

# Хакатон по метеорологии

Алексей Харламов  
Павел Остяков



НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ  
УНИВЕРСИТЕТ

# Участник 🔍

📄 Хакатон ^

Баллы

1	10011000 (pavelostaaa, MrAxcel)	84489992
2	Антихайп (cszdr, vasyarv)	82909074
3	Classics(ivanicki-i, Amir14111) 🇷🇺	80423665
4	kutcu	79701913
5	sal.vios	78750097
6	alexisonfireyar 🇷🇺	75905679
7	Ясонов Евгений	75734019
8	Антон Патрикеев	75606736
9	Aiaiai (karfly, illusionww)	74950292
10	Iviconun(FireSonics, mrk.andreev)	74761301

# Задача

- Даны три города: Москва, Санкт-Петербург и Казань
- Каждый город разбит на квадраты
- Есть данные с любительских метеостанций и данные о качестве связи у пользователей
- Для каждого квадрата и часа нужно восстановить, шёл дождь или нет
- Тренировочные данные – начало июля 2017. Тестовые – июль-август 2017.

# Что нужно было предсказать?

- Для каждого квадрата было известно precipitation - количество мм осадков, выпавших на конец часа
- Таргет: precipitation > 0.25
- Бинарная классификация
- Метрика - ROC AUC

# Данные

- Порядка 30ГБ
- 2 типа: Пользовательские и Netatmo
- Поделены по времени(час)

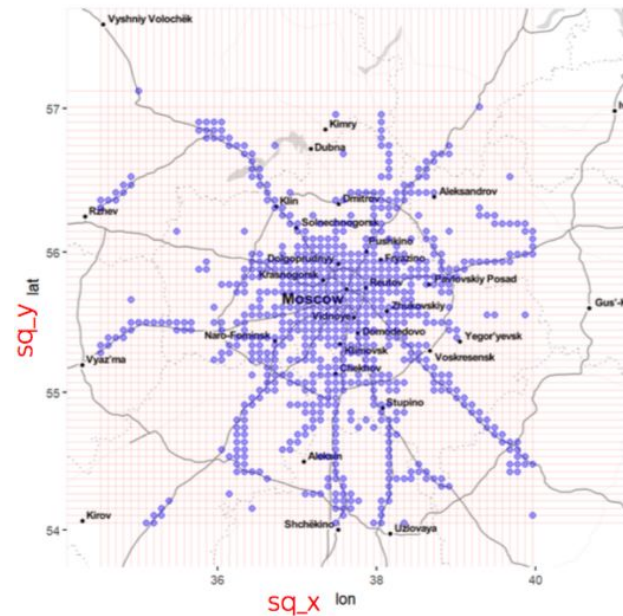


Семплированные (по Москве 0.1% пользователей),  
обфусцированные данные из Аппметрики + практически  
сырые данные любительских метеостанций  
в Москве, Санкт-Петербурге и Казани.

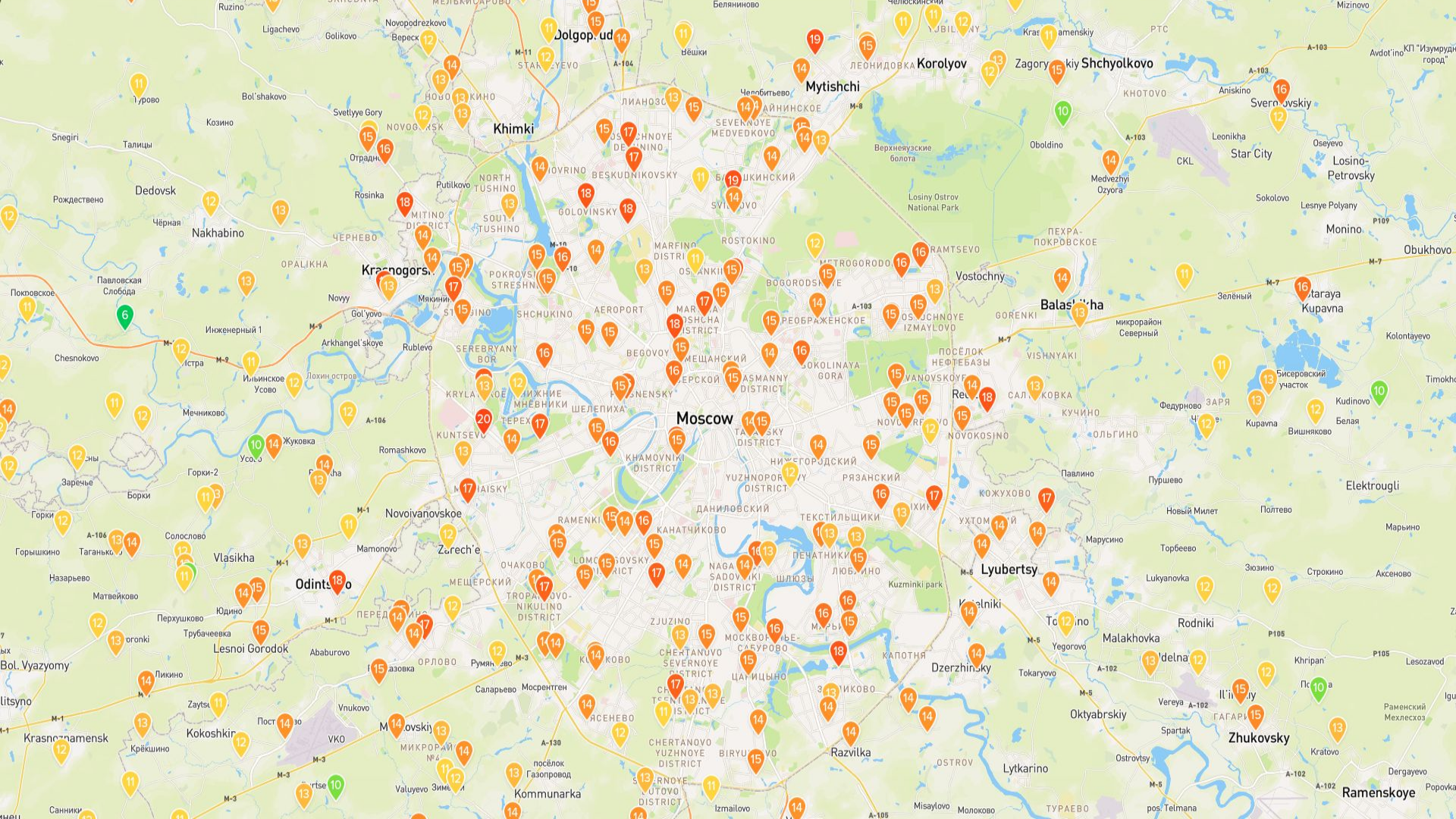


# Проблемы

- Квадраты распределены неравномерно
- Данные от пользователей очень шумные
- Много пропусков
- Данных мало







# Пользовательские факторы

- Данные о сотовом сигнале (сила сигнала, тип, оператор)
- GPS данные (высота, скорость движения, точность определения)
- Координаты пользователя
- Модель телефона



# Данные с любительских метеостанций

Метеоданные с шагом в 20 минут

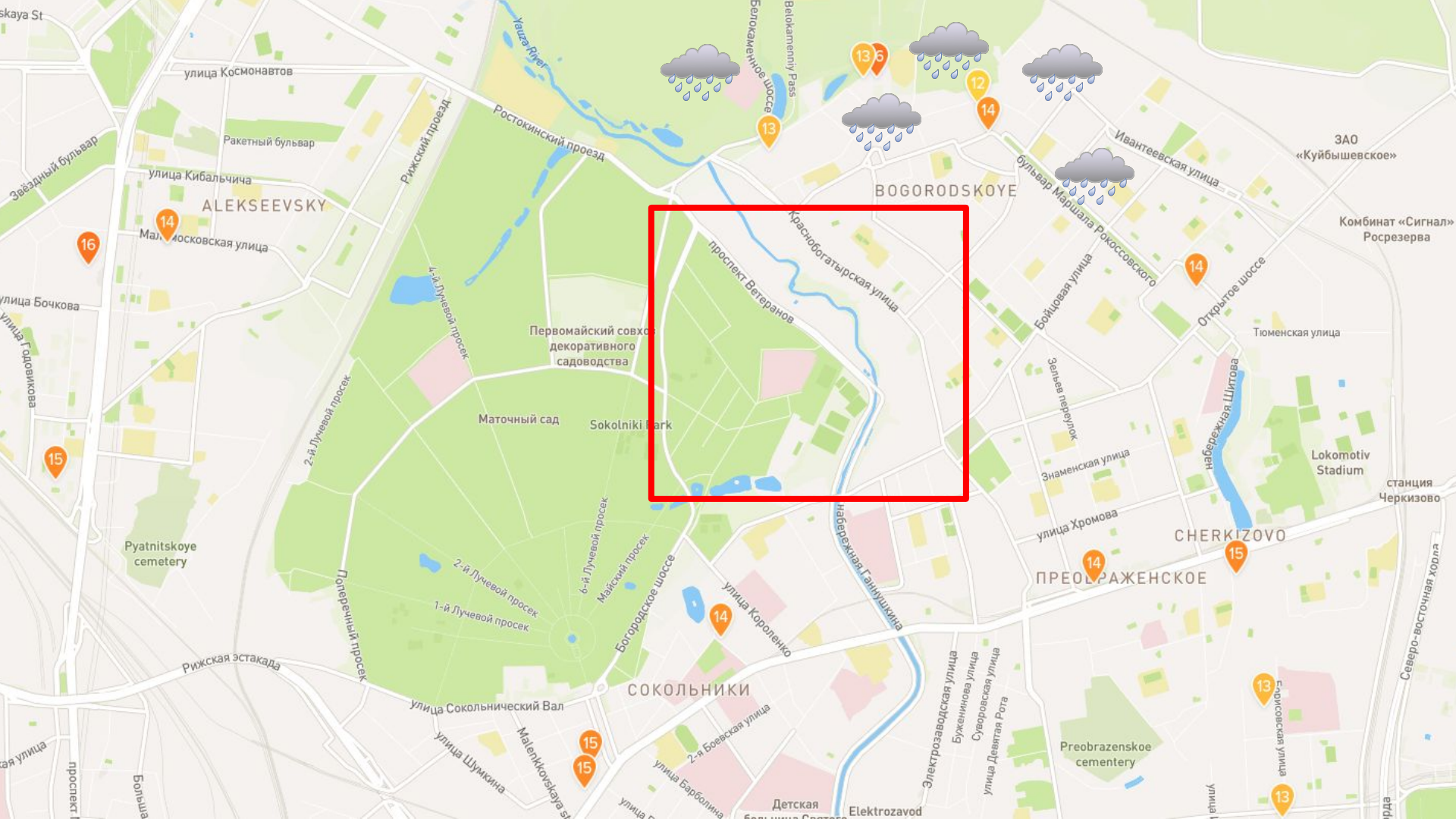
- Координаты метеостанции
- Температура
- Скорость и направление ветра
- Влажность
- Количество осадков за последний час

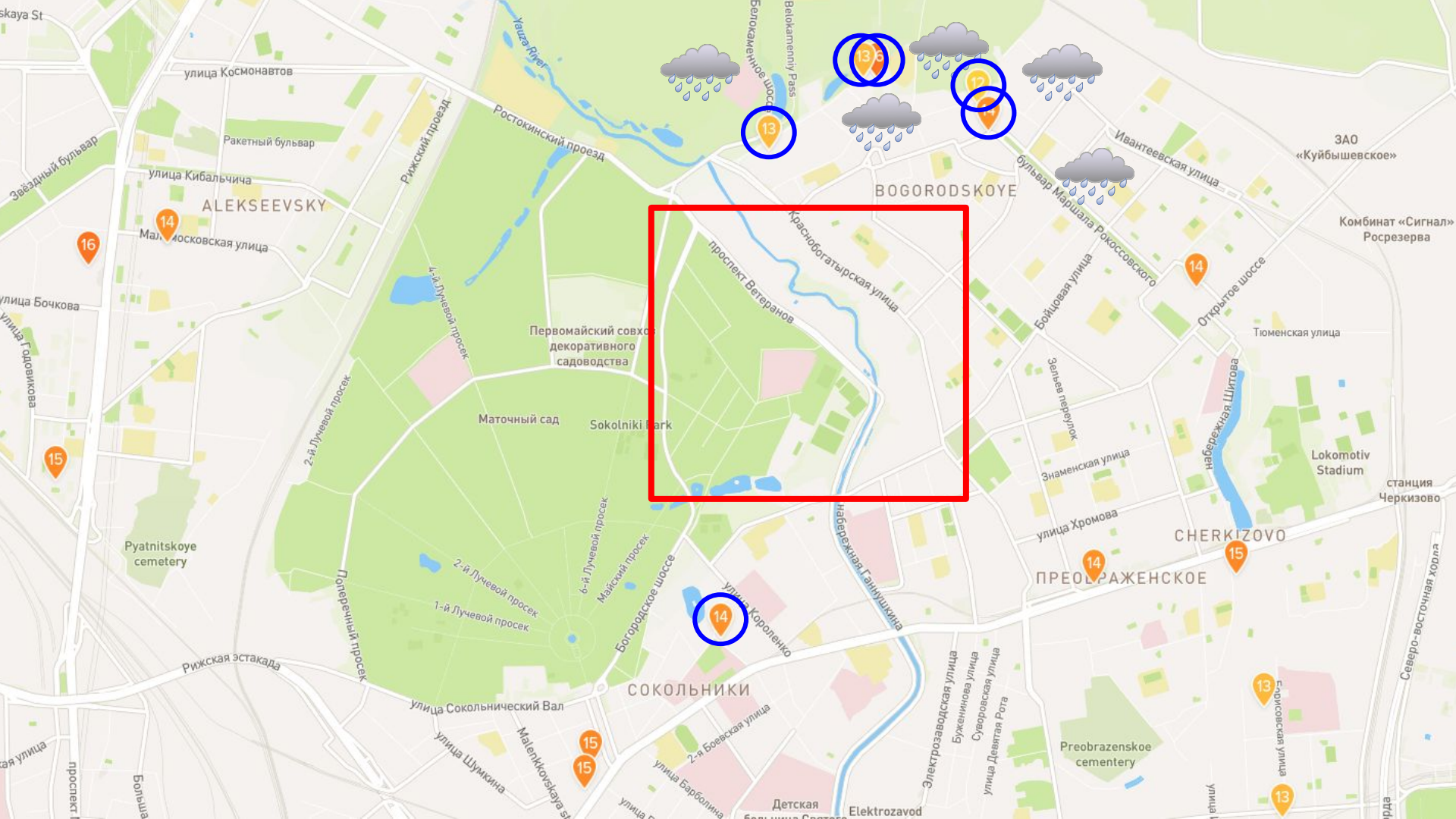
# Что сделали: факторы

- Базовые статистики для признаков(mean, var, min, max, etc.)
- 10 квантилей (0.1, 0.2, ...)
- Статистики для различных подвыборок ближайших метеостанций.
- Всего около 3000 факторов

# Что сделали: модель

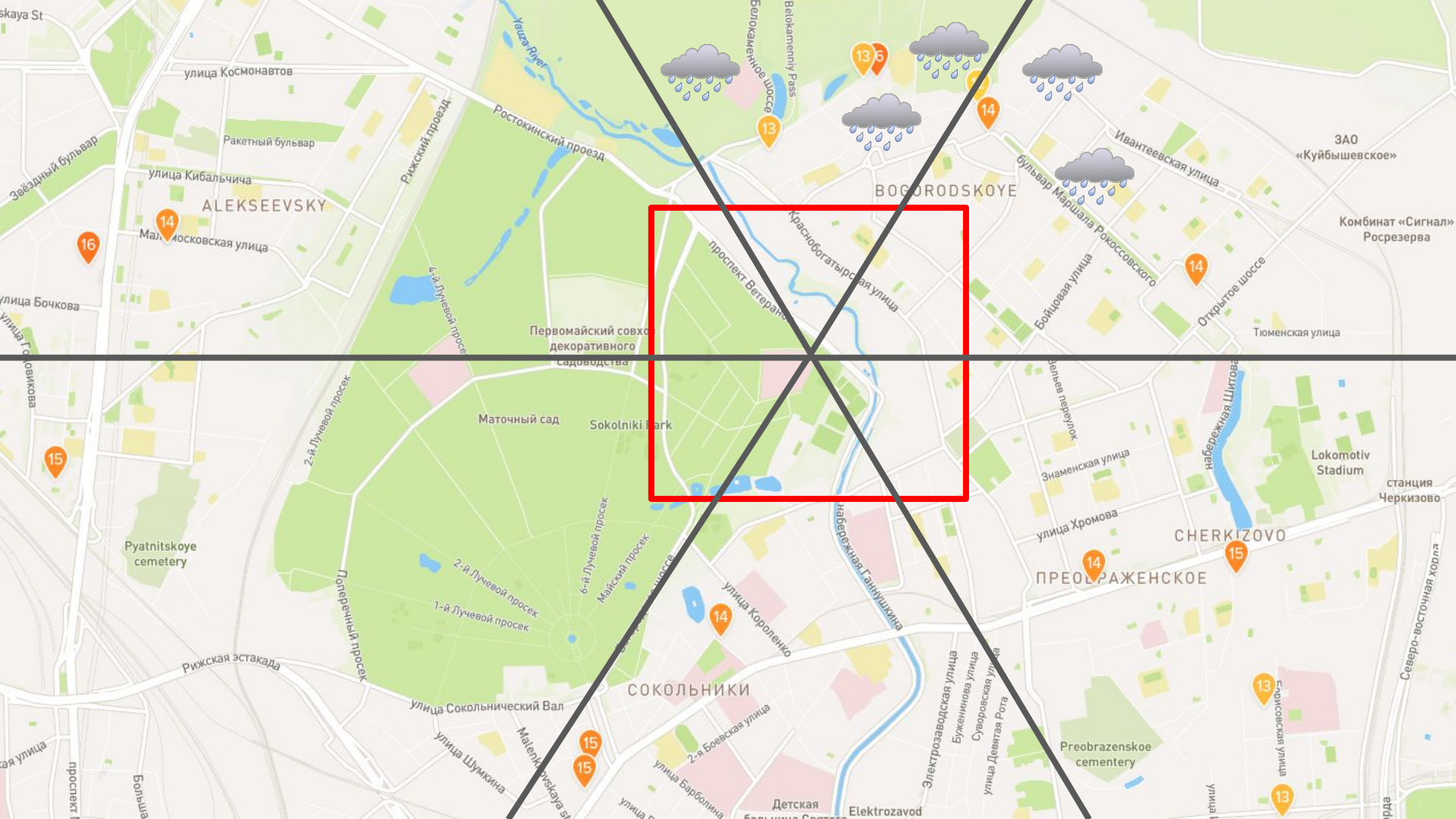
- Одна модель для всех городов.
- Catboost на дефолтных параметрах
- Получили 0.82 на лидерборде

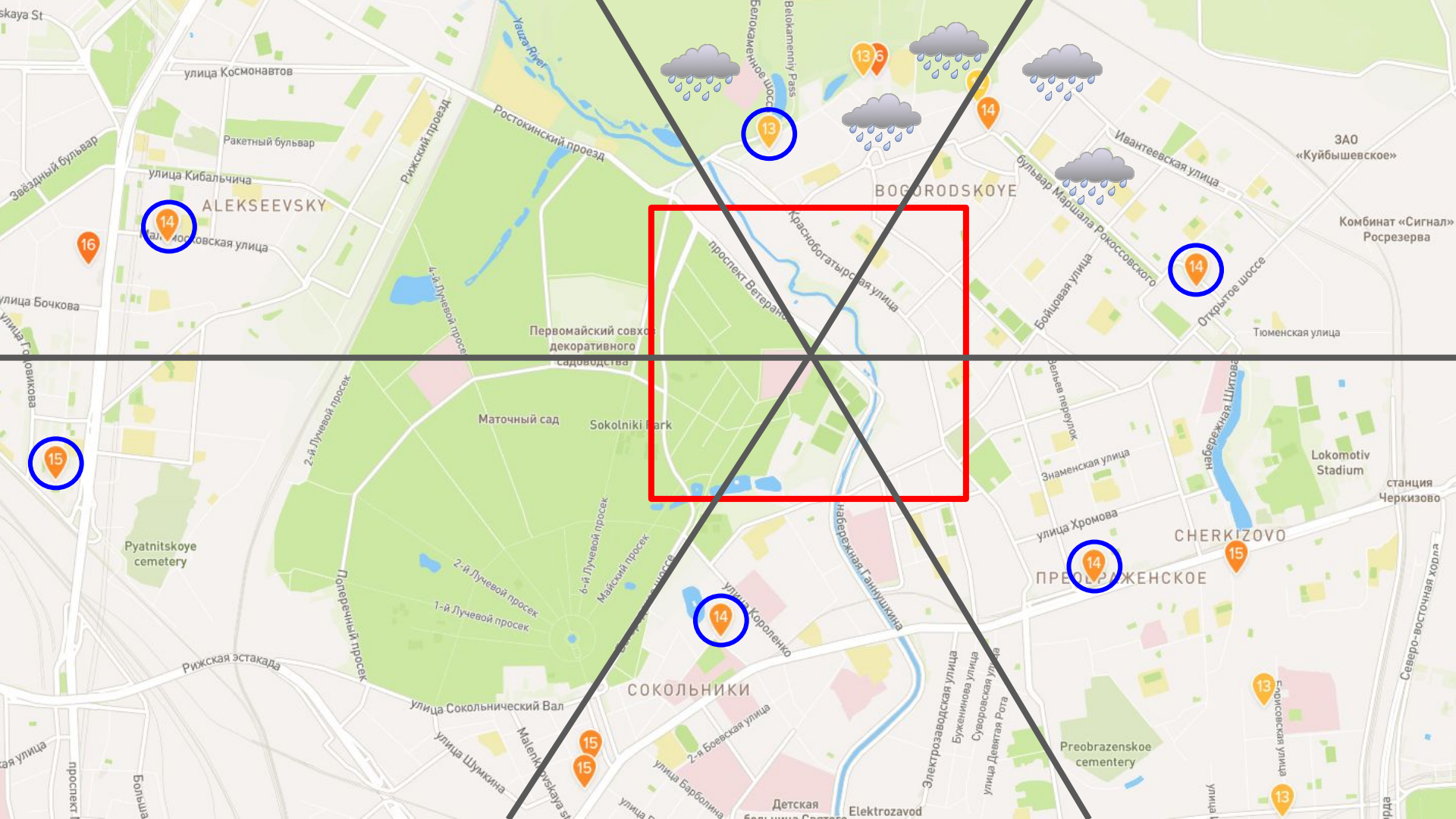




Можно добавлять не просто  
блажайшие, а брать равномерно  
вокруг







Снова считаем различные статистики,  
получаем ещё 4000 фичей и 0.842 на  
лидерборде

# Какие модели пробовали обучать?

- Линейные
- Метрические
- Нейронные сети
- Деревья
- Градиентный бустинг над решающими деревьями

Как и всегда, градиентный бустинг  
сработал, а всё остальное нет



Yandex  
CatBoost

*dmlc*  
***XGBoost***

Microsoft  
**LightGBM**





Yandex  
CatBoost

*dmlc*  
***XGBoost***

Microsoft  
**LightGBM**





Yandex  
CatBoost



*dmlc*  
***XGBoost***

Microsoft  
**LightGBM**



# Признаки внесшие наибольший вклад

- time\_of\_day
- netatmo\_pressure\_mbar
- netatmo\_sum\_rain\_1h
- netatmo\_sum\_rain\_24h
- netatmo\_humidity\_percent
- netatmo\_temperature\_c
- netatmo\_wind\_direction
- LocationSpeed

# На чем обучали?

- Два процессора Intel xeon по 14 ядер
- 256 ГБ ОЗУ
- Время генерации фичей - 2 часа
- Время обучения catboost - 10 минут

# Возможные улучшения

- Чистка данных
- Оптимизация ROC AUC directly
- Более явный учет информации с ближайших квадратов
- Большой упор на пользовательские признаки
- Явная дифференциация информации по операторам сотовой связи
- Стекинг, блендинг...
- Смена таргета (переход к задаче регрессии)

# Ссылки

Воспроизводимый код - [github.com/gamers5a/YandexMeteumSolution](https://github.com/gamers5a/YandexMeteumSolution)

Алексей Харламов - [axcel@me.com](mailto:axcel@me.com)

Павел Остяков - [pavelosta@gmail.com](mailto:pavelosta@gmail.com)