

ПРОГНОЗИРОВАНИЕ ОТТОКА ПОЛЬЗОВАТЕЛЕЙ

- » Данные разного типа
 - ▶ числовые
 - ▶ номинальные
 - ▶ порядковые
- » Временные ряды

НЕСБАЛАНСИРОВАННАЯ ВЫБОРКА

- › Доля целевого класса может быть намного меньше доли нецелевого класса (0.1% vs 99.9%)
- › Несбалансированность выборки может негативно сказаться на качестве модели
- › Важно заметить это в процессе построения модели!

- Задать веса для объектов таким образом, чтобы:
 - ▶ Скомпенсировать количество объектов меньшего класса их важностью
 - ▶ Задать стоимость ошибки классификации разного рода

- Сгенерировать больше объектов меньшего класса:
 - ▶ Дублирование объектов
 - ▶ Генерация новых объектов путем изменения некоторых признаков существующих объектов
 - ▶ Генерация новых объектов на основе нескольких существующих объектов

- Исключить из обучения объекты преобладающего класса:
 - ▶ Удаление из выборки случайных объектов преобладающего класса
 - ▶ Удаление из выборки групп схожих объектов из преобладающего класса

- › Обучение на данных, доступных НЕ только за исторический период
- › Контроль обучения на данных из будущего
- › Контроль переобучения

КРОСС-ВАЛИДАЦИЯ



- › По объектам
- › По времени

ПОДБОР ПАРАМЕТРОВ

- › Используем кросс-валидацию
- › Сразу фиксируем hold-out dataset
- › Их может быть несколько
- › Используем для финальной проверки решения

- › Одна целевая метрика
- › Оффлайн метрика совпадает или коррелирует с целевой метрикой
- › Хорошо оценивать модель «скользящим окном» по времени

- На каких группах объектов модель ошибается?
- Является ли инвестиция в дальнейшее улучшение модели экономически оправданной?

- › Какие факторы внесли наибольший вклад в модель?
- › Гипотезы относительно причин оттока пользователей?
- › Какие объекты классифицируются наиболее/наименее уверенно?
- › Какие еще данные могли бы быть полезны?

- › Как меняется качество модели во времени?
- › Как быстро она «протухает»?
- › Сколько времени занимает переобучение модели?
- › Сколько времени требуется на переключение с одной модели на другую?

- › Изменились ли данные?
- › Изменилось ли качество модели?
- › Хорошо оценивать модель с разных сторон с помощью набора метрик

- › Оценивать качество решения задачи на всех этапах
- › Заранее продумать список потенциальных «узких» мест