

Kaggle Facebook V: Predicting Check Ins



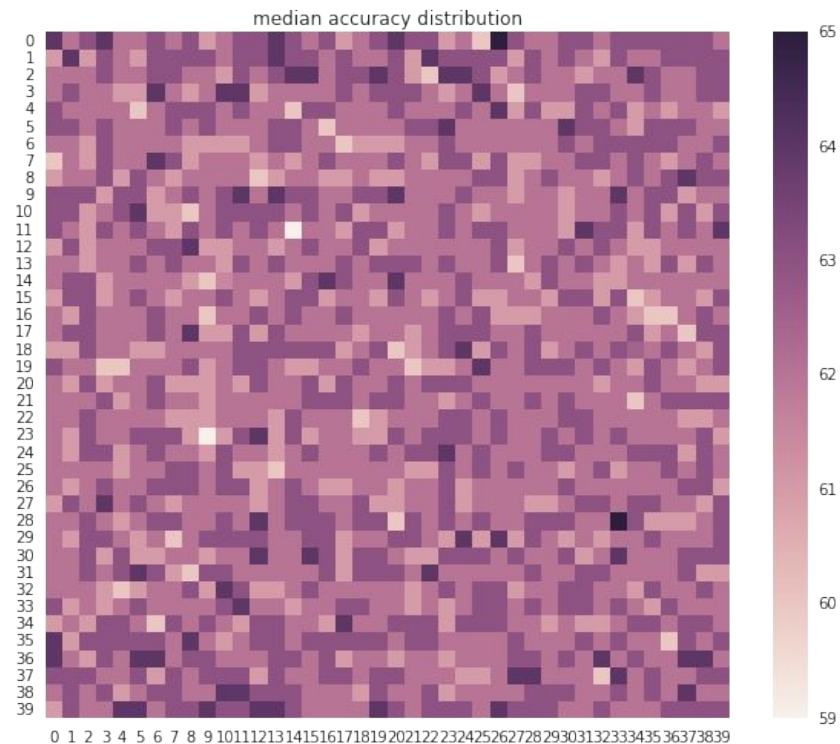
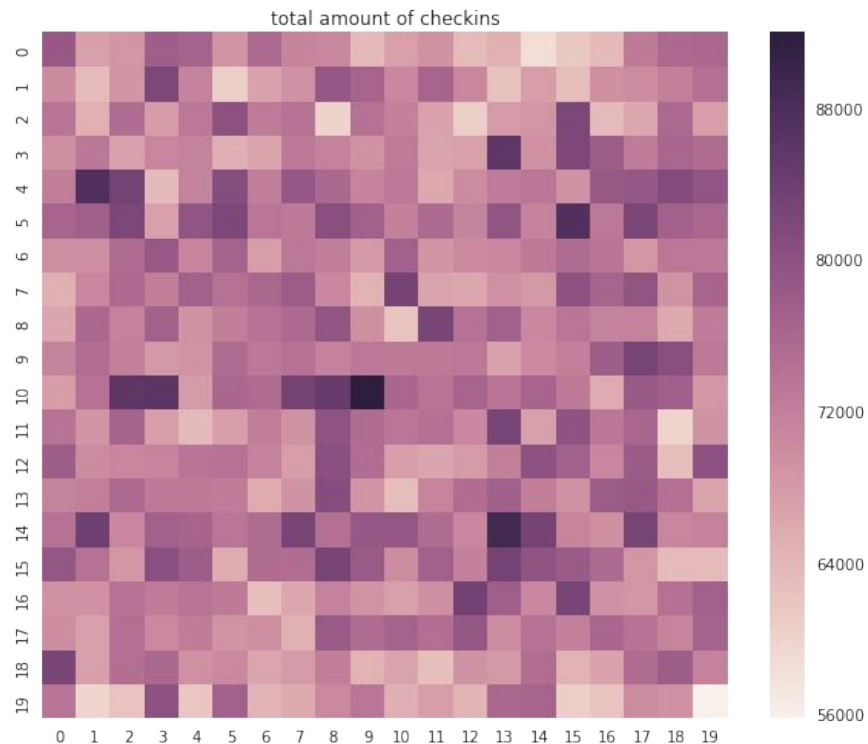
Першин Михаил

2016

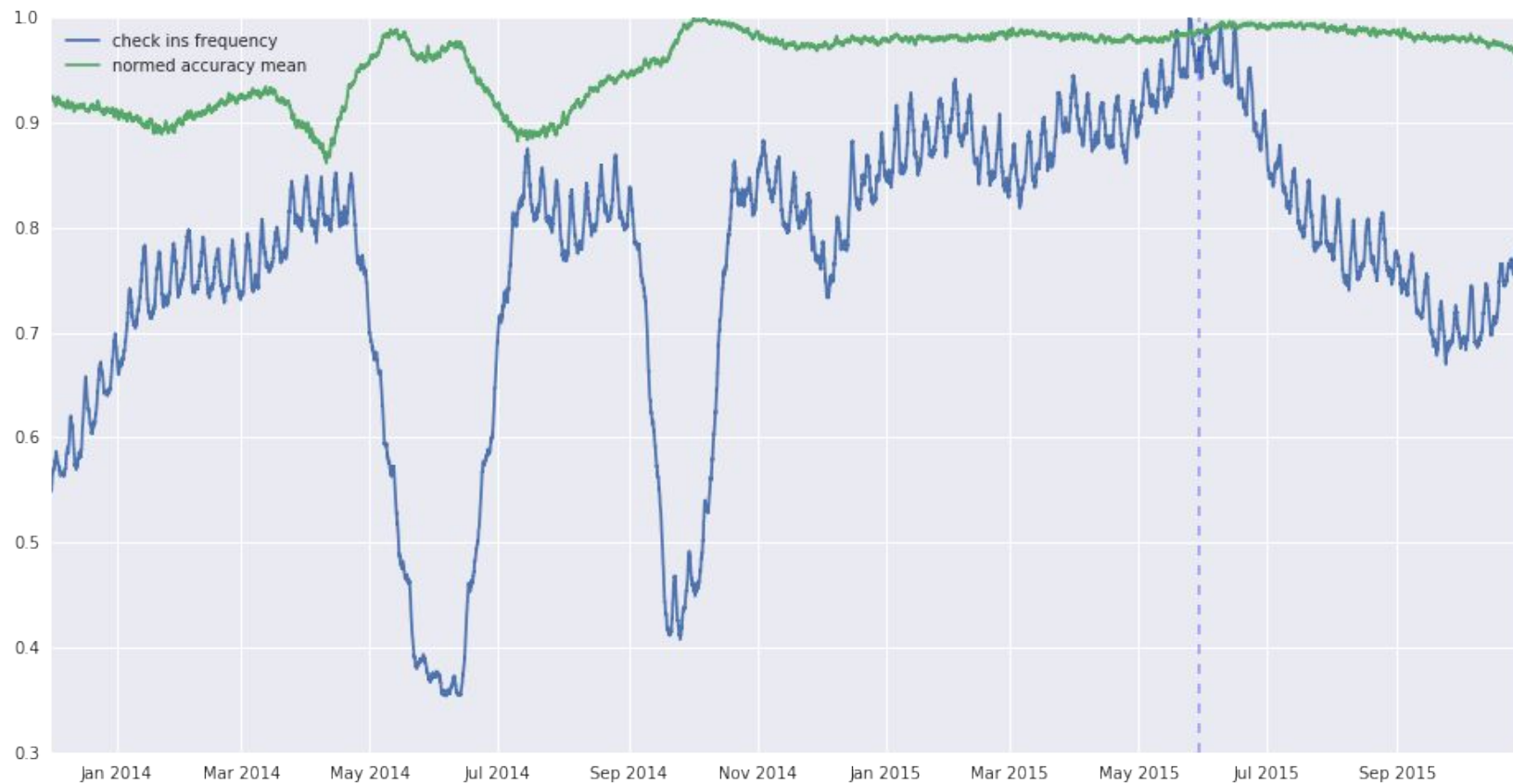
Задача

- Предсказать в каком месте пользователь совершил чек ин, основываясь на координатах, времени и точности измерения
- ~ 29 миллионов обучающих примеров, ~8 миллионов тестовых
- ~ 100 тысяч уникальных мест
- Метрика: Mean Average Precision @3 (MAP@3)

Данные



Данные



Фичи

- Координаты простые и с ними много не придумаешь
- Точность непонятная
- Время в сыром виде бесполезно - нужно выделять фичи для учёта периодичности событий

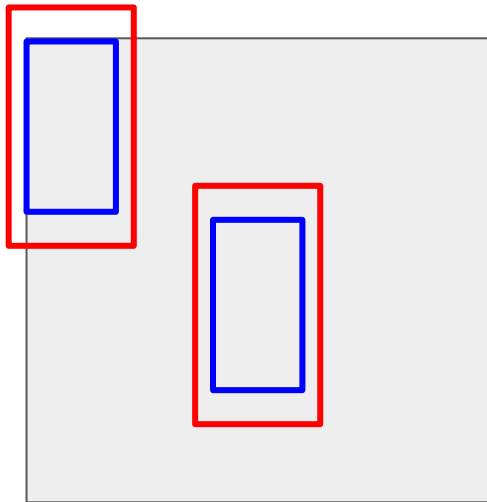
Модель 1

- Что у меня работало плохо: Naive Bayes
- Что работало неплохо: RF, KNN
- Что работало хорошо: xgboost

Модель 1

Если пытаться кормить модель данными “в лоб” - требования к ресурсам становятся абсурдными.

Что делаем: вспоминаем про природу данных и разбиваем их сеткой



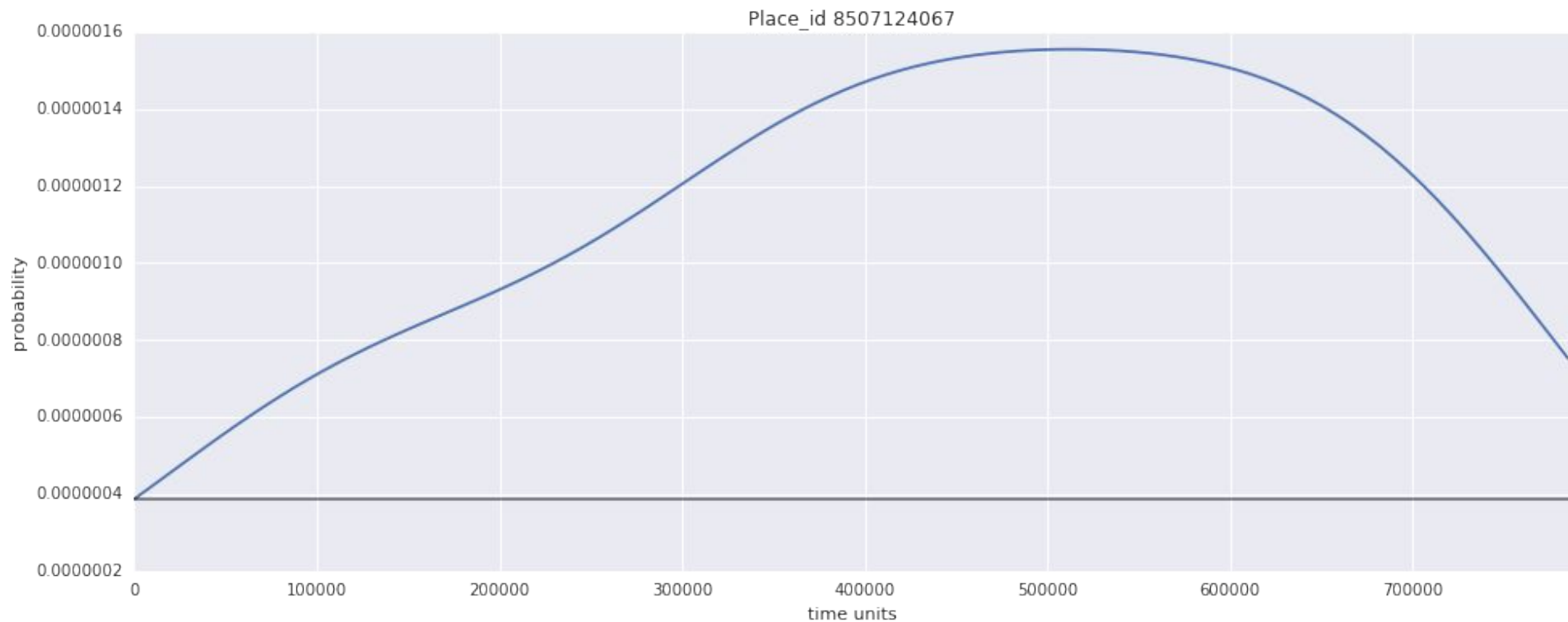
Модель 1

После аккуратного подбора всех возможных гиперпараметров:

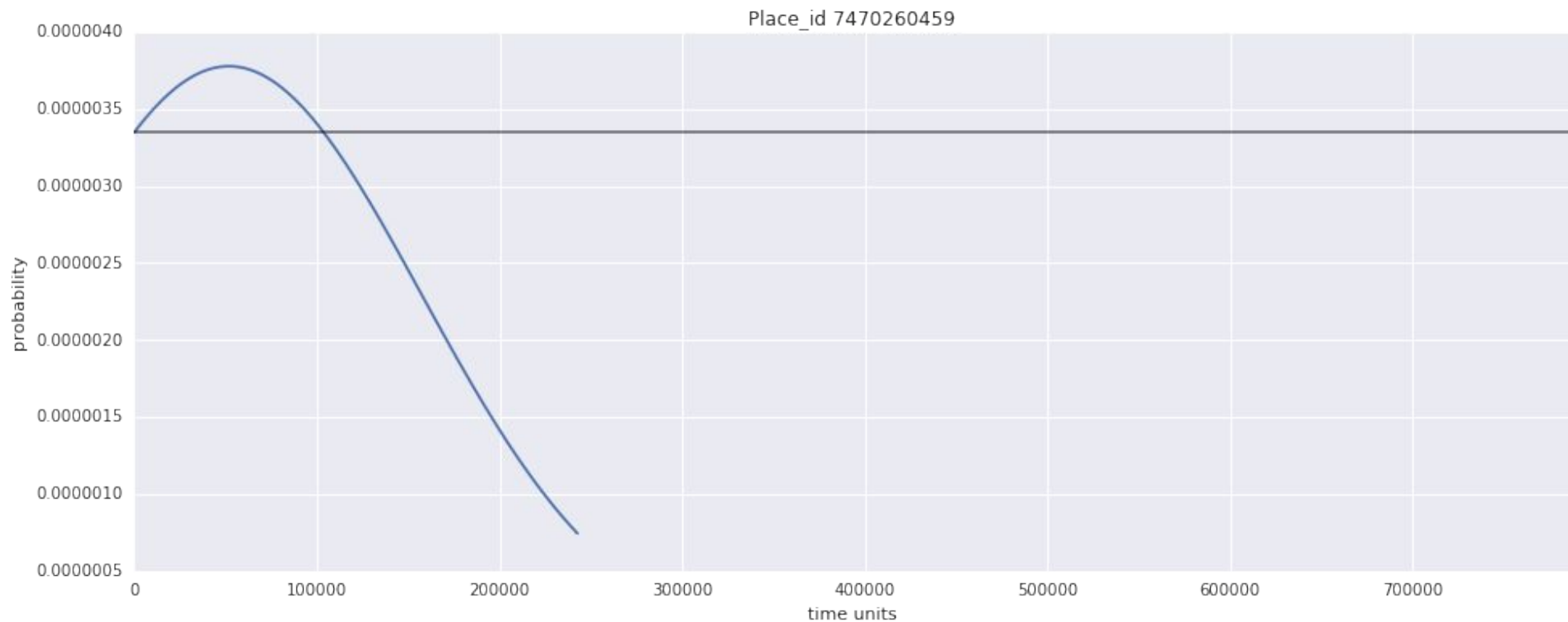
~0.589 MAP@3

~ 86е место на Private Leaderboard-е (топ 7% !)

Киллер-фича: kde



Киллер-фича: kde



Киллер-фича: kde

Зачем останавливаться, если хорошо пошло? Построим распределения для остальных фичей, в том числе и периодических!

Что получим:

~ 0.606 MAP@3

~ 12e место на Private Leaderboard-e

Как же без ансамблей

Ансамблировать что-то внутри модели слишком дорого, попробуем использовать уже посчитанные ранее файлы (а может быть даже утянем публичные)

Используя несколько файлов решений, переранжируем `place_id`: в каждом файле будем в каждой строке назначать каждому предсказанному `place_id` баллы: 3 за первое место, 2 за второе, 1 за третье. Суммируем по файлам, и делаем коллективное предсказание.

В итоге: +0.0025, 8е место.

(на модели 1 вплоть до +0.0075 для ансамбля из 10+ решений)

Спасибо за внимание!

Код появится тут: https://github.com/pershinmr/kaggle_facebook_v