

Отчет о проделанной работе

Студенты: Серков Александр, Скрипко Николай

Тип: соревнование на Kaggle

<https://www.kaggle.com/competitions/geoplant-at-paiss/overview>

Задача: предсказать, какие растения можно найти в конкретной точке в Европе.

На входе: временные ряды (климат в этой местности за последние 20 лет), картинки (спутниковые фотографии местности), табулярные данные (широта, долгота, страна, регион) и много другого.

Важно: данные делятся на две категории: РО + РА.

- РО (presence only) это единичные наблюдения – где-то что-то есть
- РА (presence-absence) это полноценные исследования территорий (список ВСЕГО, что растет где-то)

На выходе: список ID растений (всего примерно 11255 классов).

Метрика: f1 macro

Что было проделано:

- Разобраны и воспроизведены открытые baseline решения
- Написано обучение и инференс PyTorch моделей с использованием предоставленных данных. Обучен каскад моделей на разных архитектурах (все архитектуры можно найти в src.models)
- Кластеризация РО данных для использования вместе с РА данными
- Обработка данных РА: encoding текстовых фичей
- Обучена модель для adaptive k (предсказывает, СКОЛЬКО растений нужно предсказать, дальше берется top k)

Что будет сделано

- Объединение каскада моделей в одну e2e модель
- Использование РО данных в обучении (сейчас использовались только РА, РО были кластеризованы лишь подготовлены – кластеризованы для использования)