

ЦЕНТР
АНАЛИЗА
ДАННЫХ



Feature Engineering, Ensembles and Medicine Cases

19.11.16

Data analysis center

<https://analysiscenter.ru/>

<https://analysiscenter.github.io/about/>

<https://www.facebook.com/analysiscenter.ru>



ДЕПАРТАМЕНТ
ЗДРАВООХРАНЕНИЯ
ГОРОДА МОСКВЫ



Департамент
информационных
технологий
города Москвы





Feature Engineering – Алексей Ушаков

Ensembles – Кирилл Емельянов



Case: Прогнозирование спроса на терапевтическую помощь Дмитрий Подвязников

Case: Балансировка терапевтических участков
Наталия Амелина



Feature Engineering – Алексей Ушаков



Чек-лист аналитика:

- Сформулировать корректную проблему
- Выбрать модель
- Выбрать алгоритм обучения
- Найти данные
- **Придумать/создать признаки, правильно представить, отобрать релевантные**

Feature (признак) – информация, потенциально полезная для предсказания

Мы не можем формально определить, что такое Feature Engineering, но... для успешной реализации нужно:

- Знать свои данные (область, проблематику, особенности)
- Понимать особенности задачи (регрессия/классификация, метрика и т.д.)
- Очень много экспериментировать



Но всему есть предел:

1. Рано или поздно скорость станет проблемой
2. Overfitting

Поэтому нужно выбирать только «полезные», «релевантные» признаки

Как научиться:

- Проводить больше экспериментов
- Читать статьи/скрипты под конкретные проблемы
- Придумывать задачу, «а как бы я сделал?»

О чем подумать, работая с данными:

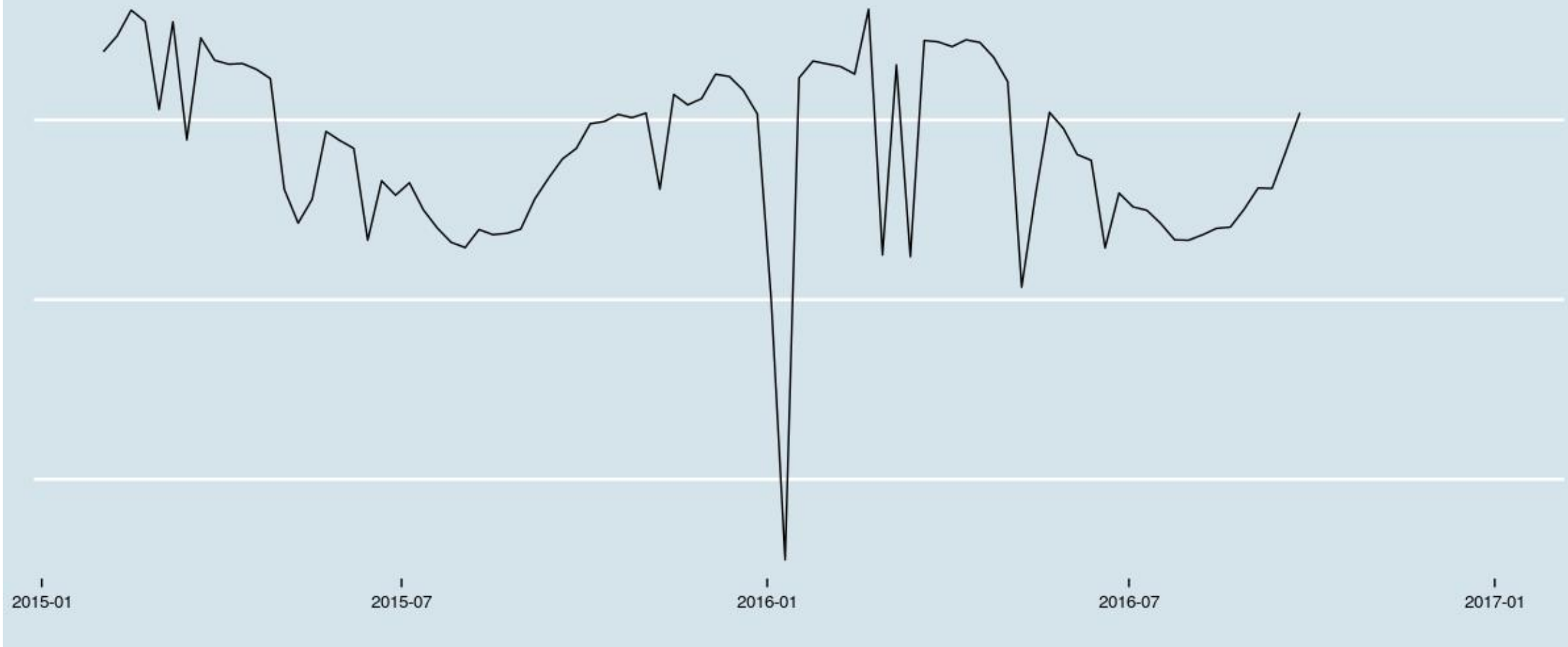
- Есть ли N/A, Missing values?
- Какое распределение у признака? Насколько сильно оно зависит от подвыборки? От времени/типа сегментации?
- Есть ли дубликаты или пересекающиеся данные?
- Все ли данные дискретны/непрерывны?

Хорошая статья:

“A Few Useful Things to Know about Machine Learning” Pedro Domingos

Спрос на услуги терапевтов в Москве

Какие признаки нужны для такой задачи?



Ensembles — Кирилл Емельянов



- Что такое ансамбли и зачем их применять
- Примеры ансамблей
- Random forest
- Boosting
- Stacking

- Обучающая выборка
- Вероятностная природа данных
- Метод обучения
где A --- пространство алгоритмов.
- Функция потерь
- Функционал потерь

$$X^l = (x_i, y_i)_{i=1}^l$$

$$(x_i, y_i) \sim p(x, y)$$

$$\mu : X^l \rightarrow A \quad ,$$

$$L(z, y)$$

$$Q(\mu) = E_{X^l} E_{(x,y)} (\mu(X^l)(x) - y)^2$$

Основные подходы к построению композиции алгоритмов:

1) Простое голосование:

$$\hat{b}(x) = F(b_1(x), \dots, b_T(x)) = \frac{1}{T} \sum_{t=1}^T b_t(x)$$

2) Взвешенное голосование:

$$\hat{b}(x) = F(b_1(x), \dots, b_T(x)) = \sum_{t=1}^T \alpha_t \cdot b_t(x), \quad x \in X, \alpha_t \in \mathbf{R}$$

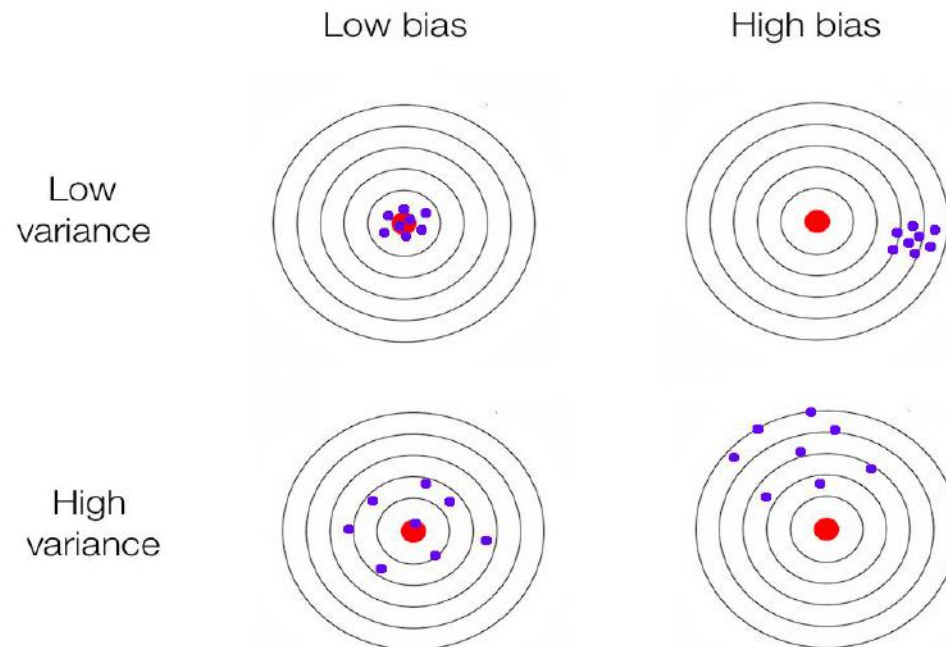
3) Смесь алгоритмов:

$$\hat{b}(x) = F(b_1(x), \dots, b_T(x)) = \sum_{t=1}^T g_t(x) \cdot b_t(x), \quad g_t(x) : X \rightarrow \mathbf{R}$$

Теорема:

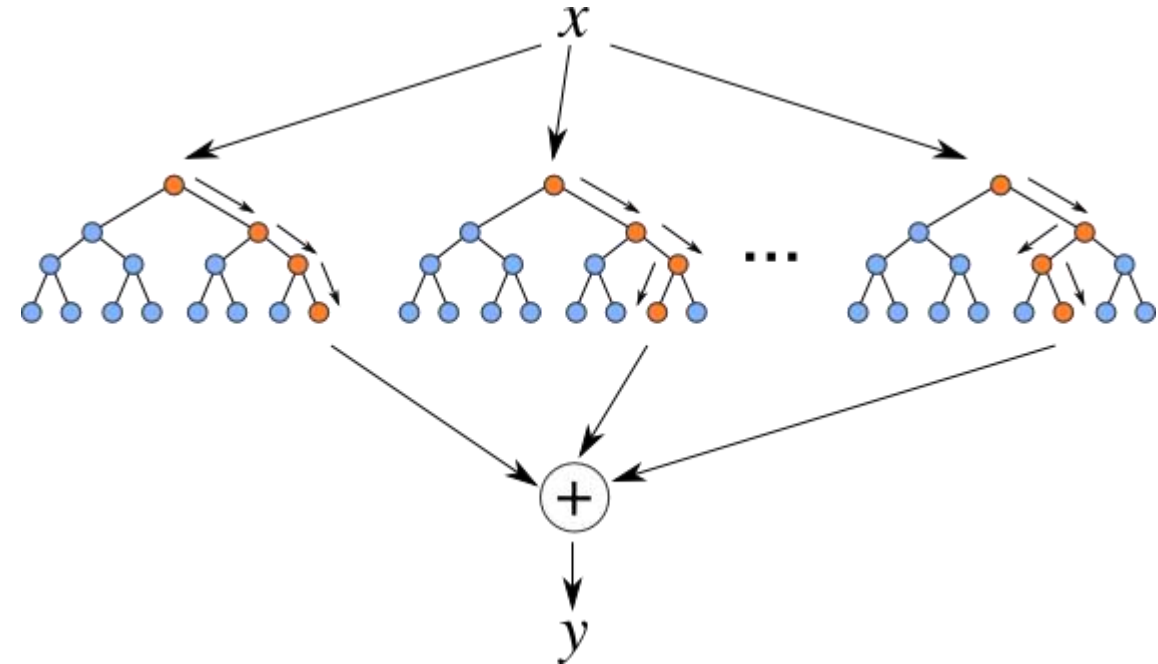
$$Q(\mu) = \underbrace{E_{(x,y)} (a^*(x) - y)^2}_{noise} + \underbrace{E_{(x,y)} (\bar{a}(x) - a^*(x))^2}_{bias} + \underbrace{E_{(x,y)} E_{X^l} (\mu(X^l)(x) - \bar{a}(x))^2}_{variance}$$

где $\bar{a}(x) = E_{X^l} (\mu(X^l)(x))$, $a^*(x) = E(y|x) = \int_Y y p(y|x) dx$



- Построение такой композиции уменьшает разброс (variance).
- Качество увеличивается, если алгоритмы композиции преимущественно независимы.
- Bagging (Bootstrap aggregation):
из обучающей выборки формируются различные обучающие подвыборки меньшего размера с помощью бутстрепа (выбор с возвращением).
- RSM (Random Subspace Method):
базовые алгоритмы обучаются на различных подмножествах признаков, выделяемых случайным образом.

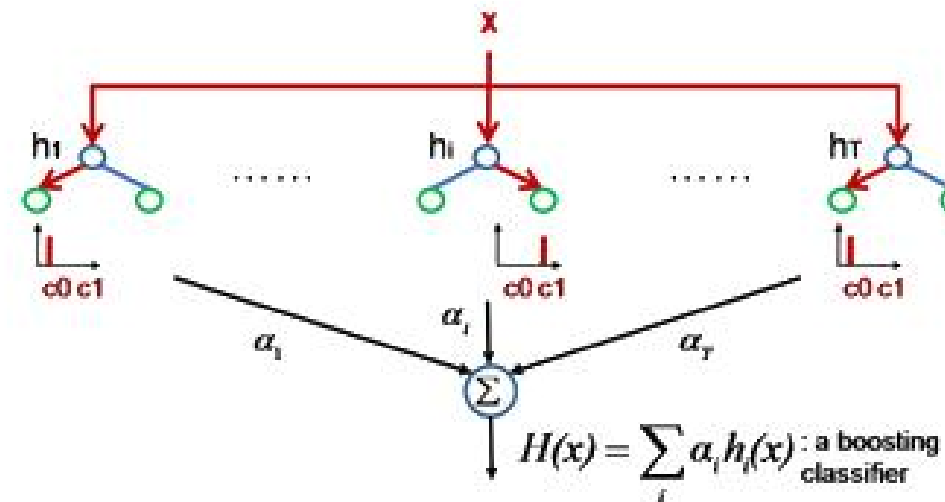
- Бэггинг над решающими деревьями.
- RSM. Признак в каждой вершине дерева выбирается из случайного подмножества признаков (k из n).
- Оптимальные значения для параметра k (эвристики):
 - $k = \lfloor \frac{n}{3} \rfloor$ для задачи регрессии
 - $k = \lfloor \sqrt{n} \rfloor$ для задачи классификации



- Комбинация слабых классификаторов, которая сама является сильным классификатором.
- Уменьшает смещение.
- Обычно в качестве слабых классификаторов выступают решающие деревья небольшой глубины (~ 5).

Примеры функций потерь:

- AdaBoost: $L(z) = e^{-z}$
- LogitBoost: $L(z) = \log_2(1 + e^{-z})$
- GentleBoost: $L(z) = (1 - z)^2$



- Произвольная дифференцируемая функция потерь.
- На каждой итерации добавляется алгоритм, который сильнее всех остальных уменьшает ошибку композиции.

- Формально, на каждой итерации алгоритма решается следующая задача:

$$\sum_{i=1}^l L(y_i, a_{N-1}(x_i) + \gamma_N b_N(x_i)) \rightarrow \min_{b_N \in B}$$

- Ответы алгоритма должны наилучшим образом приближать градиент функции потерь:

$$b_N(x) = \operatorname{argmin}_{b \in B} \sum_{i=1}^l (b(x_i) + \frac{\partial L}{\partial z} \big|_{z=a_{N-1}(x_i)})^2$$

- Вес алгоритма выбирается по аналогии с наискорейшим спуском:

$$\sum_{i=1}^l L(y_i, a_{N-1}(x_i) + \gamma_N b_N(x_i)) \rightarrow \min_{b_N \in B}$$

$$\gamma_N = \operatorname{argmin}_{\gamma \in \mathbf{R}} \sum_{i=1}^l L(y_i, a_{N-1}(x_i) + \gamma b_N(x_i))$$



- Бустинг дает лучшие результаты на обучающих выборках большого размера.
- Бэггинг и RSM более предпочтительны для коротких обучающих выборок.
- RSM позволяет сократить число признаков в случае, когда из больше чем объектов, или среди них есть много неинформативных.
- Для бэггинга возможна параллельная реализация.

При наличии построенных t различных базовых алгоритмов, вектор их предсказаний $(b_1(x), \dots, b_t(x))$ можно принять за новое признаковое описание объекта, после чего построить по ним новый мета-алгоритм a :

$$\sum_{i=1}^l L(y_i, a(b_1(x_i), \dots, b_N(x_i))) \rightarrow \min_a$$

- Обучение всех базовых алгоритмов и мета-алгоритма на одной обучающей выборке приводит к переобучению.

- Исходную выборку разбивают на K блоков

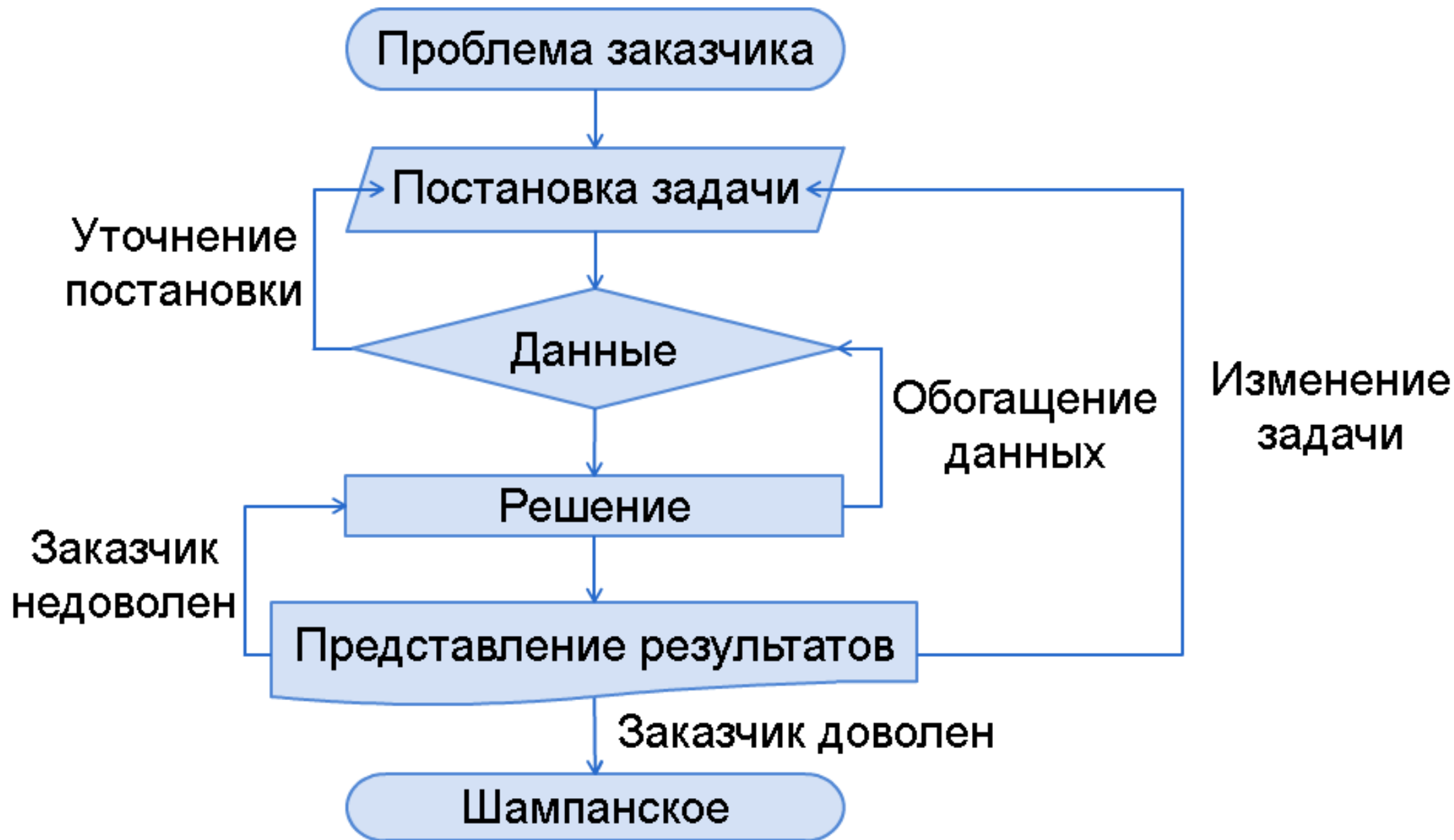
$$X = X_1 \sqcup \dots \sqcup X_k$$

- Базовые алгоритмы обучают на исходной выборке, из которой исключен один из блоков.

- Частный случай стэкинга
- Мета-алгоритм является линейной комбинацией базовых:
$$a(x) = \sum_{i=1}^k w_i b_i(x), \quad \sum_{i=1}^k w_i = 1$$
- Для обучения весов можно использовать различные линейные методы машинного обучения:
 - Логистическая регрессия
 - Линейный дискриминант Фишера
 - Метод опорных векторов (SVM)
- Если базовый входит в мета-алгоритм с отрицательным весом, то его нужно исключить из композиции

Case: Прогнозирование спроса на терапевтическую помощь Дмитрий Подвязников





Когда отпускать работников в отпуск?

Когда нанимать новых работников?

Когда отправлять работников на обучение?

Сколько заказать вакцин для прививок?

Как составить расписание работы?

Удостовериться, что всех пациентов приняли?

Как сделать так, чтобы не было очередей и все могли записаться?

Что нужно?

Прогноз числа посещений в неделю, на 12 недель вперед

Когда нужно?

Через 2 недели

...

	ID	CLINIC_ID	SUNDAY_DATE	RECEIPTS	CHRONICS	PATIENTS	VISITS	DISTRICT
0	0	22	2015-02-01	1376.0	2561.0	52307.0	1491.0	1
1	1	22	2015-02-08	1376.0	2561.0	52307.0	1713.0	1
2	2	22	2015-02-15	1376.0	2561.0	52307.0	1725.0	1
3	3	22	2015-02-22	1376.0	2561.0	52307.0	1637.0	1
4	4	22	2015-03-01	1384.0	2570.0	48437.0	1431.0	1
5	5	22	2015-03-08	1384.0	2570.0	48437.0	1661.0	1
6	6	22	2015-03-15	1384.0	2570.0	48437.0	1355.0	1
7	7	22	2015-03-22	1384.0	2570.0	48437.0	1615.0	1
8	8	22	2015-03-29	1384.0	2570.0	48437.0	1567.0	1
9	9	22	2015-04-05	1443.0	2563.0	48845.0	1485.0	1



**KEEP
CALM
AND
DO YOUR
WORK**

Рассказать про данные

Какие важные закономерности удалось найти?

Какое новое знание удалось выявить?

Объяснить вашу модель

Почему такой подход?

Почему выбрали / добавили эти фичи?

Обосновать решения

Чем обусловлен выбор метрики?

Почему такое деление на обучающую / тестовую выборки?

Case: Балансировка терапевтических участков
Наталия Амелина



Для уменьшения очередей к терапевтам, отменили принцип «участкового» врача. Пациенты могут записаться к любому доступному в поликлинике врачу



КТО теперь отвечает за пациента?

КАК оценить работу конкретного врача?

Участки не несут сейчас смысла и создают **РАЗНУЮ** нагрузку



Задача **восстановить участковую службу**:

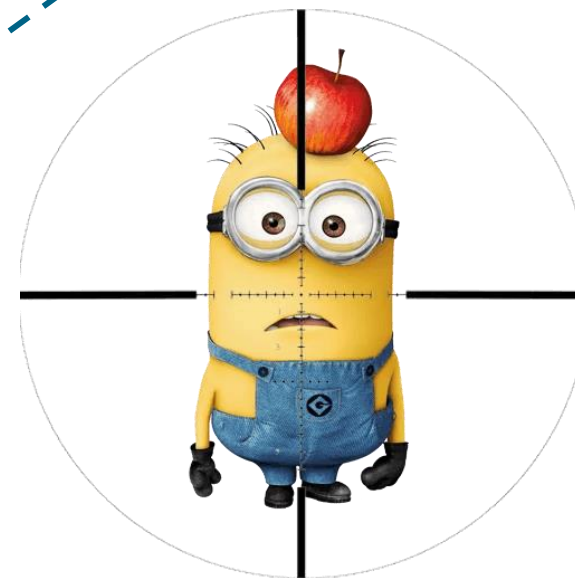
- ✓ Сопоставимая нагрузка (количество приемов) с одного участка в период
- ✓ Нагрузка остается стабильной через N лет
- ✓ Минимальное количество людей меняет участок
- ✓ Не трогать тех, кто хорошо знает своего врача (untouchable)
- ✓ * Восстановить семьи
- ✓ * Выделить хроников в отдельные участки

Как меняется с **возрастом** вероятность визита к врачу?

Лишь 34% застрахованных пришли к врачу за исследуемый период. Как работать с **неполными данными?**

В неидеальном мире нужно неидеальное решение: **когда** остановить балансировку?

1. Скоринг человека



2. Балансировка участков

Признак	Описание
untouchable	признак пациента, которого нельзя переносить с участка
pat_id	идентификатор пациента
lpu_district	район ЛПУ (поликлиники)
pers_district	район проживания пациента
MOSCOW_REGISTRATION	признак наличия московской регистрации
invalid	признак наличия инвалидности
nadomnik	признак неходячего пациента
age_calculate_y	возраст
sovmetst_deti	количество совместно проживающих детей до 14 лет
sovmetst_lgotnik	количество совместно проживающих льготников
sovmetst_pens	количество совместно проживающих пенсионеров
sphere_new_Учреждения образования	работник образования
sphere_new_Учреждения здравоохранения	работник здравоохранения
sphere_new_Учреждения соцзащиты	работник соц.защиты
sphere_ZHKN_ZKKH	работник ЖКХ
bezrab_безработный сейчас	безработный на 01.07.2016
bezrab_был безработный 2015/16	был безработным в период 01.01.2016 – 30.06.2016
bezrab_дело на рассмотрении	потенциально станет безработным
COMMON_GENDER	пол
lux1	признак материального благосостояния № 1
lux2	признак материального благосостояния № 2
FLAG_BESSR	признак бессрочного льготника
HAD_APPOINTMENTS	появлялся ли пациент хотя бы один раз за последние 2 года
from_region	региональный полис
lpu_id	идентификатор ЛПУ (поликлиники)
district_id	идентификатор участка
target	целевая переменная (количество приемов) – преобразованные

Презентация «**как заказчику = врачу**» + рабочий ноутбук до 1 декабря

Представление результатов и обсуждение 3 декабря

1. Pre-research

Что интересного есть в данных?

2. Скоринговая модель человека

Обоснование выбора модели, метрики точности

Интерпретируемая модель! Не черный ящик, а набор правил понятный врачу

3. Балансировка участков на скоринговой модели на 3 ЛПУ (13, 17, 19)

Не трогать untouchable

Обосновать принцип / функционал балансировки, метрика точности

Почему с этим решением участки не разбалансируются в будущем?

В каком моменте алгоритм останавливается — квазиоптимальное решение?



Feature Engineering – Алексей Ушаков
a.usakov@analysiscenter.ru



Ensembles – Кирилл Емельянов
k.emelyanov@analysiscenter.ru



Case: Прогнозирование спроса на терапевтическую помощь Дмитрий Подвязников, d.podvyaznikov@analysiscenter.ru

Case: Балансировка терапевтических участков
Наталия Амелина, n.amelina@analysiscenter.ru
amelinans@gmail.com

