

Прогноз спроса и прогноз пользовательских метрик

План

1. О задачах на сегодня

Задачи

**Прогноз
пользовательских
метрик**

Прогноз спроса

План

1. О задачах на сегодня

2. Метрики

Вы не можете управлять тем, что не измеряете



Метрики

**Метрики обучения
моделей**

Функции, которые оптимизируем математическими методами, для решения какой-то задачи

Метрики

**Метрики качества
моделей**

Метрики для сравнения моделей между собой

**Метрики обучения
моделей**

Функции, которые оптимизируем математическими методами, для решения какой-то задачи

Метрики

Прокси-метрики

Напрямую связаны с бизнес-метриками или влияют на них. Измеряются быстро, в пилотах.

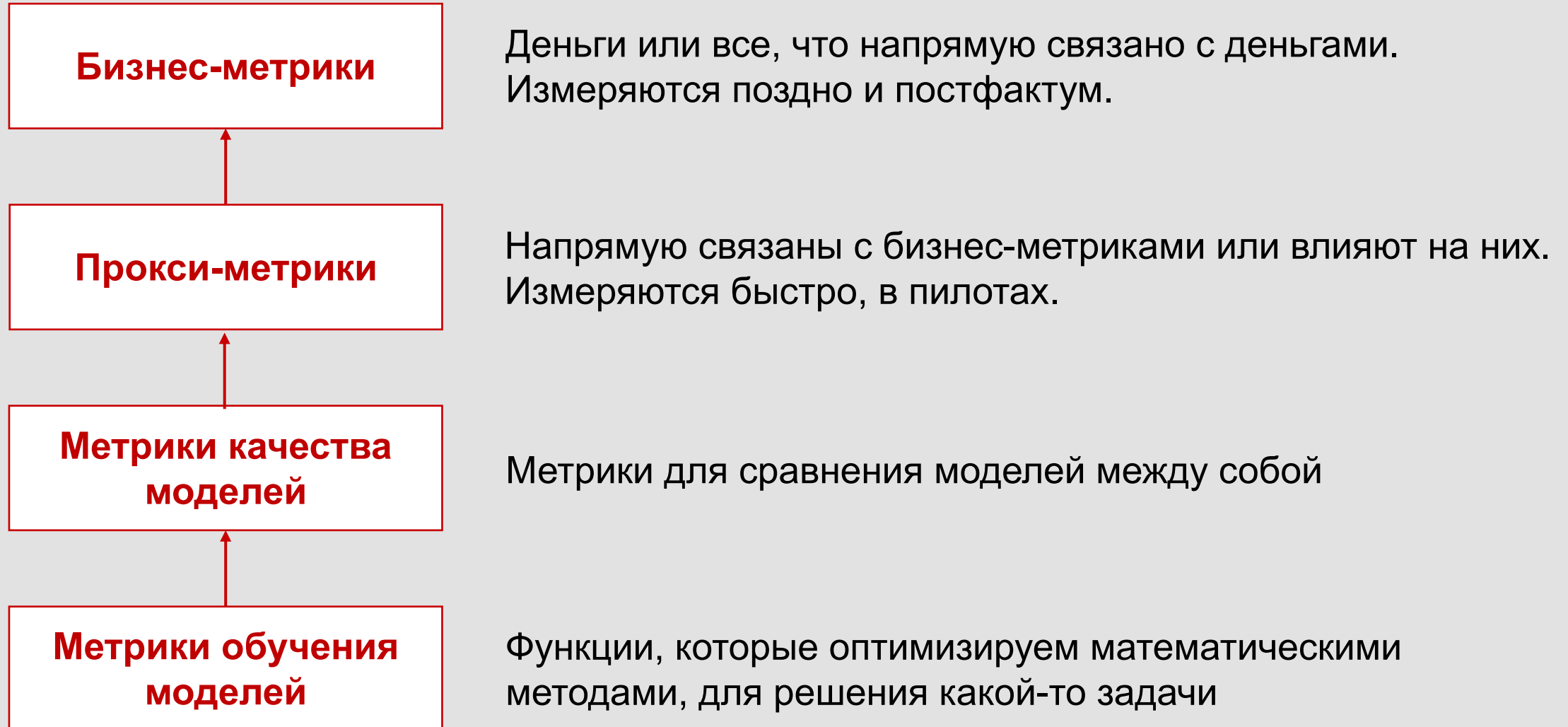
Метрики качества моделей

Метрики для сравнения моделей между собой

Метрики обучения моделей

Функции, которые оптимизируем математическими методами, для решения какой-то задачи

Метрики



Другой подход

Success

Показатели, на которое направлено изменение

Informative

Показатели, хорошо коррелируемые с целевой — могут что-то объяснить

Guardrail

Показатели, на которые направленно влияет изменение, но не являющиеся целевыми.
Рекомендуется за ними наблюдать, чтобы их не уронить и не объяснять эффект каннибализацией

Другой подход

Success

Средний чек

Informative

Добавление товара в
корзину

Guardrail

Время от входа в корзину до ее
прохождения

Ловушка менеджера

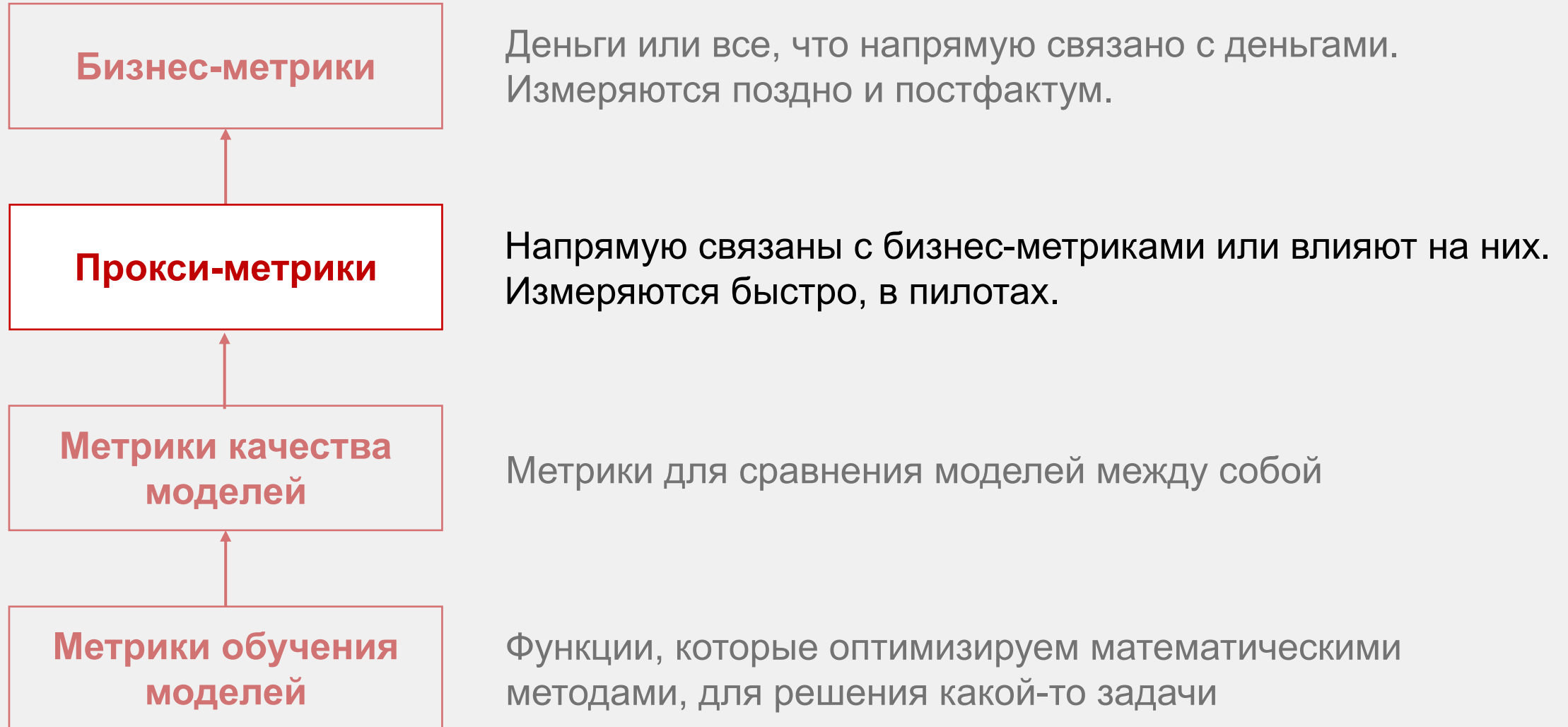
Метрики тщеславия

«Нашим продуктом
пользуются более
миллиона людей»

Метрики качества

В прошлом году
пользовались 2 млн

Метрики



Клиентские метрики в ритейле

Клиентские метрики в ритейле

- Средний PTO на клиента
- Частота покупок клиента
- LTV
customer retention rate
- PTO в период действия кампании
- PTO при совершении целевого действия в период действия кампании
- Кол-во чеков в период действия кампании
- Средний чек в период действия кампании
- доп PTO на человека в кампании
- доп PTO по кампании
- чистый отклик в покупку в период действия кампании
- чистый отклик в целевое действие в период действия кампании
- валовый доход по кампании

План

1. О задачах на сегодня
2. Метрики
- 3. Прогноз клиентских метрик**

**Зачем прогнозировать метрики
пользователей?**

Зачем прогнозировать метрики пользователей?

- опрозрачивание работы и отчеты
- отслеживание жизненного цикла
- маркетинг:
 - выбор механик
 - использование более точных порогов в конкретных механиках

Постановка задачи с т.з. ML

- **Объект x** прогноза
- **Целевая переменная (таргет) y**

Постановка задачи с т.з. ML

- **Объект x** прогноза

клиент программы лояльности (есть уникальный идентификатор)

- **Целевая переменная (таргет) y**

значение метрики в период времени, например:

- количество чеков в календарную неделю
- товарооборот в календарную неделю

Признаки?

Признаки

- Интересуемая величина – временной ряд
- Различные лаги, функции над лагами – основной источник фичей
- Количество покупок за 7/14/21/28 дней
- Товарооборот за 7/14/21/28 дней
- Доля промотоваров за 7/14/21/28 дней
- Количество контактов с пользователем за 7/14/21/28 дней
- ...

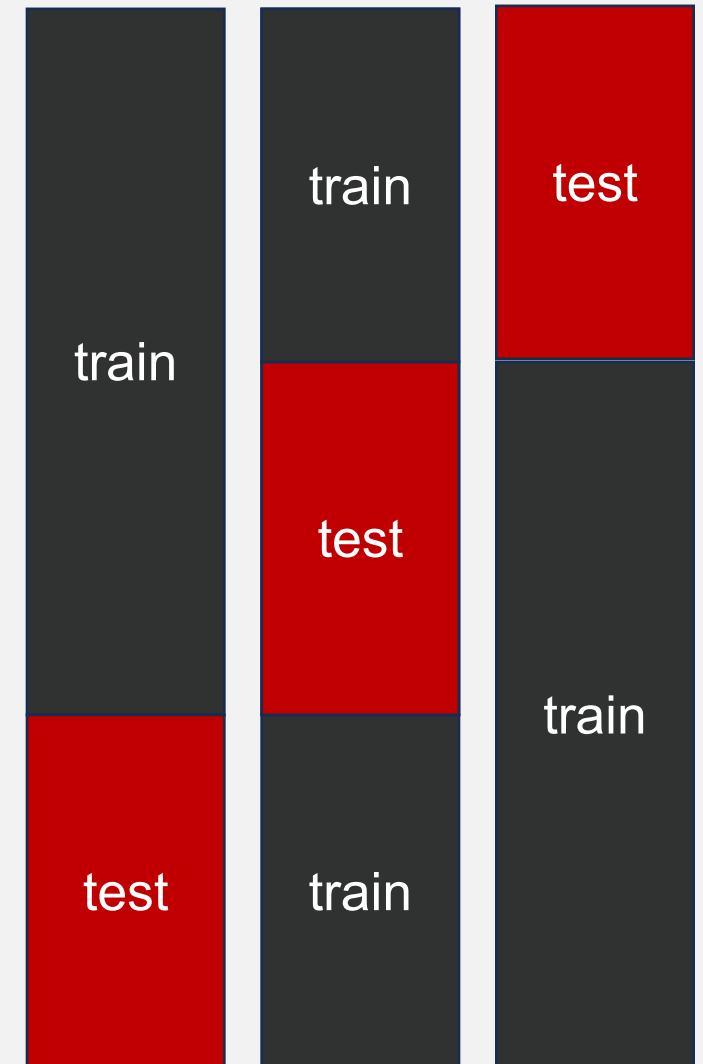
Метрики качества?

Метрики качества?

- В общем случае, все, что уместно для регрессии
- В частном случае, выбор пал на MAE

Валидация алгоритмов

- K-Fold-кросс-валидация
- Разбиваем данные так, чтобы:
 - каждое наблюдение хотя бы 1 раз в тесте
 - каждое наблюдение хотя бы 1 раз в обучении
- k – количество блоков
- Обычно $k = \{3, 5, 7, 10\}$

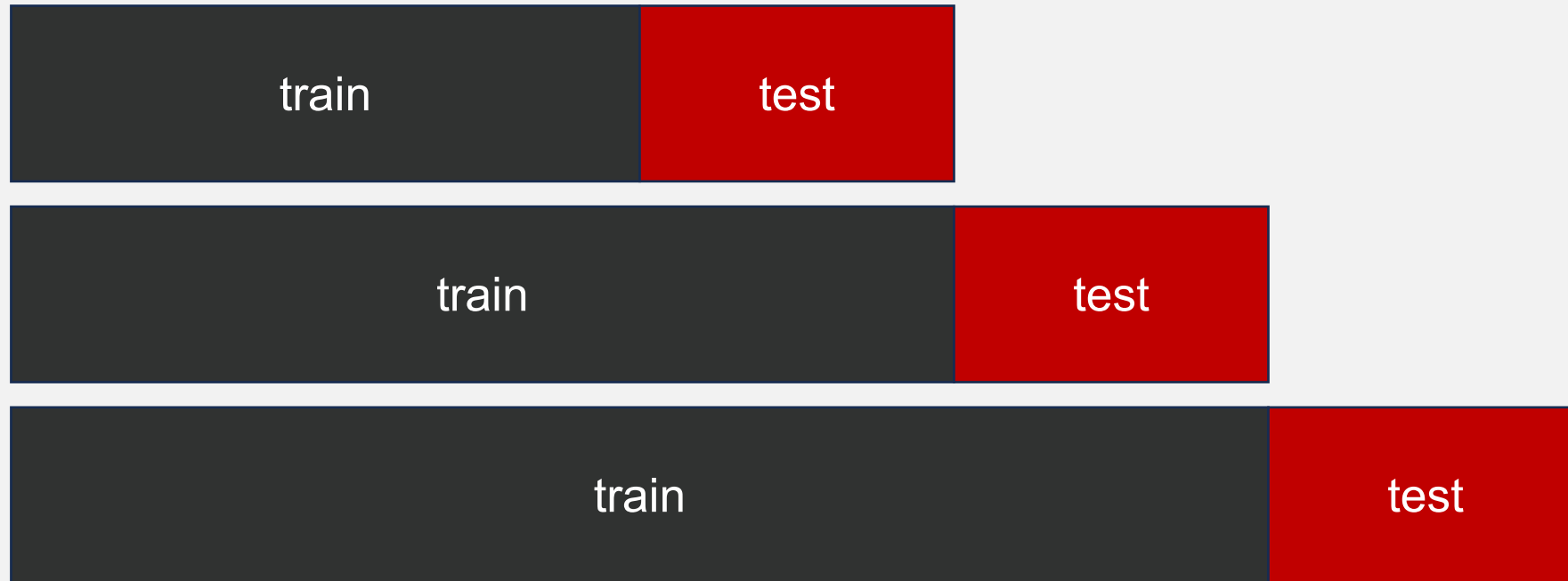


Валидация

- K-Fold-кросс-валидация
- Разбиваем данные на k частей
 - каждое разбиение хотя бы раз в тесте
 - каждое разбиение хотя бы раз в обучении
- k — количество разбиений
- Обычно $k = 10$



Time Series Cross Validation



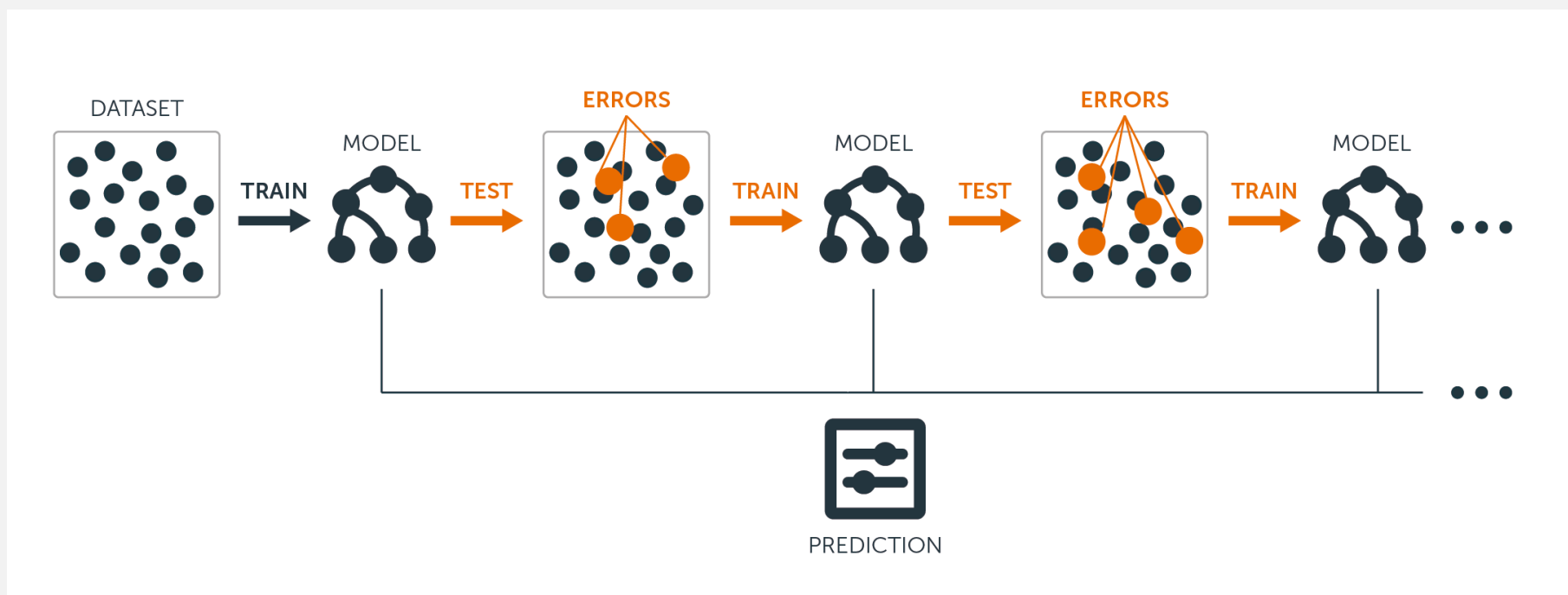
Walk-Forward Cross Validation



Алгоритм?

Алгоритм?

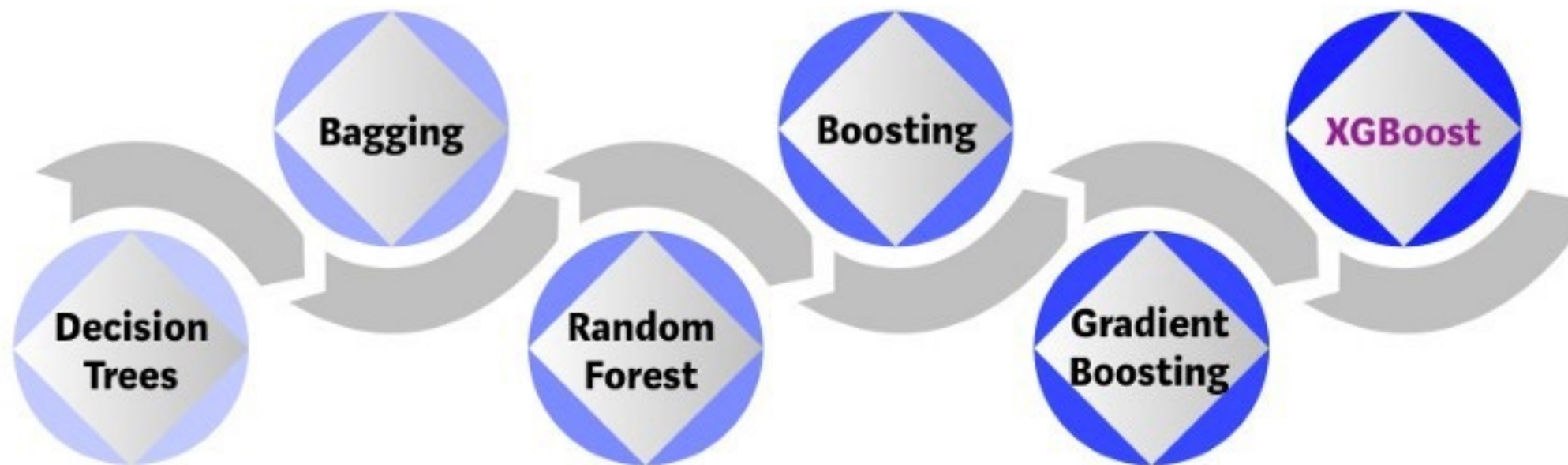
- Бустинги – основной выбор



Bootstrap aggregating or Bagging is a ensemble meta-algorithm combining predictions from multiple- decision trees through a majority voting mechanism

Models are built sequentially by minimizing the errors from previous models while increasing (or boosting) influence of high-performing models

Optimized Gradient Boosting algorithm through parallel processing, tree-pruning, handling missing values and regularization to avoid overfitting/bias

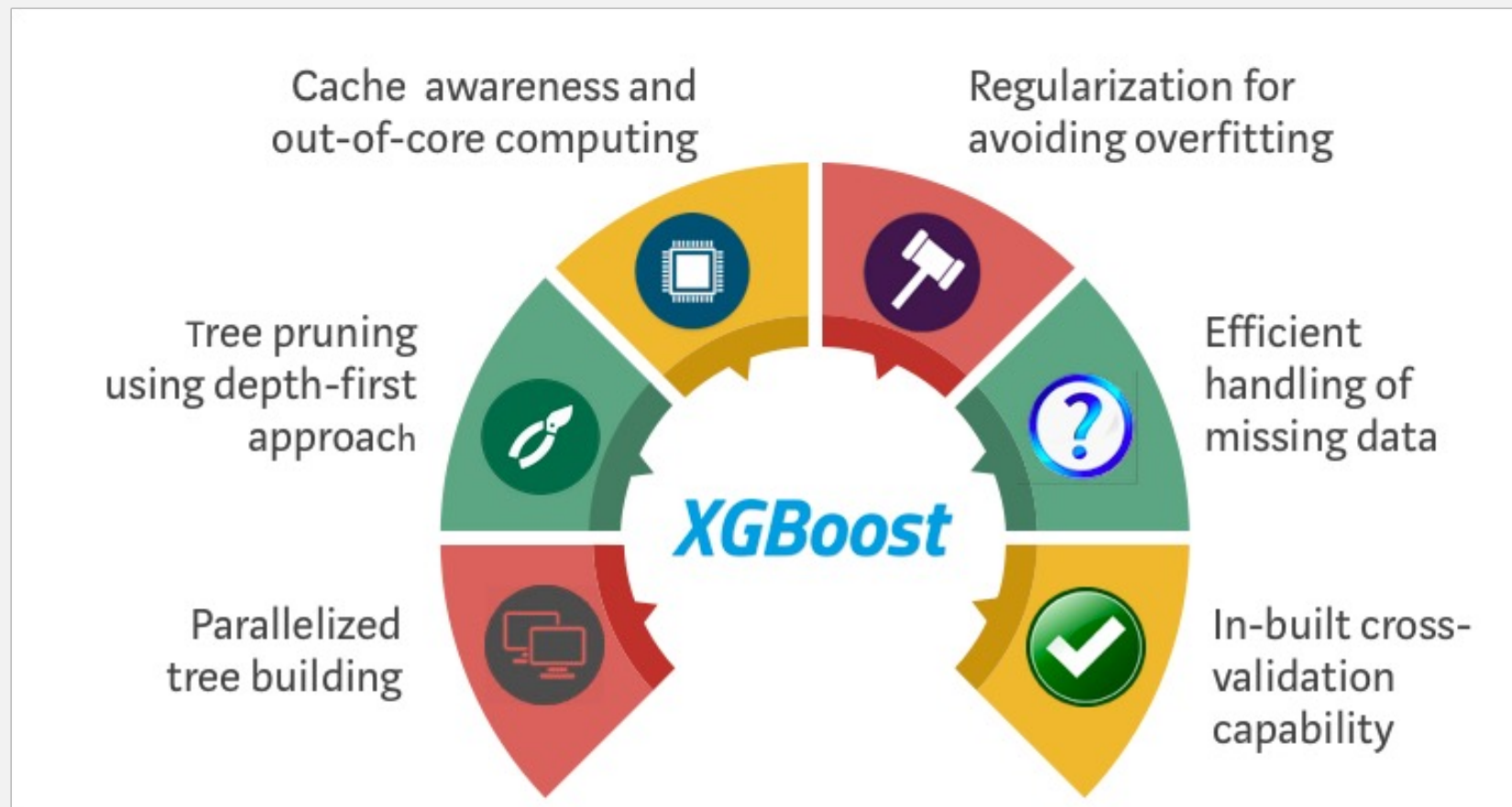


A graphical representation of possible solutions to a decision based on certain conditions

Bagging-based algorithm where only a subset of features are selected at random to build a forest or collection of decision trees

Gradient Boosting employs gradient descent algorithm to minimize errors in sequential models

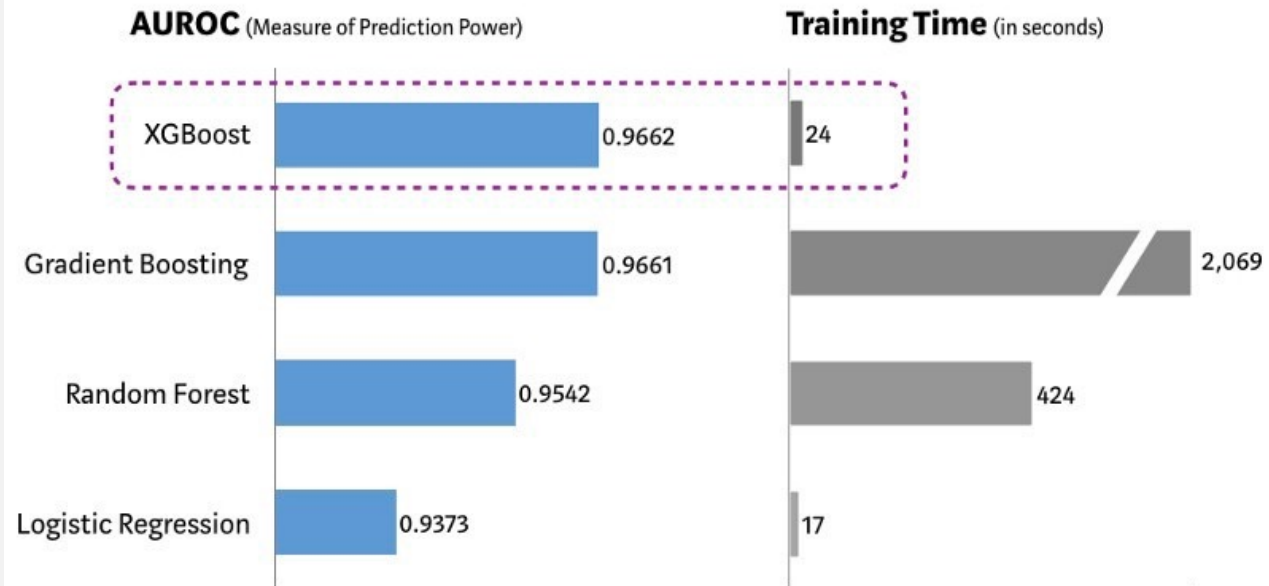
XGBoost



XGBoost

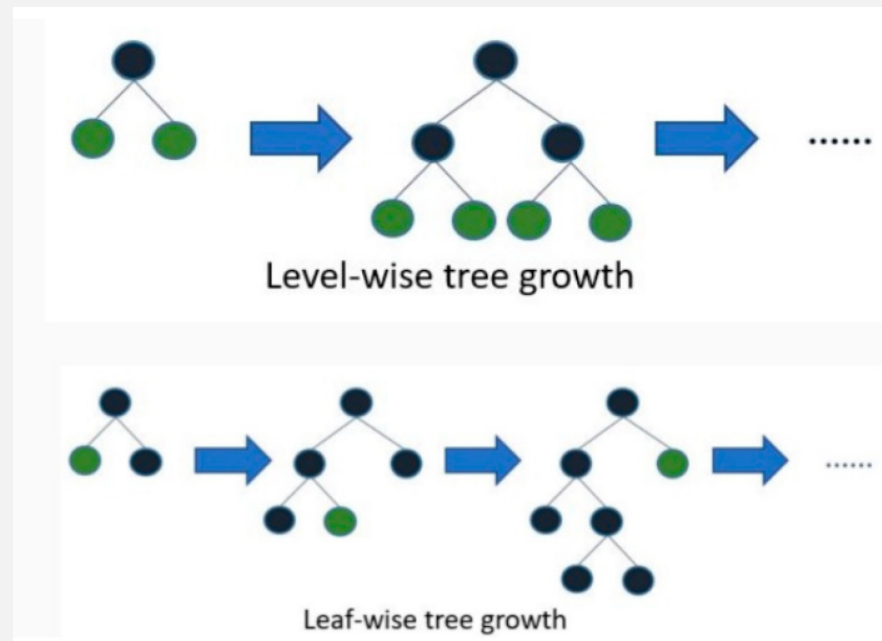
Performance Comparison using SKLearn's 'Make_Classification' Dataset

(5 Fold Cross Validation, 1MM randomly generated data sample, 20 features)



О построении вершин

- Level-wise: дерево строится рекурсивно до тех пор, пока не достигнута максимальная глубина
- Leaf-wise: среди текущих листьев выбирается тот, чьё разбиение сильнее всего уменьшает ошибку

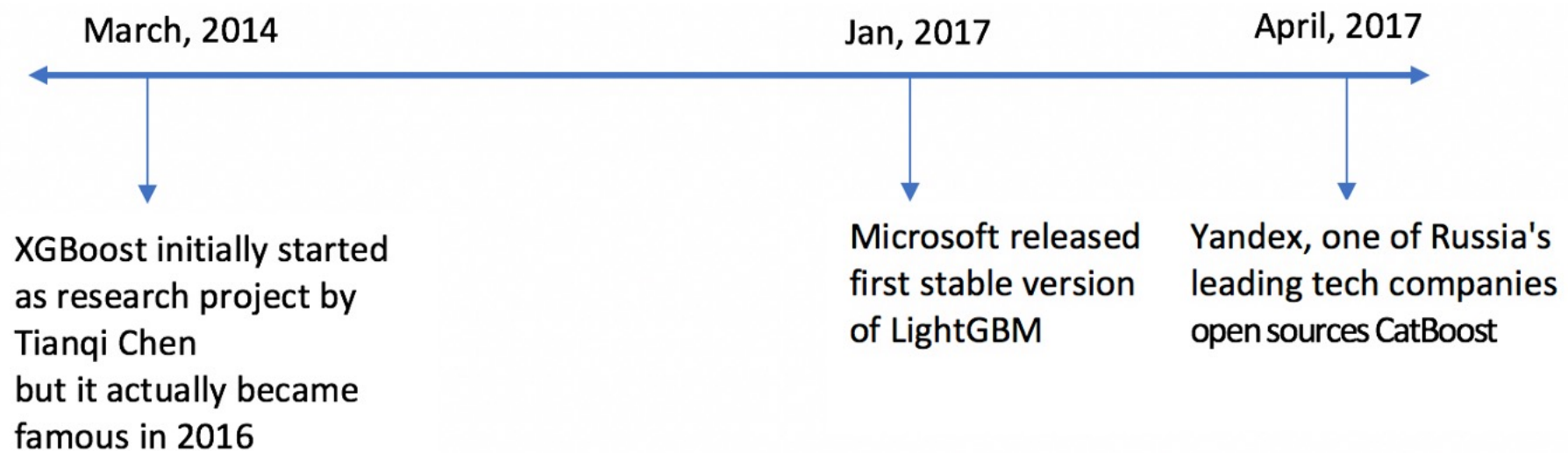


XGBoost vs LightGBM

- XGBoost разветвляет один уровень одновременно, LightGBM – одну вершину
- Разработчики XGBoost добавили эту опцию в свою реализацию, но XGBoost LightGBM быстрее в 1.3 — 1.5 раза, чем XGB.

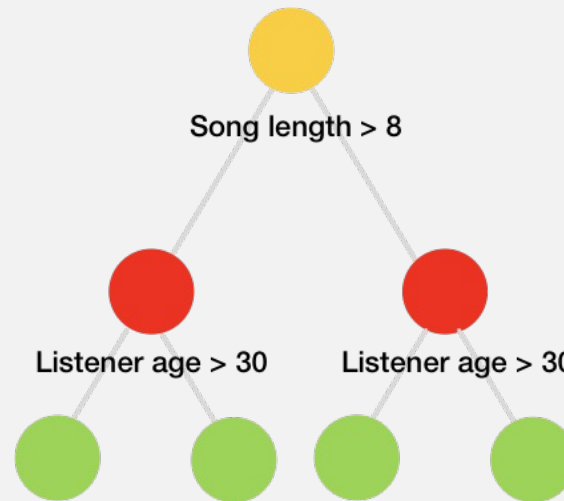
XGBoost vs LightGBM

- XGBoost разветвляет один уровень одновременно, LightGBM – одну вершину
- Разработчики XGBoost добавили эту опцию в свою реализацию, но XGBoost LightGBM быстрее в 1.3 — 1.5 раза, чем XGB.



CatBoost

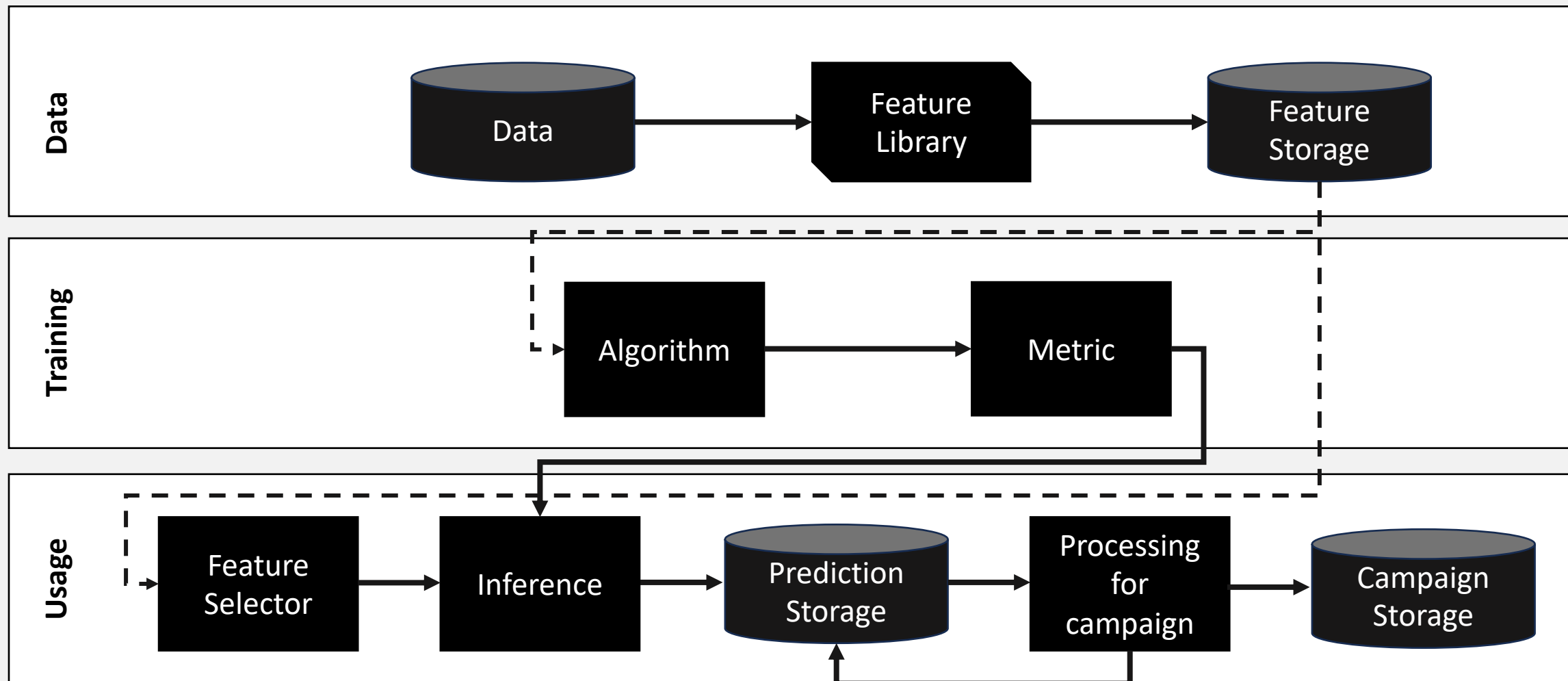
- Oblivious decision trees
- Ограничение: на одном уровне дерева используется один и тот же предикат



	CatBoost	LightGBM	XGBoost
Developer	Yandex	Microsoft	DMLC
Release Year	2017	2016	2014
Tree Symmetry	Symmetric	Asymmetric Leaf-wise tree growth	Asymmetric Level-wise tree growth
Splitting Method	Greedy method	Gradient-based One-Side Sampling (GOSS)	Pre-sorted and histogram-based algorithm
Type of Boosting	Ordered	-	-
Numerical Columns	Support	Support	Support
Categorical Columns	Support Perform one-hot encoding (default) Transforming categorical to numerical columns by border, bucket, binarized target mean value, counter methods available	Support, but must use numerical columns Can interpret ordinal category	Supports, but must use numerical columns Cannot interpret ordinal category, users must convert to one-hot encoding, label encoding or mean encoding
Text Columns	Support Support Bag-of-Words, Naïve-Bayes or BM-25 to calculate numerical features from text data	Do not support	Do not support
Missing values	Handle missing value Interpret as NaN (default) Possible to interpret as error, or processed as minimum or maximum values	Handle missing value Interpret as NaN (default) or zero Assign missing values to side that reduces loss the most in each split	Handle missing value Interpret as NaN (tree booster) or zero (linear booster) Assign missing values to side that reduces loss the most in each split

Function	XGBoost	CatBoost	Light GBM
Important parameters which control overfitting	<ol style="list-style-type: none"> 1. learning_rate or eta – optimal values lie between 0.01-0.2 2. max_depth 3. min_child_weight: similar to min_child leaf; default is 1 	<ol style="list-style-type: none"> 1. Learning_rate 2. Depth - value can be any integer up to 16. Recommended - [1 to 10] 3. No such feature like min_child_weight 4. l2-leaf-reg: L2 regularization coefficient. Used for leaf value calculation (any positive integer allowed) 	<ol style="list-style-type: none"> 1. learning_rate 2. max_depth: default is 20. Important to note that tree still grows leaf-wise. Hence it is important to tune num_leaves (number of leaves in a tree) which should be smaller than $2^{(\text{max_depth})}$. It is a very important parameter for LGBM 3. min_data_in_leaf: default=20, alias= min_data, min_child_samples
Parameters for categorical values	Not Available	<ol style="list-style-type: none"> 1. cat_features: It denotes the index of categorical features 2. one_hot_max_size: Use one-hot encoding for all features with number of different values less than or equal to the given parameter value (max – 255) 	<ol style="list-style-type: none"> 1. categorical_feature: specify the categorical features we want to use for training our model
Parameters for controlling speed	<ol style="list-style-type: none"> 1. colsample_bytree: subsample ratio of columns 2. subsample: subsample ratio of the training instance 3. n_estimators: maximum number of decision trees; high value can lead to overfitting 	<ol style="list-style-type: none"> 1. rsm: Random subspace method. The percentage of features to use at each split selection 2. No such parameter to subset data 3. iterations: maximum number of trees that can be built; high value can lead to overfitting 	<ol style="list-style-type: none"> 1. feature_fraction: fraction of features to be taken for each iteration 2. bagging_fraction: data to be used for each iteration and is generally used to speed up the training and avoid overfitting 3. num_iterations: number of boosting iterations to be performed; default=100

Итого про схему прогноза



План

1. О задачах на сегодня
2. Метрики
3. Прогноз клиентских метрик
- 4. Прогноз спроса**

Зачем прогнозировать спрос?

Зачем прогнозировать спрос?

- пополнение товарных запасов \Rightarrow контроль затарки
- уменьшение списаний \Rightarrow сокращение списаний
- управление остатками и управление доступностью витрины
- автоматизация закупок

О задаче спроса

- **Объект x** прогноза
- **Целевая переменная (таргет) y**

О задаче спроса

- **Объект x** прогноза
 - магазин (даркстор) – товар
- **Целевая переменная (таргет) y**
 - продажи в штуках
 - продажи в деньгах
 - прогноз скорости роста количества заказов (посмотреть можно [тут](#))

Метрики успеха?

Метрики успеха

- оффлайн проверки качества прогнозов
- сокращение списаний
- доступность витрины
- ...

Признаки?

Признаки

- лаги продаж товара
- лаги продаж категории
- данные о магазине
- данные о внешнем мире
- промо

Метрики?

Метрики?

- MAPE
- WAPE

$$WAPE = \frac{\sum_{i,t} |y_{i,t} - \hat{y}_{i,t}|}{\sum_{i,t} |y_{i,t}|}$$

Неочевидные моменты про метрики

- WARE может дискриминировать редко покупаемые товары и предлагать не заказывать их

Неочевидные моменты про метрики

- WARE может дискриминировать редко покупаемые товары и предлагать не заказывать их
- Не всегда важен качественный прогноз по всей сети, сколь важны по самым маржинальным / любимым / флагманским продуктам

Алгоритмы

- Опять-таки, подходит все, что позволит спрогнозировать временной ряд: SARIMA, Prophet, LightGBM/CatBoost, LSTM

Алгоритмы

- Подходит все, что позволит спрогнозировать временной ряд: SARIMA, Prophet, LightGBM/CatBoost, LSTM
- На практике – снова бустинги:
 - учет огромного числа факторов
 - масштабируемость
 - победа в честных проверках перфоманса

Что может помешать?

Что может помешать?

- Проблема «холодного старта»: новые товары/магазины, недостаточно предыстории самого объекта
- Недостаточно наблюдений для редких категорий: мало данных, редко продаются
- Сезонность
- Рост бизнеса: компания масштабируется, продажи растут
- Ловушка заниженной доступности: мало продали, потому что мало привезли, мало продаем
- Перепрогноз / затарка: много заказали, не продали, большие скидки

Ожидание

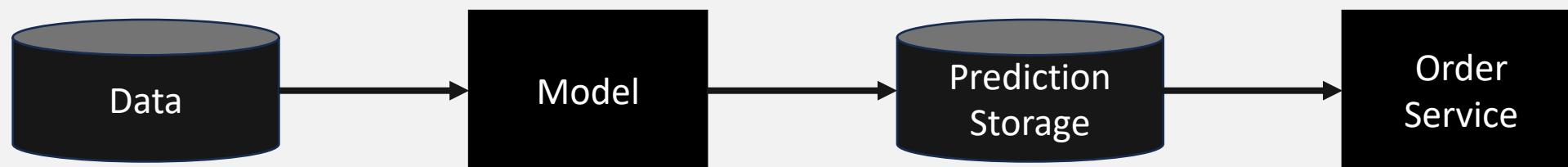
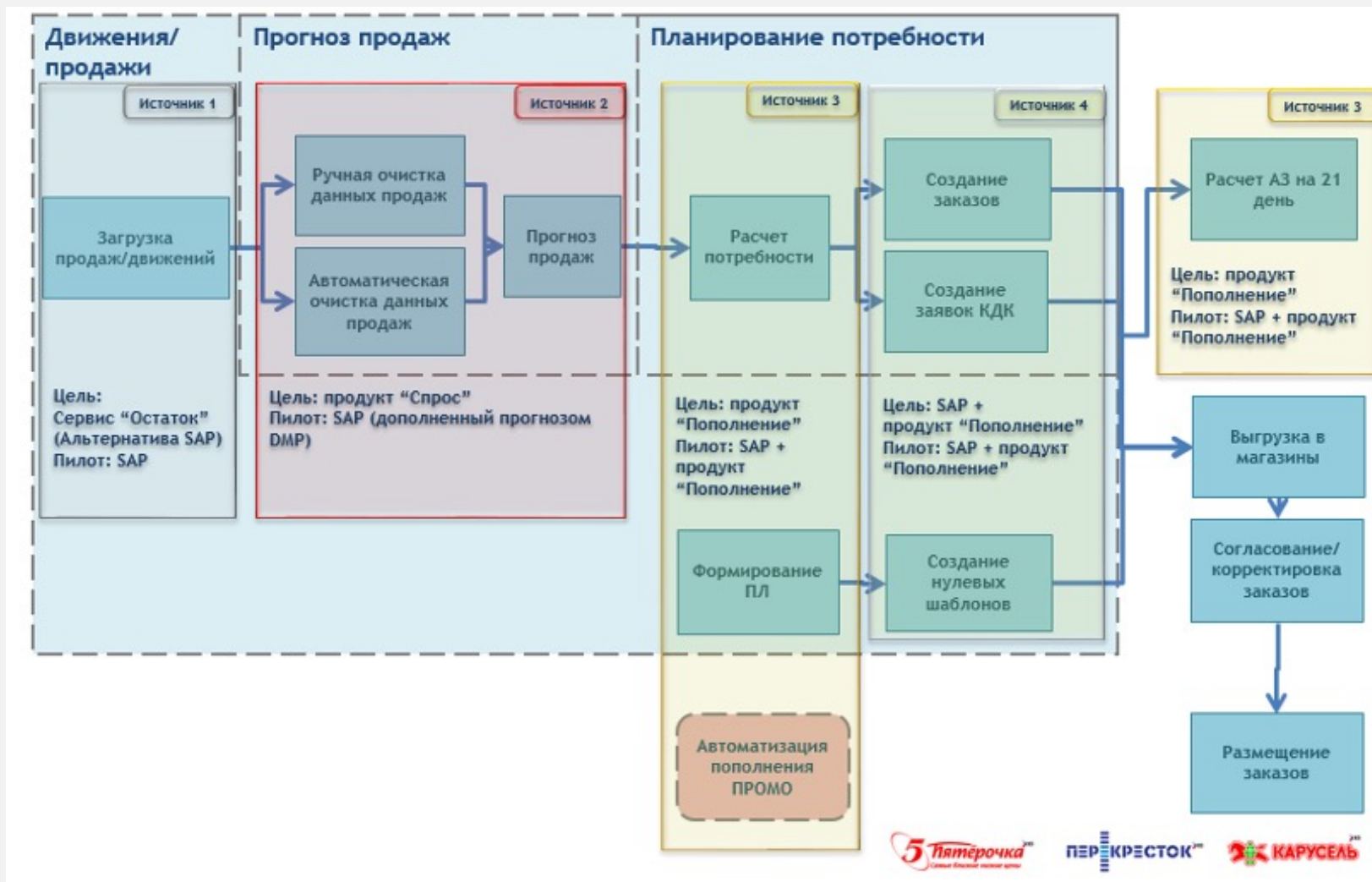
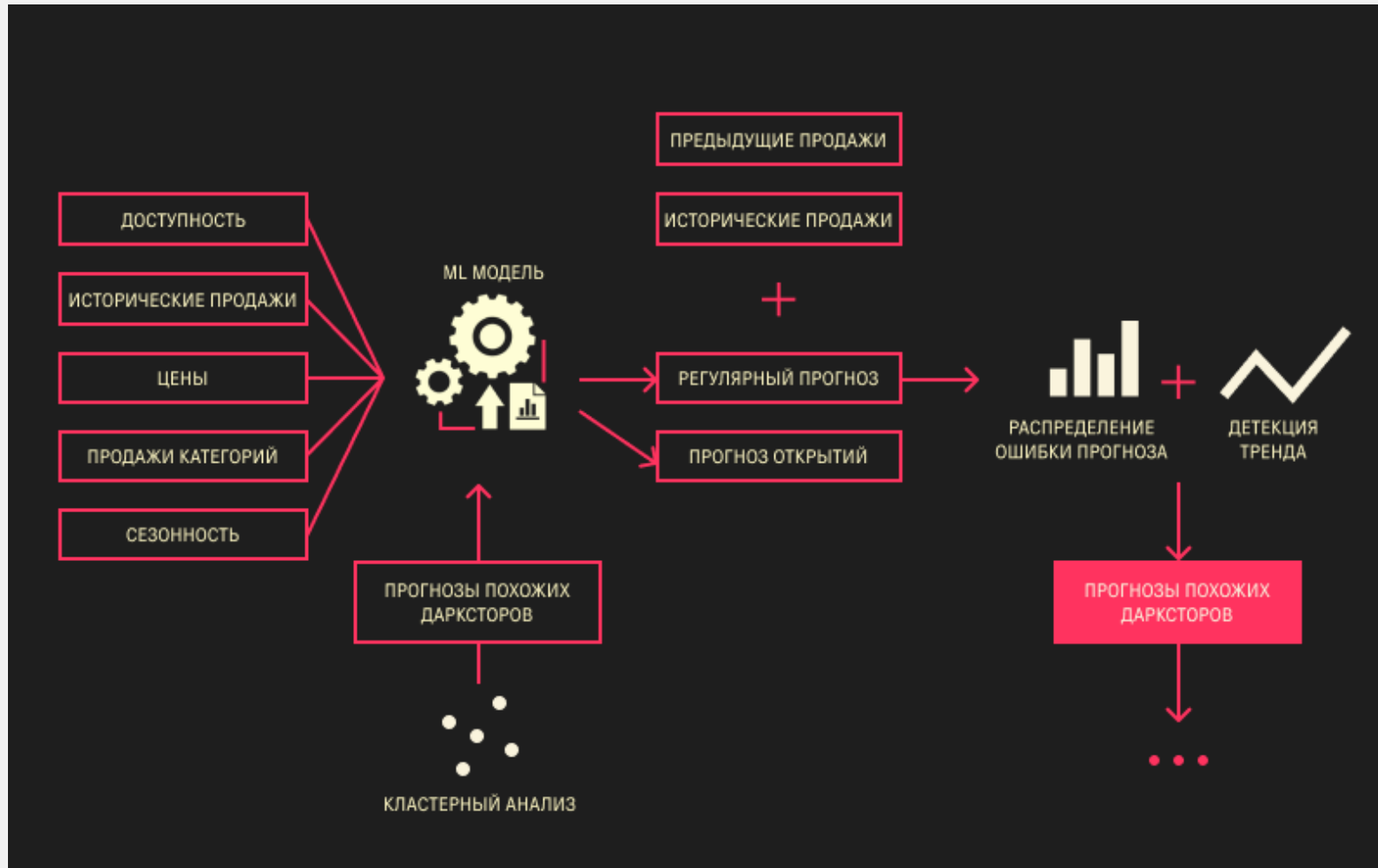


Схема работы в Х5



[Ссылка](#)

Схема работы модуля в Самокате



[Ссылка](#)