

Tinkoff Data Science Challenge

Восстановление позиции мерчанта

Иван Бендына

Постановка задачи

Мерчант — место, принимающее платежи с использованием банковской карты.

Банк хочет знать геолокацию мерчантов.

Постановка задачи

Даны чекины пользователей.

Чекин это отправка своих геокоординат в небольшой промежуток времени от покупки.

Постановка задачи

Данные очень шумные.

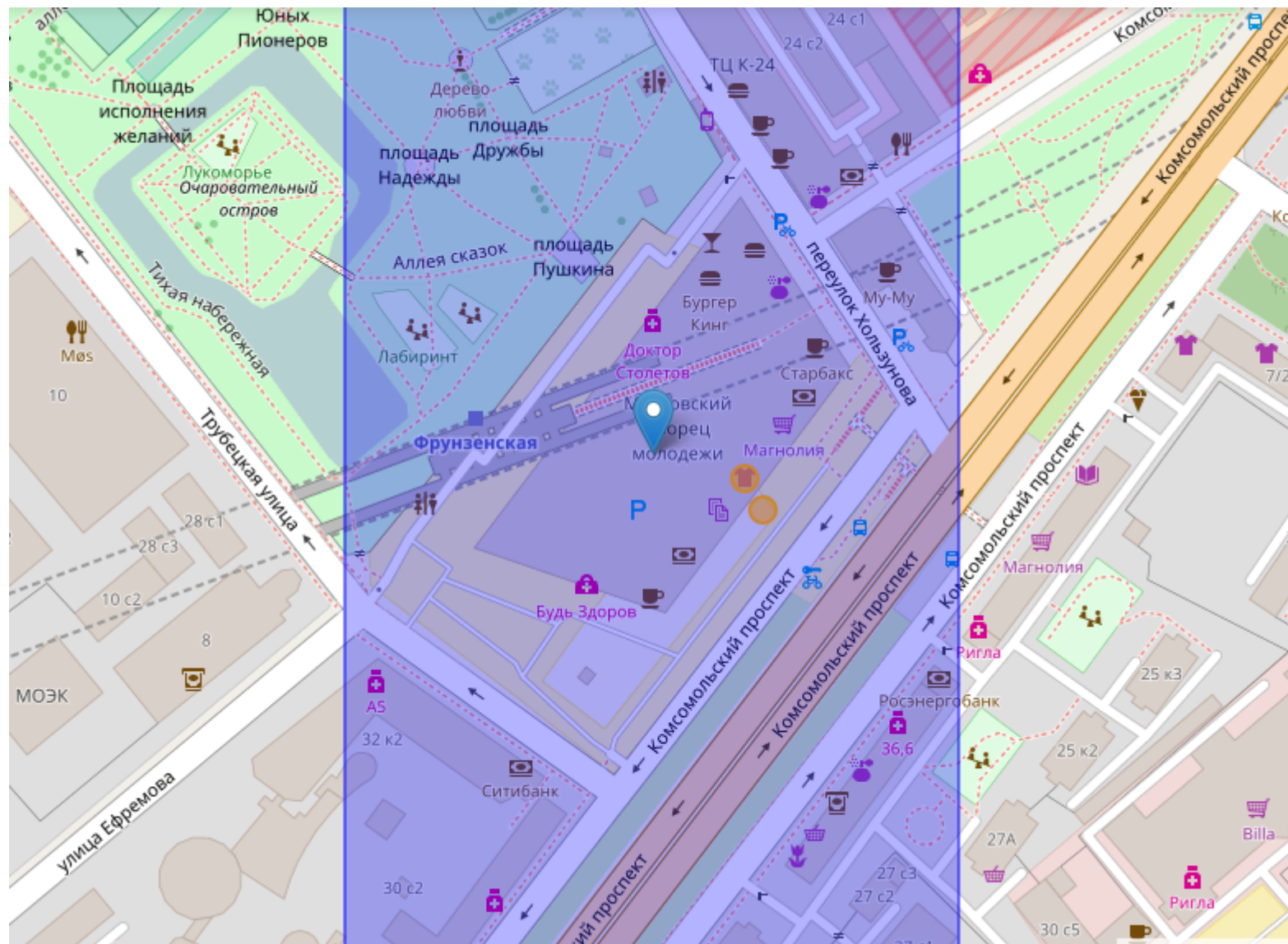
Возможные причины:

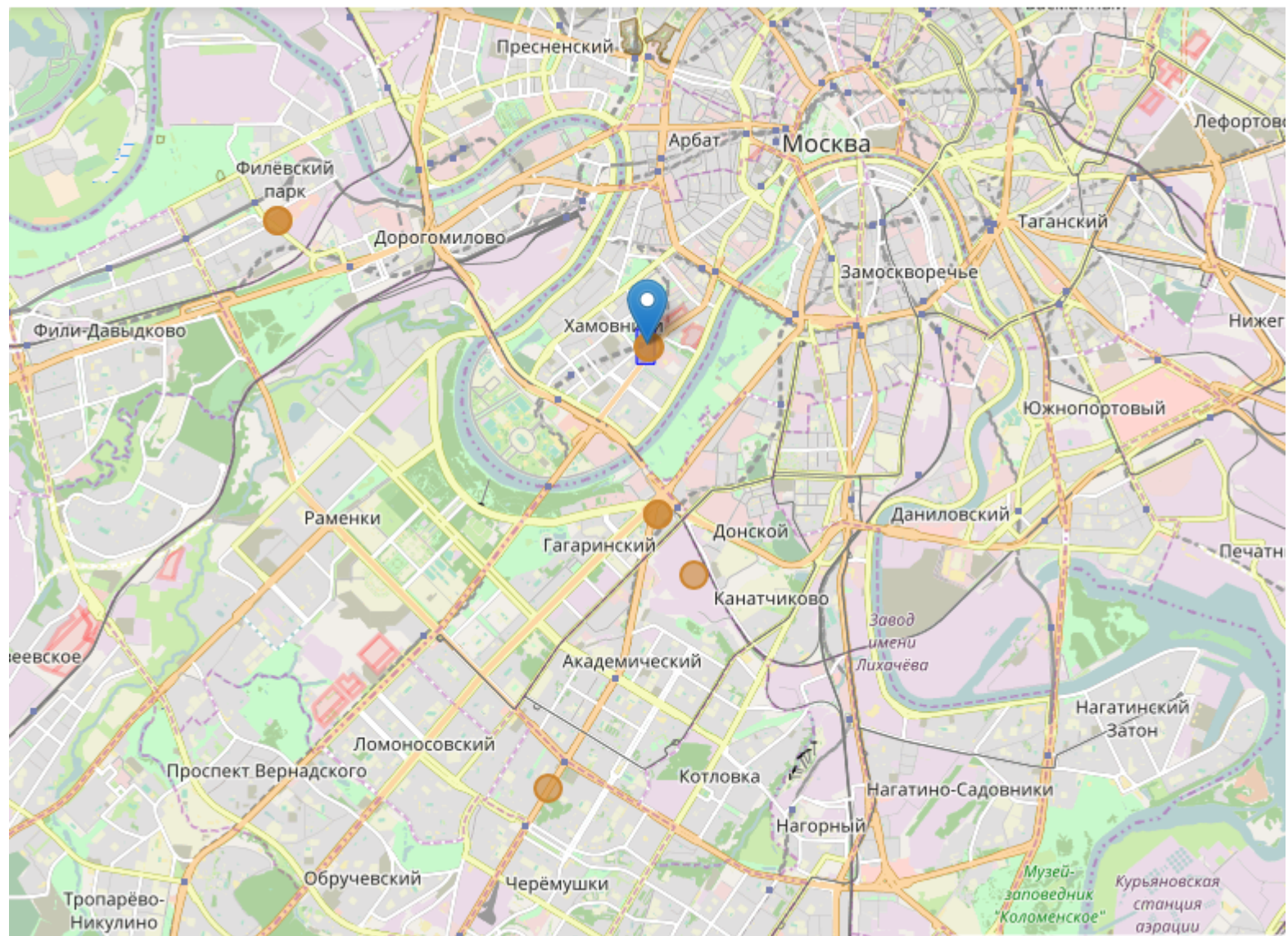
- не обновились GPS координаты на телефоне
- неточности геолокационного алгоритма (wi-fi и вышки)
- мерчант может быть только онлайн
- пользователь успел далеко отъехать
- ошибки ввода, и.т.д

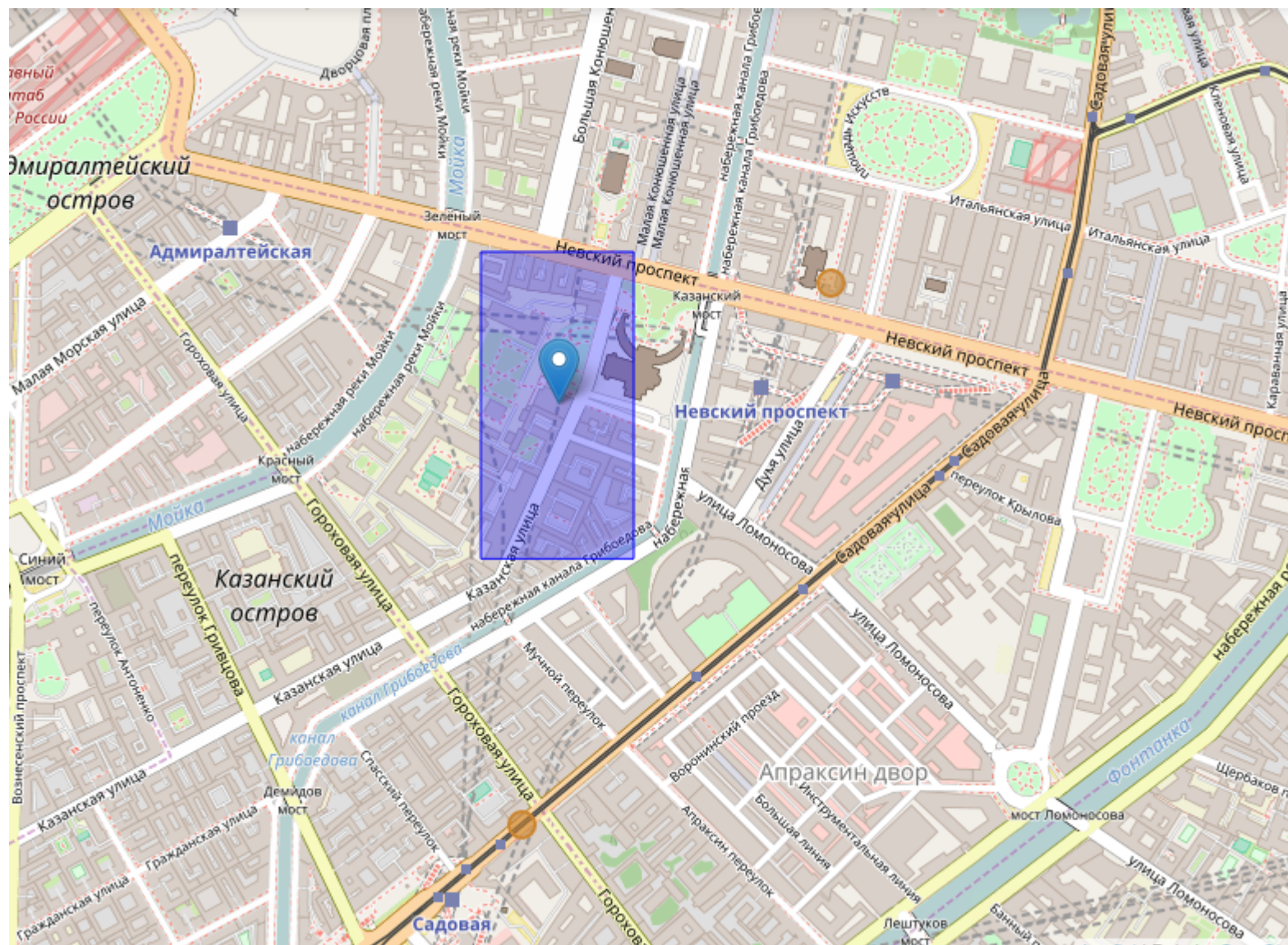
Постановка задачи

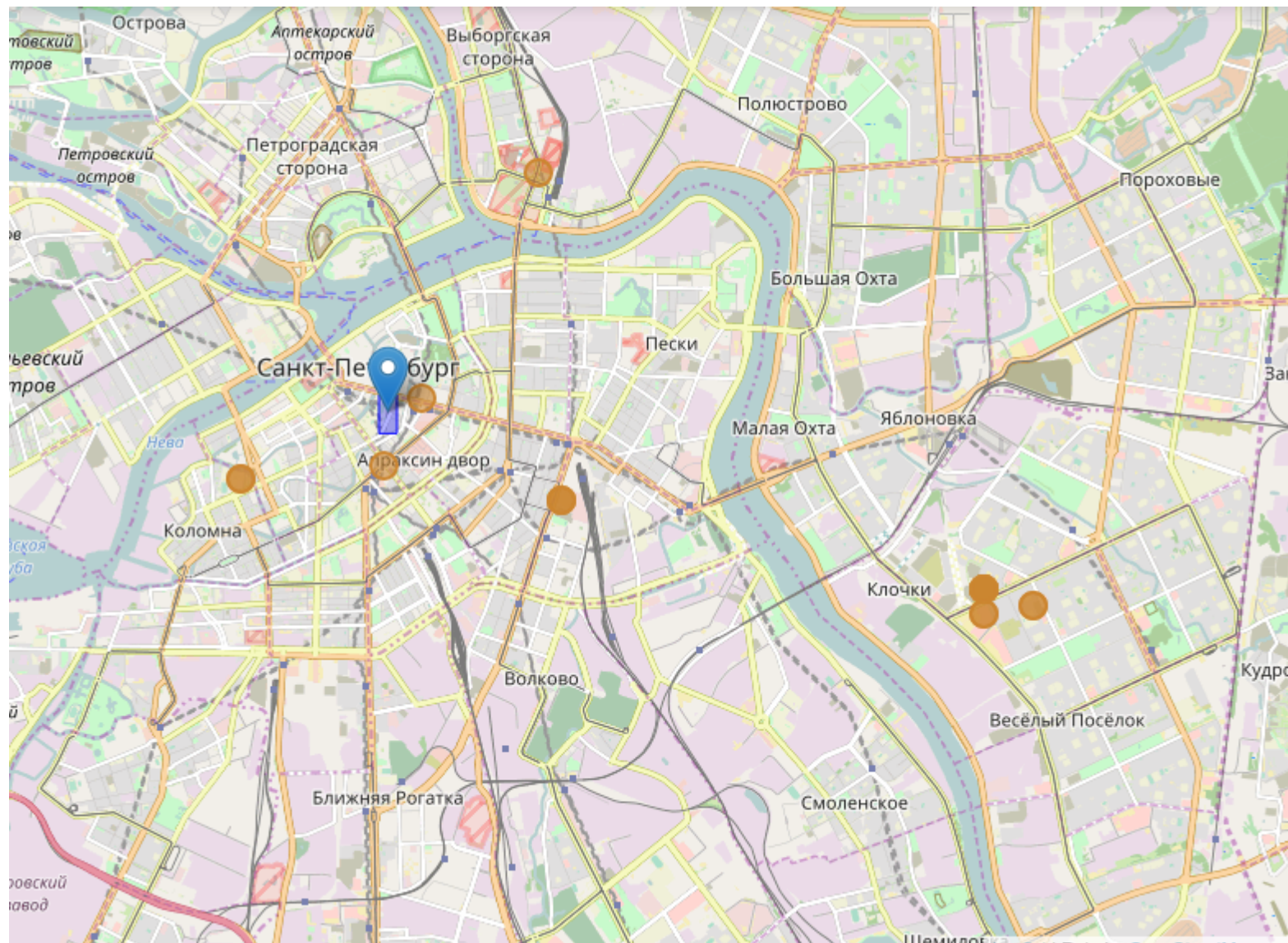
Метрика: accuracy.

Ответ правильный, если координаты отличаются не более чем на 0.002 от настоящих по широте и долготе.









Фильтрация

Частые точки чекинов:

55.75034704,	37.62385111	(27955)
55.75034704,	37.6235111	(722)
55.750347,	37.623851	(471)
55.750347043,	37.623851113	(150)
0.0,	0.0	(25018)
51.17889992,	-1.82639994	(5420)
51.1789,	-1.8264	(228)
51.172,	-1.2634	(177)

Фильтрация

Фильтруем точки чекинов не в России.

Геоданные OSM

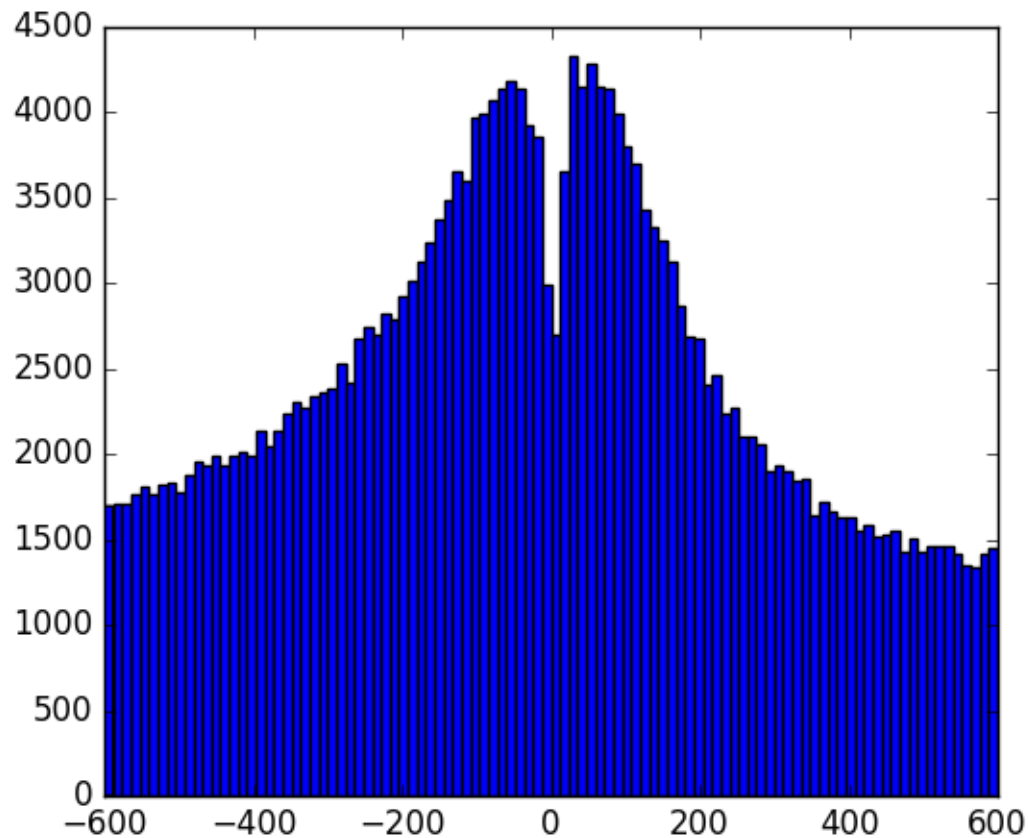
Nominatim reverse geocoder – адрес по координатам.

Можно запустить локально.

Данные в формате osm.pbf <http://download.geofabrik.de>

Россия 2Gb, развернуть на ноутбуке около 10 часов.

Распределение разницы времени транзакции и чекина



Решение

Создаем матрицу объекты-признаки, где объекты это чекины.
Таргет – это попадает ли чекин в прямоугольник.

Оптимальное решение: 0.47

Признаки

Признаки основанные на расстоянии от этого чекина до других чекинов этого же мерчанта.

- среднее расстояние
- среднее расстояние до уникальных
- перцентиль по среднему расстоянию

Признаки

Признаки основанные на расположенных рядом известных мерчантах.

- минимальное расстояние
- количество в прямоугольнике 0.002 (0.004, 0.006)

Признаки

Признаки основанные на зданиях.

- площадь здания
- количество магазинов
- суммарная/максимальная площадь зданий в радиусе 200 метров

Признаки

Признаки основанные на времени.

- разница между чекином и транзакцией
- была ли одна транзакция с разными координатами чекинов?

Модель

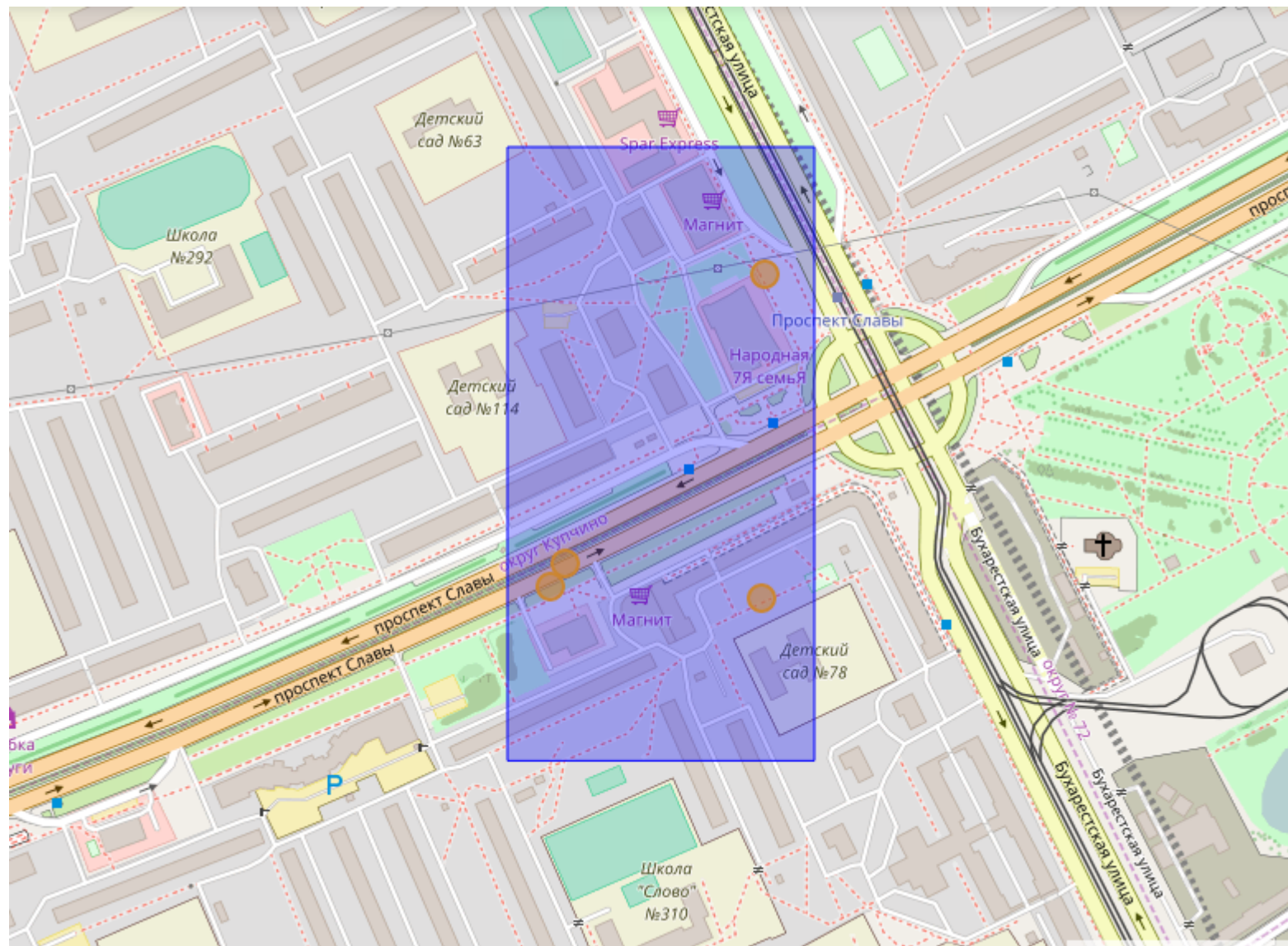
Все эти признаки в xgboost — в качестве предсказания чекин с наибольшей вероятностью.

LB score **0.360**

Объединение транзакций

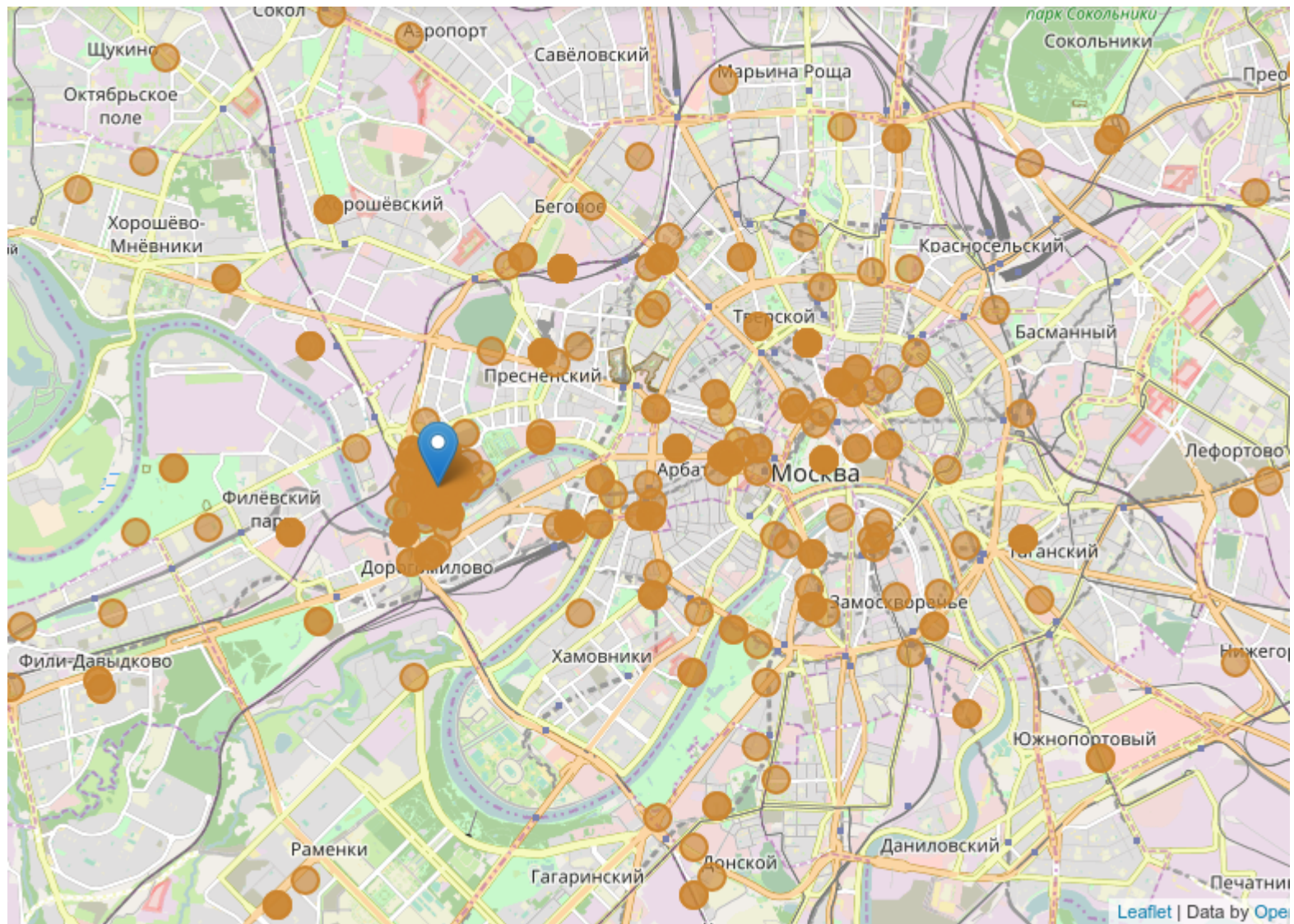
Если вероятности нескольких транзакций достаточно большие и они находятся недалеко, то в качестве предсказания – центр описанного прямоугольника этих транзакций.

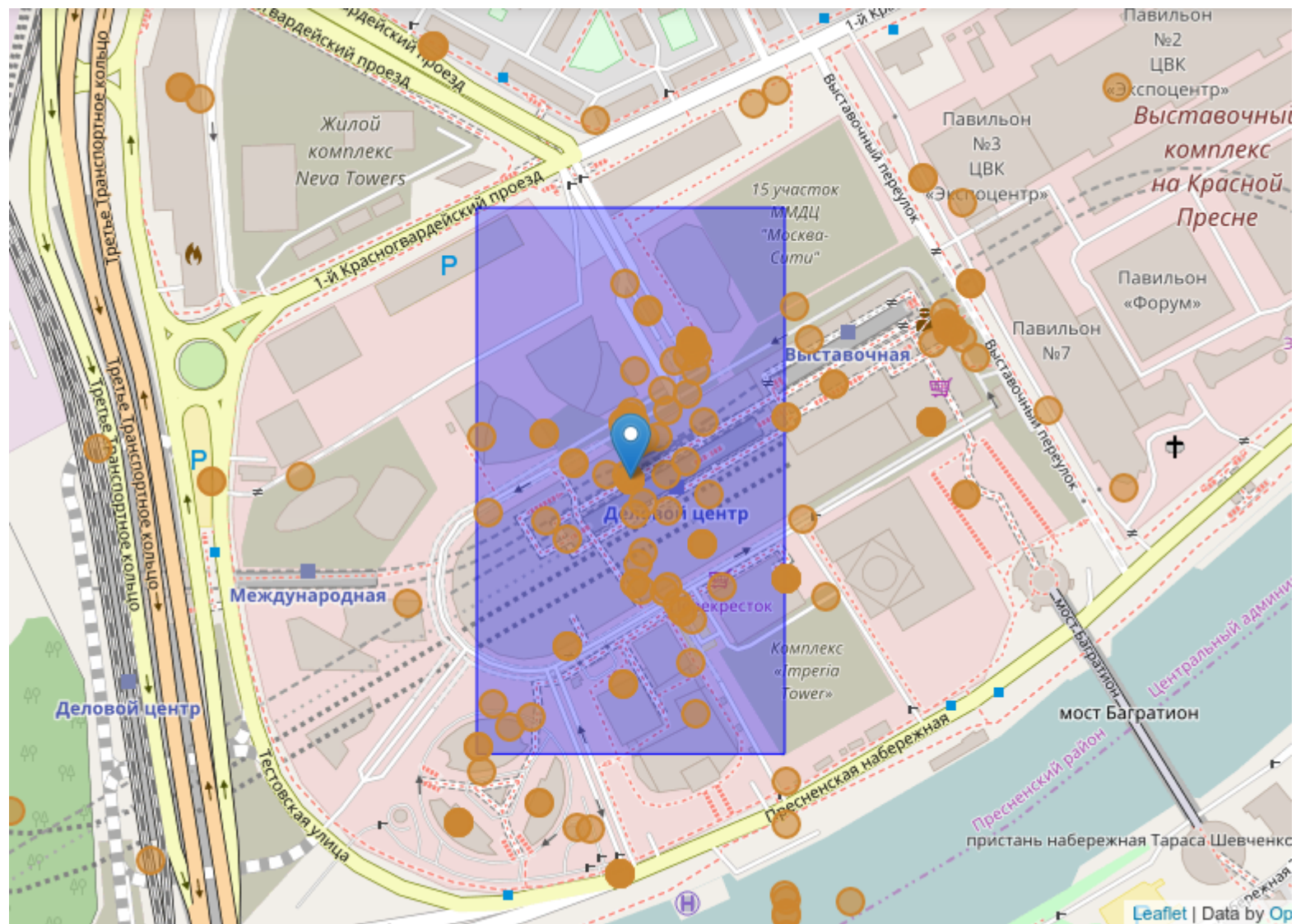
LB score **0.365**

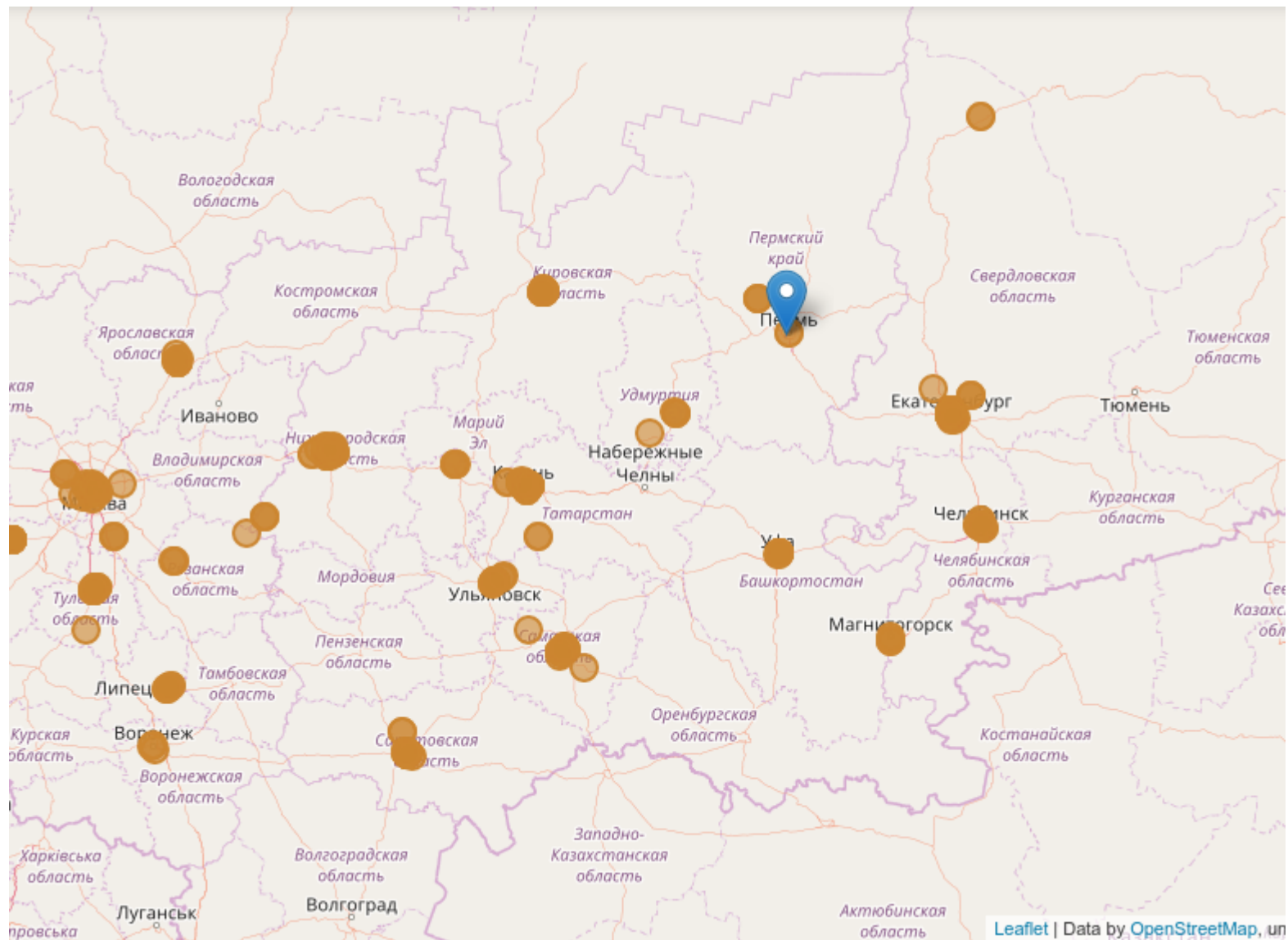


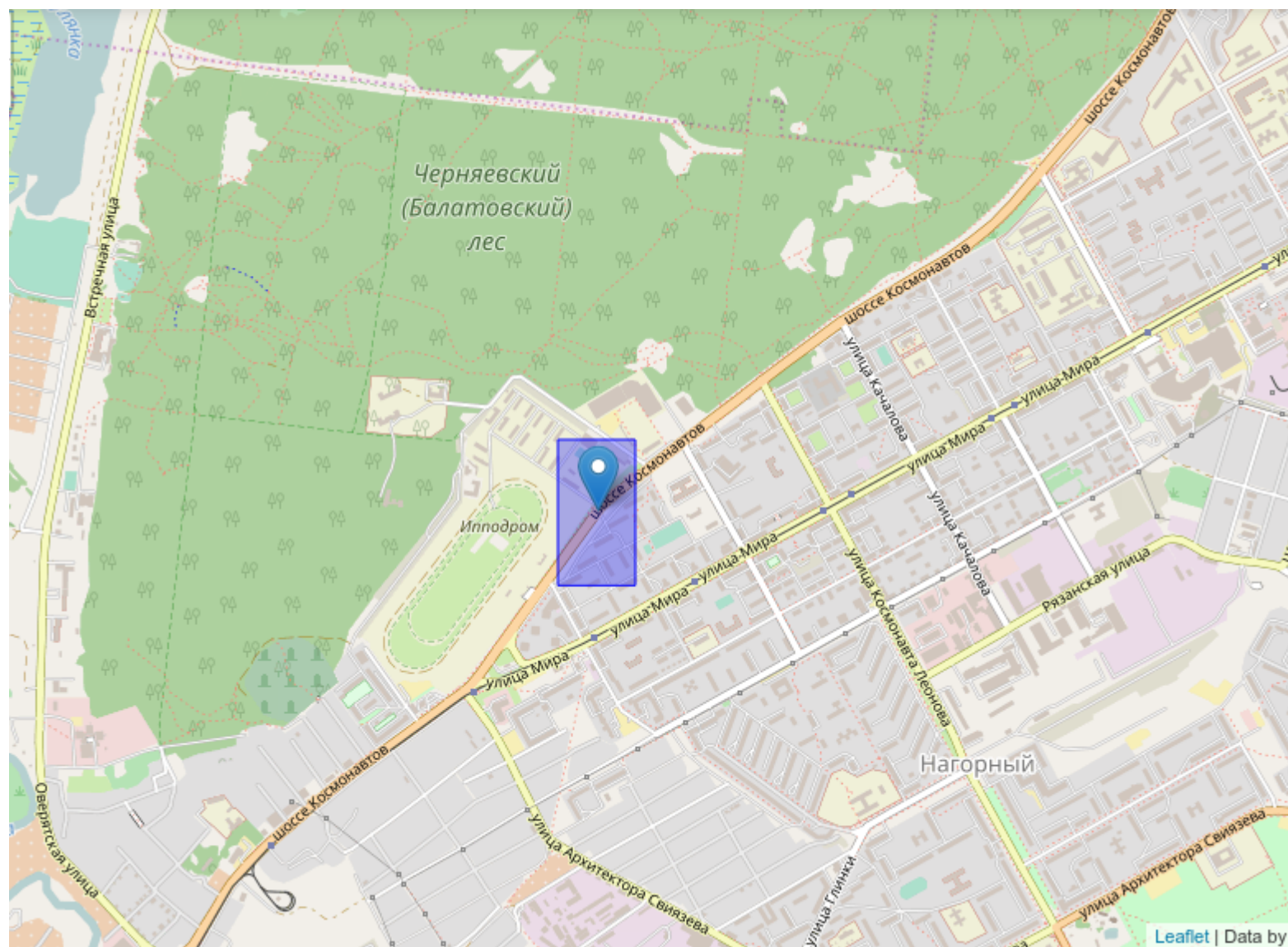
Смещение предсказания

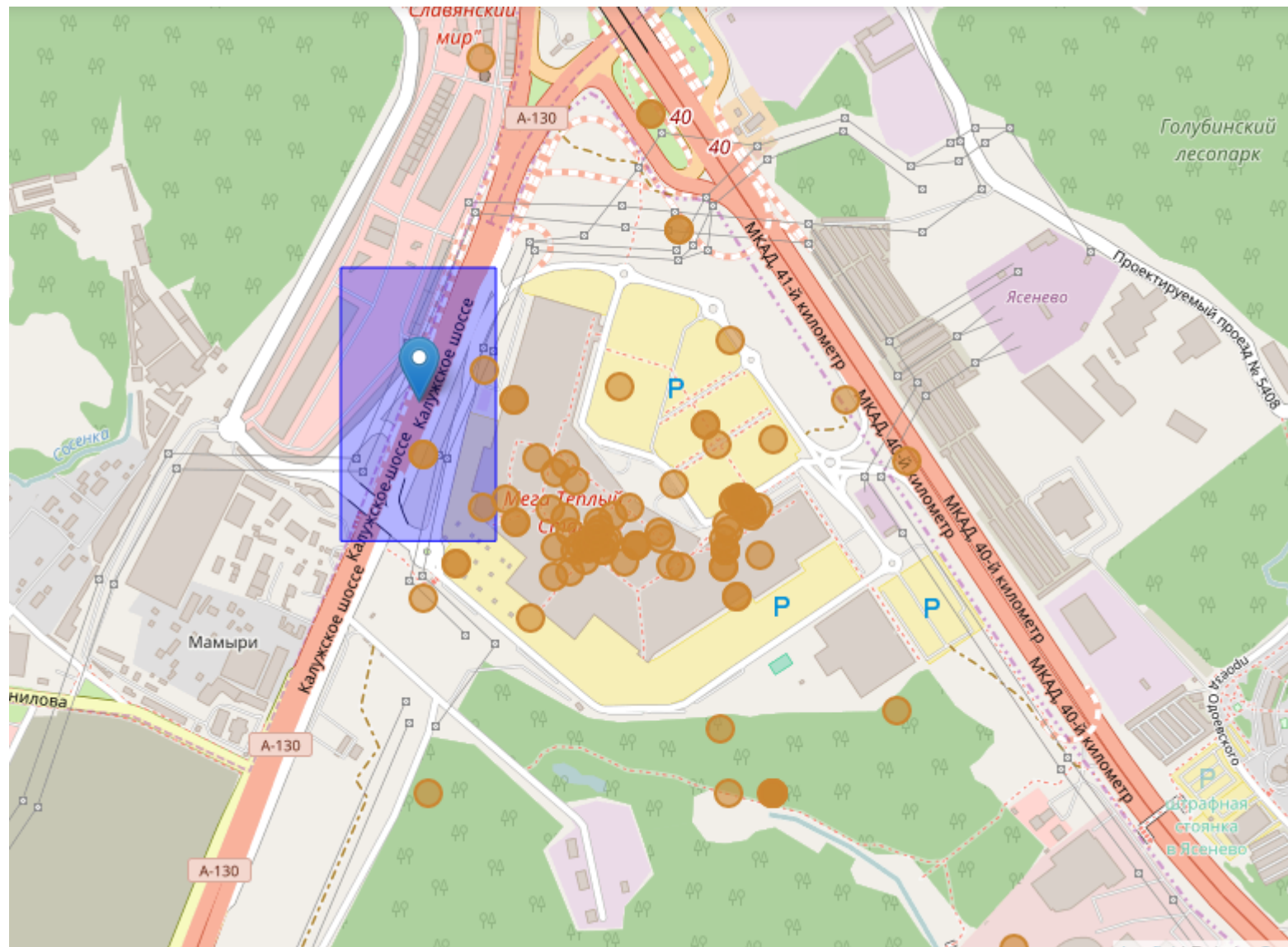
Если на кросс-валидации проверить точность для прямоугольника 0.0025, то она возрастает на 0.02. То есть часто мы не угадываем совсем немного.

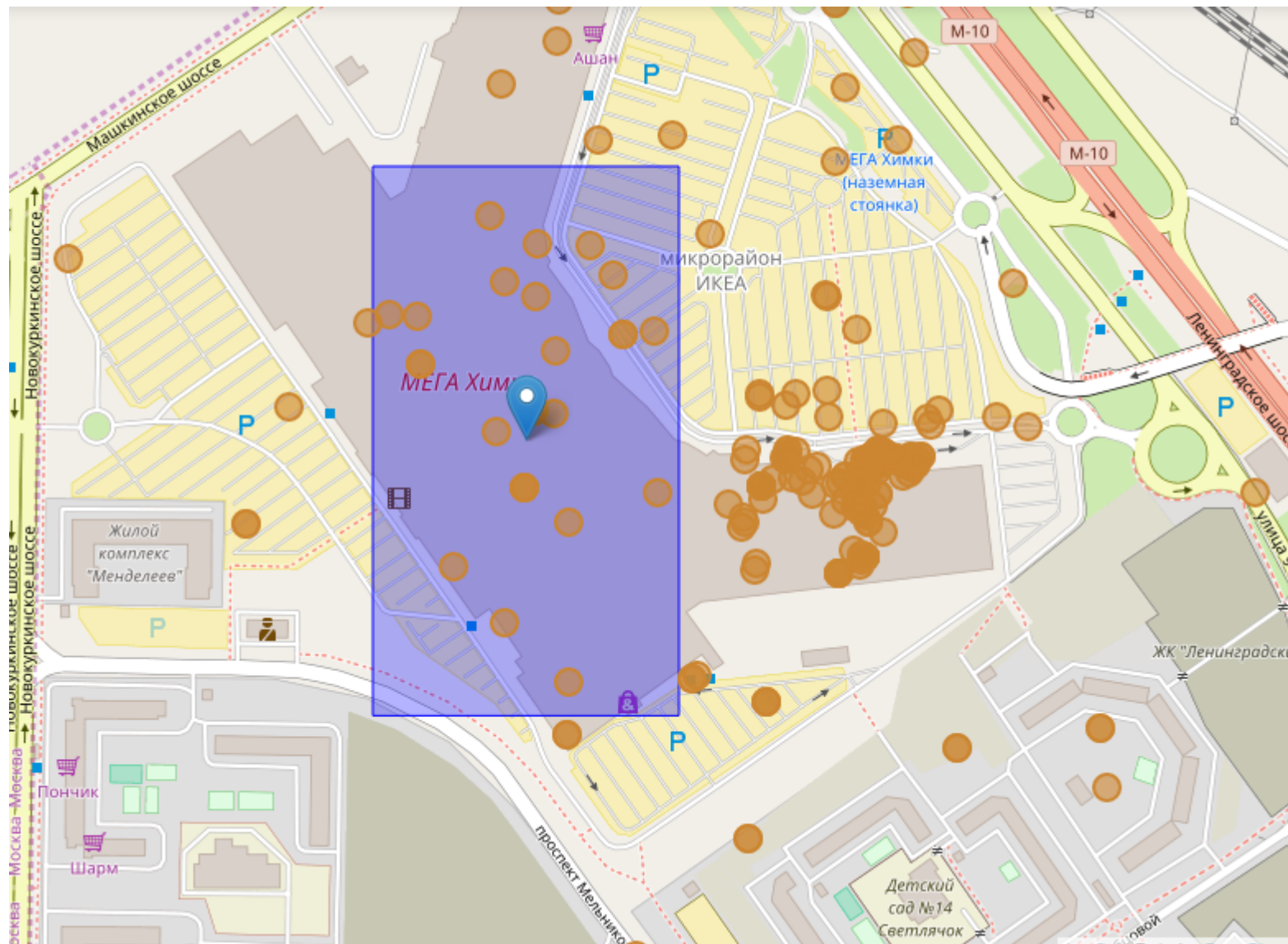


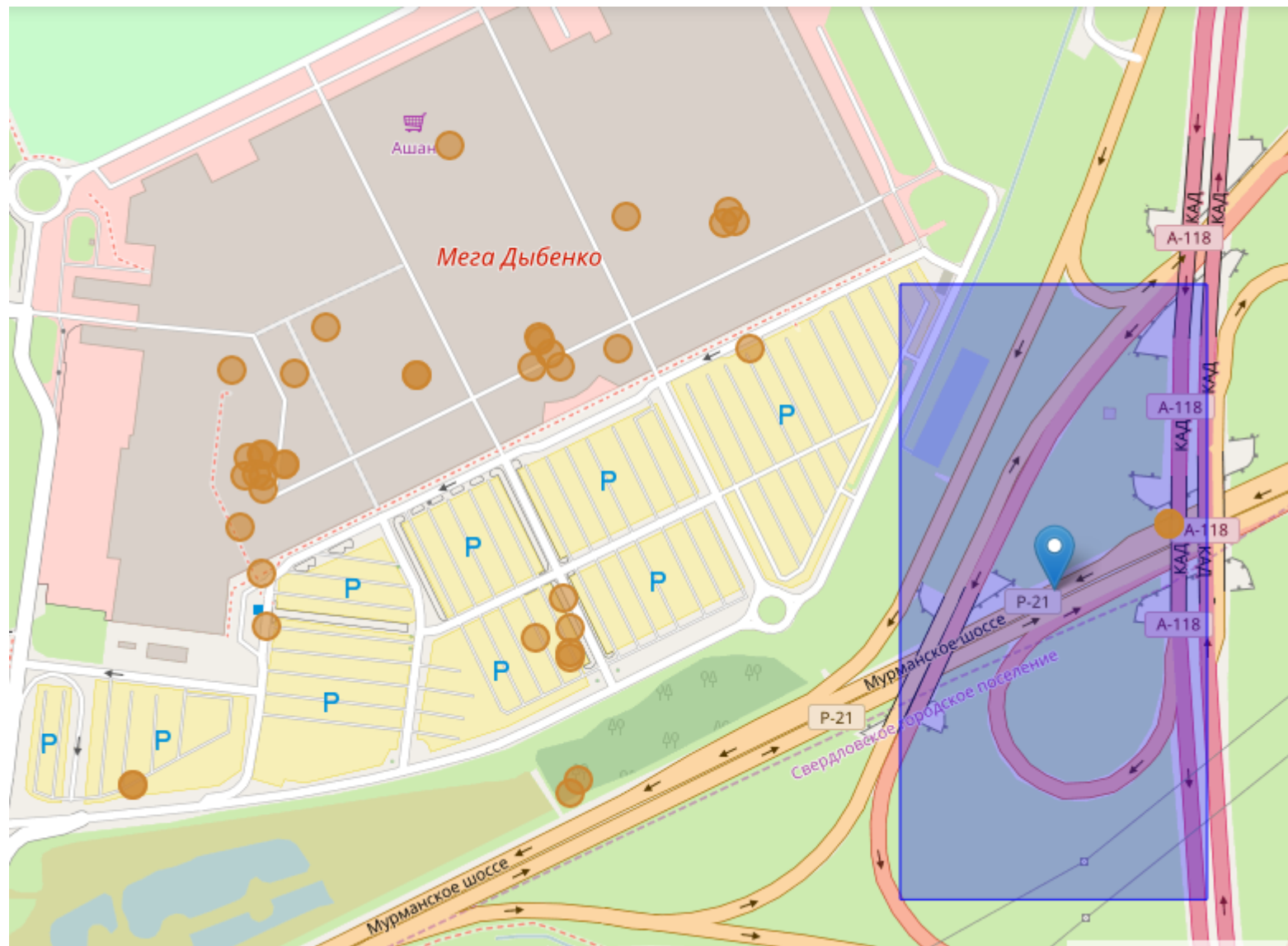












Смещение предсказания

Добавляем вручную три последних.

LB score **+0.002**

Смещение предсказания

Делаем смещение к известным мерчантам, чтобы в прямоугольнике было наибольшее количество мерчантов.

Смещение не более чем на 0.001.

LB score **0.394**

Что можно попробовать

- регрессия, например, расстояние от чекина до мерчанта
- ранжирование
- смещение с учетом магазинов, больших ТЦ

Ссылки

- https://boosters.pro/champ_3 страница соревнования
- <https://github.com/CapitanKK/tinkoff-challenge-2> мое решение

Вопросы