилось инициализировать модель {component_name}: {reason}")
```

критерий 10

Синтаксический сахар питона

Использование моржовых операторов в регулярных выражениях re.match для group

```
def _check_content_type(self, line: str) -> Tuple[str, str]:
    if segmentation := re.match(r"\d+(:_\d+)*>\s*(.*)$", line):
        return 'segmented', segmentation.group(1)
    if glossing := re.match(r"\d+(:_\d+)*<\s*(.*)$", line):
        return 'glossed', glossing.group(1)

    if line.startswith("#"):
        meta = line.lstrip("# \t")
        return 'metadata', meta
    if translation_line := re.match(r"(\d+(:_\d+)*)=(.*)$", line):
        num, translation = translation_line.groups()
        self._add_line('line_num', num)
        return 'translation', translation
```

Match-Case statements

```
for char, l in zip(ex, preds):
    # print(char, l)
    match l:
        case 1:
            if prev_label is None or prev_label == 0:
                final_string += char
                prev_label = l
            else:
                final_string += f'_{char}'
        case 2:
            final_string += char
            prev_label = l
            if prev_label is None or prev_label == 0:
                prev_label = 1
        case 0:
            final_string += '\t'
            prev_label = l
final_data.append(final_string)
```

критерий 11

Распаковка словарей конфигов через **

```
@classmethod
def from_yaml(cls, path: str) -> 'InferenceConfig':
    """Загрузить конфигурацию из YAML"""
    with open(path) as f:
        config_dict = yaml.safe_load(f)

    task = config_dict['task']
    model_dict = config_dict['model']

    if task == 'segmentation':
        model = SegmentationConfig(**model_dict)
    else:
        model = ModelConfig(**model_dict)

    return cls(
        name=config_dict['name'],
        task=task,
        model=model,
        training=TrainingConfig(**config_dict.get('training', {})),
        tokenizer=TokenizerConfig(**config_dict.get('tokenizer', {})))
    )
```

Использование f-strings Типизация кода

ЮНИТ-ТЕСТЫ

```
tests > test_tokenizer.py > TestCharTokenizer > test_get_tokens_ids_unknown_char
9
10 @pytest.fixture
11 def tokenizer():
12     """Фикстура для создания экземпляра токенизатора"""
13     return CharTokenizer(symbols)
14
15 class TestCharTokenizer:
16     """Тесты для класса CharTokenizer"""
17
18     def test_initialization(self, tokenizer):
19         """Тест инициализации токенизатора"""
20         assert tokenizer is not None
21         assert hasattr(tokenizer, 'id2char')
22         assert hasattr(tokenizer, 'char2id')
23         assert isinstance(tokenizer.id2char, dict)
24         assert isinstance(tokenizer.char2id, dict)
25
26     def test_tokenize_simple_word(self, tokenizer):
27         """Тест токенизации простого слова"""
28         result = tokenizer.tokenize(['слово'])
29         assert isinstance(result, list)
30         assert len(result) == 1
31         assert isinstance(result[0], list)
32
33     def test_tokenize_word_with_hyphen(self, tokenizer):
34         """Тест токенизации слова с дефисом"""
35         result = tokenizer.tokenize(['слово-морфема-ещё-одна'])
36         assert isinstance(result, list)
37         assert len(result) == 1
38         # Дефис должен быть удален при токенизации
39         chars = ''.join(result[0])
40         assert '-' not in chars or '<PUNC>' in result[0]
41
42     def test_tokenize_multiple_words(self, tokenizer):
43         """Тест токенизации нескольких слов"""
44         words = ['слово', 'второе', 'третье']
45         result = tokenizer.tokenize(words)
46         assert len(result) == 3
47         assert all(isinstance(chars, list) for chars in result)
```

Структура

```
backend/
  — src/
    — core/
      — base_classes.py      # Базовые классы
      — config.py            # Конфигурации
      — initializing.py      # Регистрация компонентов
      — lemma_inserter.py
        — data/
      — models/
      — encoders/
      — pipelines/
      — training/
      — inference/
        — tokenization/
      — vocabularies/
    — requirements.txt
```

Структура классов

- есть разделение на абстрактные и конкретные классы
- используется паттерн для регистрации компонентов → при расширение пайплайна ничего не сломается
- есть типизация с использованием type hints + jaxtyping

Шаблоны проектирования

Template Method (Шаблонный метод)

```
class SystemPipeline(ABC):
    @abstractmethod
    def run(self, inputs: Sequence[Any]) -> Sequence[Any]:
        raise NotImplementedError
```

Registry Pattern (Реестр)

```
MODEL_REGISTRY: Dict[str, Type[BasicNeuralClassifier]] = {}
DATAMODULE_REGISTRY: Dict[str, Type[DataModule]] = {}

@register_model("lemma_tagger")
class LemmaAffixTagger(BasicNeuralClassifier):
    ...
```

Factory
Method
Pattern?

Финальный продукт

1

Glossing Tool
Полуавтоматическое глоссирование с возможностью редактирования

1 — 2 — 3 — 4

Введите предложение

Ввести вручную Из репозитория

ныух яӈгүр чи лаӈр ӈыӈдлю

Далее

2

Назад

Сегментация предложения

Совет: Кликните на морфему для редактирования. Нажмите кнопку → между морфемами для объединения.

ны-ух яӈ-гү-р| чи лаӈр ӈыӈ-д-лю

3

ныух

ГЛОССЫ ДЛЯ МОРФЕМ

Морфема: ныух
сегодня

яӈ-гү-р

яӈ гү ӈ

ГЛОССЫ ДЛЯ МОРФЕМ

Морфема: яӈ
✓ — выберите —
что.делать
как.делать
Морфема: гү
— выберите —

4

Результат глоссирования

ОРИГИНАЛ

ныух яӈгүр чи лаӈр ӈыӈдлю

СЕГМЕНТАЦИЯ

ныух яӈ-гү-р чи лаӈр ӈыӈ-д-лю

ГЛОССЫ

сегодня как.делать-CAUS-CONV.3.SG ты UNK искать-IND-INDEF

ПЕРЕВОД (РЕДАКТИРУЕМЫЙ)

Сегодня как ты охотишься на нерпу|