

# Домашние работы MLOps

girafe-ai, MIPT, fall 2025

## Общая информация

В качестве практической части курса вам нужно будет оформить одну задачу машинного обучения в репозиторий индустриального уровня со всеми необходимыми компонентами (они будут описаны в заданиях ниже). Смысл курса состоит в том, чтобы вы были способны самостоятельно произвести все необходимые операции и получить конечный результат, а не просто цифры и графики в jupyter notebook-е.

Вам нужно будет выбрать задачу, описать её в предложенном формате, создать репозиторий и работать над ним постепенно, выполняя выданные задания неделя за неделей с помощью коммитов в основную ветку этого репозитория. Ссылку на репозиторий вы передадите однажды, а далее по ней мы будем проверять наличие сделанных заданий в указанные даты.

## Дедлайны и баллы

Задания с критериями и баллами будут публиковаться в этом документе по мере того, как мы с вами будем проходить темы в курсе.

Проверка работ также будет происходить поэтапно и публиковаться в канале курса.

Для сдачи работ будет два дедлайна:

- Мягкий: примерно неделя от момента выдачи
- Жёсткий: примерно неделя от мягкого

До мягкого дедлайна вы получите полные баллы. После мягкого дедлайна, но до жёсткого будет применён понижающий коэффициент 0.7 на заработанные вами баллы. После жёсткого коэффициент будет понижен до 0.5 и такие работы будут проверены к экзамену. Крайний срок сдачи работ - за две недели до даты экзамена. Работы, присланные позже, будут учитываться только на пересдаче. Это, в частности, означает, что если в семестре ничего не делать, а прийти только на экзамен (даже со сделанной работой), то можно получить только баллы за экзамен и по итогу сдачи отправиться на пересдачу.

Во второй части курса будет проведено peer review ваших проектов, об этом объявим отдельно. Экспертная оценка (ручная) будет проведена после полной реализации пайплайна обучения и peer review и будет являться финальной точкой проверки.

Экзамен будет проходить устно с опорой на проект (будут задаваться вопросы по нему, почему делали так, а не иначе, какие могут быть варианты), по формату это обычный устный экзамен, программу опубликуем ближе к концу когда пройдём большую часть.

Тесты по теории будут проходить на основе нашего телеграм-бота, там необходимо будет ответить на несложные вопросы по материалам занятий. Один тест занимает не более 15 минут. Чтобы подготовиться к нему, полезно повторить соответствующее занятие. Это является стимулом чтобы смотреть лекции вовремя (а не в последний день перед экзаменом). Результаты этих тестов будут учтены в финальной оценке (см ниже). В конце вас ждёт финальный тест, там будут вопросы по всем занятиям, он займет примерно 2 часа. Это можно считать вариантом письменного экзамена.

Составные части оценки:

- 40% - оценка проекта по критериям
- 10% - peer review
- 30% - экзамен
- 20% - тесты по теории

Границы оценок:

- 90% - отлично
- 75% - хорошо
- 60% - удовлетворительно

## Task 1: Project proposal

Вам необходимо выбрать тему проекта для дальнейшей реализации.  
Описать умеренно детально, что вы собираетесь делать.

Выбор проекта обсуждался на первых занятиях, возможно стоит пересмотреть этот блок чтобы найти ответы на свои вопросы, некоторые ключевые моменты освещены ниже.

Шаблон описания проекта будет создан для каждого персонально. Вам нужно будет взять свой шаблон и заполнить его информацией о своём проекте. Далее внутри этого документа мы с помощью комментариев сможем обсуждать содержание вашей работы (если это потребуется).

Как работать с шаблоном:

- Заполнить название проекта и своё ФИ
  - название должно быть человекочитаемым и понятным неспециалисту в этой области
  - Как правило названия длиннее 5 слов воспринимаются тяжело
- В документе нужно сохранить
  - структуру заголовков, их размер и начертание
  - оформление основного текста: начертание шрифта, размер шрифта
  - можно использовать картинки, выделения, подчёркивания, цвета и проч.Главное не ломайте структуру документа - разбираться в неструктурированном потоке мыслей нам очень сложно.

- Внутри каждого заголовка вместо текста в кавычках нужно написать ваш текст для этого раздела.
  - Необходимо именно ваше описание, как вы это видите и понимаете.
  - Например, просто оставить ссылку на датасет будет недостаточно, нужно самостоятельно описать этот датасет в разрезе количества объектов, объёма, времени создания, его особенностей и проч.
- Во всех разделах нужно приводить ссылки на те или иные утверждения, компоненты, библиотеки, которые вы используете.

Требования к проекту работы:

- Над проектом должен работать один человек (не играем в поезда с вагонами)
- У каждого студента должен быть уникальный проект
  - Данных и моделей сейчас чрезвычайно много
    - как их выбрать - обсуждалось на занятиях и есть раздел ниже
- Использовать современные индустриальные библиотеки машинного обучения
  - Не sklearn
  - Рекомендуем pytorch
- Сейчас используем только нейросети потому что они используют все обсуждаемые в курсе инструменты
  - Если вы очень хотите взять что-то иное, можно обсудить это в чате или лично после занятия
- Проект должен включать обучение модели
  - Обвязка над API OpenAI или аналога не будет использовать большую часть инструментов, которые изучаются в курсе.
  - Если вы хотите сделать что-то на основе промтov, то можно, как вариант, использовать оптимизаторы промтov, но проще выбрать тренируемую модель.
- Данные должны быть реальными и не слишком маленькими (меньше 10Мб)
  - [Репозиторий UCI](#) не подойдёт
  - Синтетические данные не подойдут
  - Не слишком большие или сложные (так скорее всего будет тяжело вам)
- Нельзя брать уже оформленные репозитории
  - репозитории прошлых лет
  - репозитории от победителей соревнований на kaggle (готовые оформленные)
  - вместо этого, например, возьмите за основу решение, оформленное в jupyter notebook, их очень много

В качестве источников вдохновения для выбора темы можем посоветовать:

- Вашу дипломную работу
- Соревнования на [kaggle.com](#)
  - Работы победителей
- Поиск Гугла по датасетам [datasetsearch.research.google.com](#)

Примеры задач для проектов (себе их брать уже нельзя):

- <https://www.kaggle.com/competitions/learning-agency-lab-automated-essay-scoring-2/discussion/517014>
- <https://www.kaggle.com/competitions/cidaut-ai-fake-scene-classification-2024/overview>
- <https://www.kaggle.com/datasets/saurabhshahane/fake-news-classification/data>
- <https://www.kaggle.com/datasets/alxmamaev/flowers-recognition>
- <https://github.com/pmernyei/wiki-cs-dataset>
- <https://www.kaggle.com/datasets/cubeai/skin-cancer-classification-for-yolov8>
- <https://www.kaggle.com/datasets/meowmeowmeowmeow/gtsrb-german-traffic-sign>
- <https://www.kaggle.com/competitions/llm-detect-ai-generated-text/data>

Примеры сделанных проектов:

- <https://github.com/s-a-v-a-n-n-a/cat-breed-detector>
- <https://github.com/Jatana/ielts-band-predictor>
- <https://github.com/arunashamil/ag-news>
- <https://github.com/kostyamyasso2002/MLOps-road-segmentation>
- <https://github.com/RenataKostolina/simpsons/>
- <https://github.com/tenebrissilvam/CRISPR-off-t-engine>
- <https://github.com/NovikovIE/mushroom-classification>

Дисклеймер: приведённые проекты не являются идеальными, то есть не нужно повторять каждую деталь оттуда. Они приведены для понимания, как может выглядеть финальный результат и что примерно ожидается от вас. Пожалуйста, следуйте актуальным материалам и заданиям курса.

Если ваш проект не будет соответствовать критериям, его нужно будет доработать согласно нашим комментария потому что далее над ним нужно будет работать до конца курса. Обычно таких людей единицы, переживать не стоит 😊.

Дедлайны:

- Мягкий: 23 ноября, 23.59 Мск
- Жёсткий: 30 ноября, 23.59 Мск

## Task 2: Training code

Необходимо создать открытый репозиторий на github, в нём создать пакет на Python, который будет являться валидным с точки зрения пакетирования и решать вашу задачу (из предыдущей части).

В репозитории должны быть реализованы (детали ниже в разделах):

- загрузка данных

- препроцессинг данных (если необходимо)
- тренировка
- подготовка к продакшенну
- инференс

Как мы будем проверять вашу работу (обобщённые шаги):

1. клонируем репозиторий
2. создаём новый чистый virtualenv (по инструкции)
3. устанавливаем зависимости (по инструкции)
4. `pre-commit install` - ожидаем успешную установку
5. `pre-commit run -a` - ожидаем успешный результат
6. запускаем тренировку
7. запускаем систему предсказания

Ниже описаны обязательные части вашего проекта: файлы и описание требований к основным этапам работы.

## Список необходимых условий

Для получения оценки нужно выполнить все из них. Работы не удовлетворяющие им будут оценены в 0 баллов.

- репозиторий доступен по предоставленной ссылке
- сдаваемая версия кода находится в основной ветке репозитория (master или main).  
Остальные ветки не учитываются.
- тема проекта соответствует заявленной в первом задании
- реализованы (хотя бы в базовом виде) все основные разделы задания (помечены звёздочкой в заголовках)

## README.md (10 баллов) \*

Фактически это инструкция по включению в проект вашего нового члена команды. Для этого вам понадобится описать две части:

- смысловое содержание проекта
- технические детали работы с ним

Первая часть это описание вашего проекта из предыдущего задания. Базово нужно его скопировать оттуда. Необходимо сохранить корректное визуальное представление (деление на заголовки разного уровня, структуру и проч). Возможно стоит поправить какие-то детали на основе того, что вы уже узнали.

Вторая часть касается технической стороны дела:

## Setup

Должен присутствовать раздел `Setup`, где описывается процедура настройки окружения вашего проекта. Мы обсуждали множество различных вариантов (`poetry`, `conda`, `uv` и множество других развлечений). Ваша инструкция должна привести к возможности продолжить разработку, а также запустить обучение и предсказание вашей модели. То есть должны быть настроены все необходимые инструменты.

## Train

Должен присутствовать раздел `Train`, в котором рассказано, как запустить тренировку вашей модели. Если у вас есть несколько этапов (загрузка данных, `preprocessing`, несколько вариантов модели и проч), нужно описать каждый из них. Обязательно привести команды, которыми нужно запускать то или иное действие потому как мы обсуждали разные варианты работы с `CLI`)

## Production preparation

Опишите шаги подготовки натренированной модели к работе, что для этого нужно сделать. Сюда могут входить перевод в `onnx`, `tensorrt`, etc.

Также в этом разделе можно описать комплектацию поставки вашей модели (какие артефакты, модули нужны для запуска).

## Infer

Смысл тот же, что и у `Train`, но тут должно быть описано, как после тренировки запустить модель на новых данных. Также нужно описать формат таких данных, дать пример (можно в виде артефакта в вашем `data storage`).

Код предсказания должен зависеть от минимального количества зависимостей, поэтому его скорее всего не стоит реализовывать в одном файле с `Train` процедурой.

## Dependencies (5 баллов) \*

Зависимости должны быть под управлением `poetry` или `uv` (на выбор), представлены в `pyproject.toml`, также не забудьте добавить `poetry.lock` или `lock` файл `uv`.

Зависимости должны успешно устанавливаться, с ними должен успешно запускаться код проекта.

Не должно быть лишних (не используемых в проекте) зависимостей.

## Code quality tools (10 баллов) \*

Необходимо использовать `pre-commit` с базовыми хуками `pre-commit`, `black`, `isort`, `flake8`, `prettier` (для не `Python` файлов).

Все хуки должны быть соответствующе настроены.

Также необходимо, чтобы запуск `pre-commit run -a` не выдавал ошибок (зелёные результаты).

## Training framework (5 баллов) \*

Базовой опцией является PyTorch Lightning, также вы можете воспользоваться другими фреймворками, которые мы обсуждали, но в них меньше возможностей.

Обучение должно быть сделано с помощью выбранного фреймворка, используя доступные в нём возможности (не нужно строить велосипеды).

## Data management (10 баллов) \*

Необходимо использовать dvc для хранения данных.

В качестве хранилища можно использовать гугл диск, либо s3 (убедитесь, что оно доступно), либо локальное хранилище.

Скачивание данных с помощью dvc необходимо встроить в имеющиеся команды `train` и `infer`, для этого у dvc есть `python api` (на крайний случай можно использовать CLI из Python кода).

Если вы используете локальное хранилище, необходимо написать функцию `download_data()`, которая скачивает ваши данные из открытых источников.

## Hydra (10 баллов) \*

Переведите основные гиперпараметры препроцессинга, обучения и постпроцессинга в yaml конфиги hydra. Сами конфиги лучше всего расположить в папке `configs` в корне репозитория.

Конфиги должны быть сгруппированы по тем операциям, которые вы проводите (например свой файл для препроцессинга, другой для основной модели, третий для сервиса и т.д.). Используйте иерархические конфиги.

В коде не должно быть магических констант. Их можно вынести либо в отдельный Python файл (если они глобальные и их не предполагается менять), либо под управление hydra.

## Logging (5 баллов) \*

Необходимо добавить логирование ваших основных метрик и функций потерь (всего не менее 3 графиков). Также в эксперимент записывать использованные гиперпараметры и версию кода (`git commit id`). Считайте, что сервер mlflow уже поднят по адресу `127.0.0.1:8080` (вы можете поднимать его локально для тестов по этому же адресу). Адрес добавьте в поле конфига.

Графики и логи положите в отдельную директорию в репозитории с названием `plots`.

Можно использовать дополнительные системы логирования (например wandb), если вам это нужно.

## Model production packaging (10 баллов)

Переведите вашу модель в onnx (5 баллов). Также можно перевести пайплайны

препроцессинга и постпроцессинга (если они не тривиальны в вашем проекте).

Эту часть можно сделать прямо в конце пайплайна обучения, либо как отдельную команду, тогда опишите это в README.

Переведите вашу модель в TensorRT (5 баллов). Оптимально это сделать из готового файла onnx, либо можно непосредственно из модели. Оформите такой перевод в виде shell файла или python cli команды. Если вы используете пример данных для такого перевода, их тоже нужно добавить в dvc.

## Inference server (10 баллов)

Можно реализовать двумя подходами: с помощью MLflow Serving (макс 5 баллов) или Trion Inference Server (макс 10 баллов).

В разделе Infer README нужно описать, как его поднять и использовать.

## Список типичных ошибок

### Что нужно делать (-3 балла)

- Не должно быть исполняемого кода на уровне файла. Используйте в таких случаях  
`if __name__ == '__main__':`
  - В частности нельзя объявлять переменные (кроме констант) на верхнем уровне файла (не внутри функции или класса)
- Не использовать `warnings.filterwarnings("ignore")`. Никогда не делайте этого в ~~предакшн~~ любых проектах. Это огромный задел на отстреливание себе ноги. Люди пишут предостережения чтобы вас предостеречь - сделайте необходимые для этого действия
- Нельзя сохранять данные в гит!!!!!!!!!!!!!! То есть файлы .json, .csv, .h5 и проч. То же касается файлов натренированных моделей (.cbm, .pth, .xyz, .onnx, etc).
- Название репозитория (в т.ч. url) должно отражать смысл вашего проекта (e.g. `cats-vs-dogs`), а не название курса (e.g. `mlops_homework`). Использовать – [dash] в качестве разделителя (не \_ [underscore])
- Нужно назвать питон пакет (aka папка с вашим кодом) по правилам питона (`snake_case`), а не любым другим способом (e.g. `MYopsTools`)
- Так же питон пакет необходимо называть согласно смыслу вашего проекта, не `src` или `my_project`

- Используйте дефолтный `.gitignore` для Питона (а не пустой), его дополняйте необходимыми вам путями. Дефолтный конфиг гуглится и даже предлагается вам гитхабом при создании репозитория.
- Файлы с кодом должны называться в `snake_case`, не в `CamelCase` (e.g. `Dataset.py`)
- Не использовать `argparse`, вместо него `fire`, `click`, `hydra` или другой инструмент для cli.
- Не использовать однобуквенные переменные кроме `i`, `j`, `k`. Вместо этого давайте переменным семантически наполненные имена (`data` или `features` вместо `X` и т.д.)

### Что делать желательно (+2 балла)

- Под вызовом `if __name__ == '__main__'`: лучше вызывать ровно одну функцию (можно её назвать `main` или как-то ещё), а не писать всю логику непосредственно под `if`-ом.
- Использовать `fire` совместно с `hydra` через `compose api`
- Сделать одну входную точку `commands.py`, где вызывать соответствующие функции из файлов
- Используйте `pathlib` вместо `os.path` - важа жизнь заиграет совершенно другими красками

Список ошибок не является полным. Способов сделать странные вещи очень много, поэтому старайтесь применять знания, полученные в курсе, а также ваше чувство прекрасного.

P.S. если что-то непонятно или вы хотите сделать как-то иначе - пишите в чат, обсудим. Данные правила не являются догмами (хотя в этом ДЗ необходимо их выполнить), а скорее набором хороших техник, которые, впрочем, нужно адаптировать под задачу и ситуацию в целом.