

# Автоматизированные методы машинного обучения

Алексей Зайцев  
Руководитель  
лаборатории  
Сколтех

**Skoltech**

A black and white photograph of a modern, multi-story building with a complex, angular roofline and large windows, identified as the Skoltech building.

# Про что мы сегодня поговорим

1. Этапы анализа данных. На каких из них может помочь AutoML
2. Существующие решения для автомл: Optuna, AutoML by Sber, AutoGluon by Amazon
3. Простой AutoML, как использовать из коробки. Grid search, Optuna
4. Правильная валидация моделей

# Этапы анализа данных

---



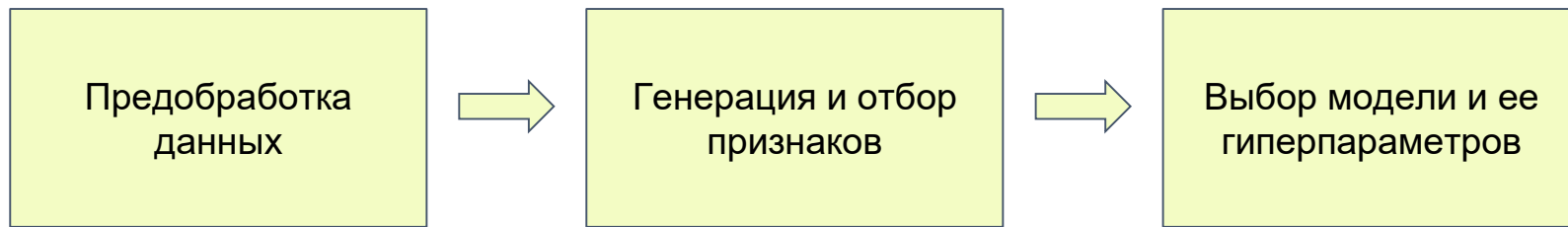
# Постановка задачи

Создать как можно более качественную модель машинного обучения в рамках ограничений.

Качество: в реальной жизни

Ограничения: время работы, срок создания модели, класс моделей, доступные данные

# Конфигурация решения для автоматизированного машинного обучения AutoML

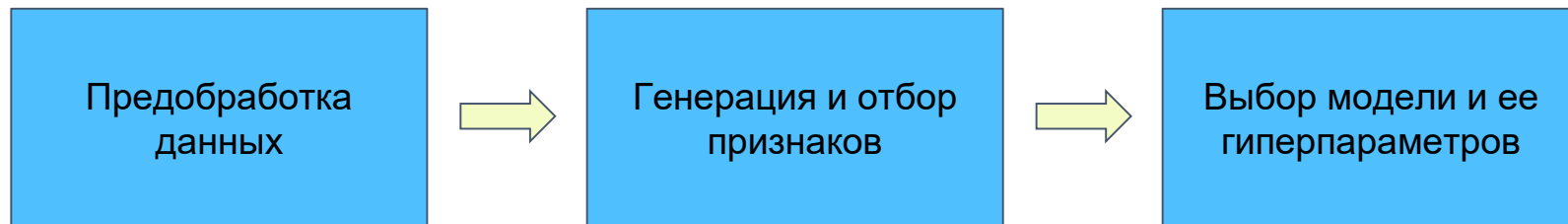


- очистка данных
- детектирование выбросов
- заполнение пропусков

- генерация признаков
- представление категориальных признаков
- отбор признаков

- валидация моделей
- Байесовская оптимизация
- сравнение моделей

# Конфигурация решения для автоматизированного машинного обучения AutoML



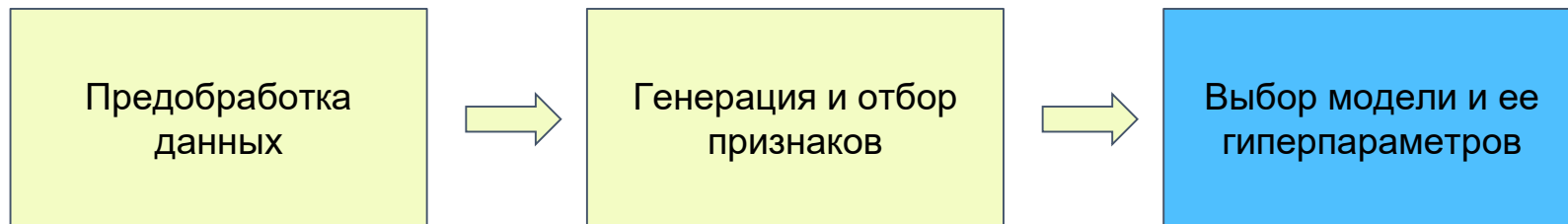
- очистка данных
- детектирование выбросов
- заполнение пропусков

- генерация признаков
- представление категориальных признаков
- отбор признаков

- валидация моделей
- Байесовская оптимизация
- сравнение моделей

На всех этапах применимы подходы автоматизированного машинного обучения: AutoML, AutoGluon

# Конфигурация решения для автоматизированного машинного обучения AutoML



- очистка данных
- детектирование выбросов
- заполнение пропусков

- генерация признаков
- представление категориальных признаков
- отбор признаков

- валидация моделей
- Байесовская оптимизация
- сравнение моделей

На этапе выбора модели могут помочь подходы для оптимизации гиперпараметров: Optuna, Hyperopt

# **Существующие решения для AutoML**





# Существующие решения общего назначения

Библиотека	Разработчик	Ссылка	Функциональность	Ориентация на банковские приложения	Автоматизация всех этапов
AutoGluon	Amazon	<a href="https://github.com/awsmlabs/autogluon">https://github.com/awsmlabs/autogluon</a>	Работа с разными данными	нет	нет
LightAutoML	Sberbank	<a href="https://github.com/sberbank-ai-lab/LightAutoML">https://github.com/sberbank-ai-lab/LightAutoML</a>	Работа с табличными данными	да	нет
NNI	Microsoft	<a href="https://github.com/microsoft/nni">https://github.com/microsoft/nni</a>	Работа с разными данными	нет	нет
H2O	H2O	<a href="https://github.com/h2oai/h2o-3">https://github.com/h2oai/h2o-3</a>	Все	нет	да
Auto sklearn	Auto sklearn	<a href="https://github.com/automl/auto-sklearn">https://github.com/automl/auto-sklearn</a>	Работа с табличными данными	нет	нет

# Существующие решения общего назначения

Библиотека	Разработчик	Ссылка	Функциональность
Optuna	Optuna	<a href="https://github.com/optuna/optuna">https://github.com/optuna/optuna</a>	Оптимизация гиперпараметров
Hyperopt	Hyperopt	<a href="https://github.com/hyperopt/hyperopt">https://github.com/hyperopt/hyperopt</a>	Оптимизация гиперпараметров
FeatureTools	Alteryx	<a href="https://github.com/alteryx/featuretools">https://github.com/alteryx/featuretools</a>	Работа с признаками

# Основные библиотеки для работы

Выберем джентельменский набор:

- Optuna для оптимизации гиперпараметров
- AutoML by Sber для оптимизации всего
- AutoGluon by Amazon для оптимизации всего



OPTUNA



AutoGluon

# Optuna

Основные свойства:

- Используются для оптимизации гиперпараметров
- Легковесная и понятная в использовании
- Способна к параллелизации
- Доступна визуализация как dashboard
- Эффективные алгоритмы поиска, способна к ранней остановке

Способна работать со следующими модулями:

- Catboost, LightGBM, XGBoost
- Pytorch, Keras, TensorFlow
- Chainer, MXNet
- FastAI V1, V2
- AllenNLP
- Catalyst



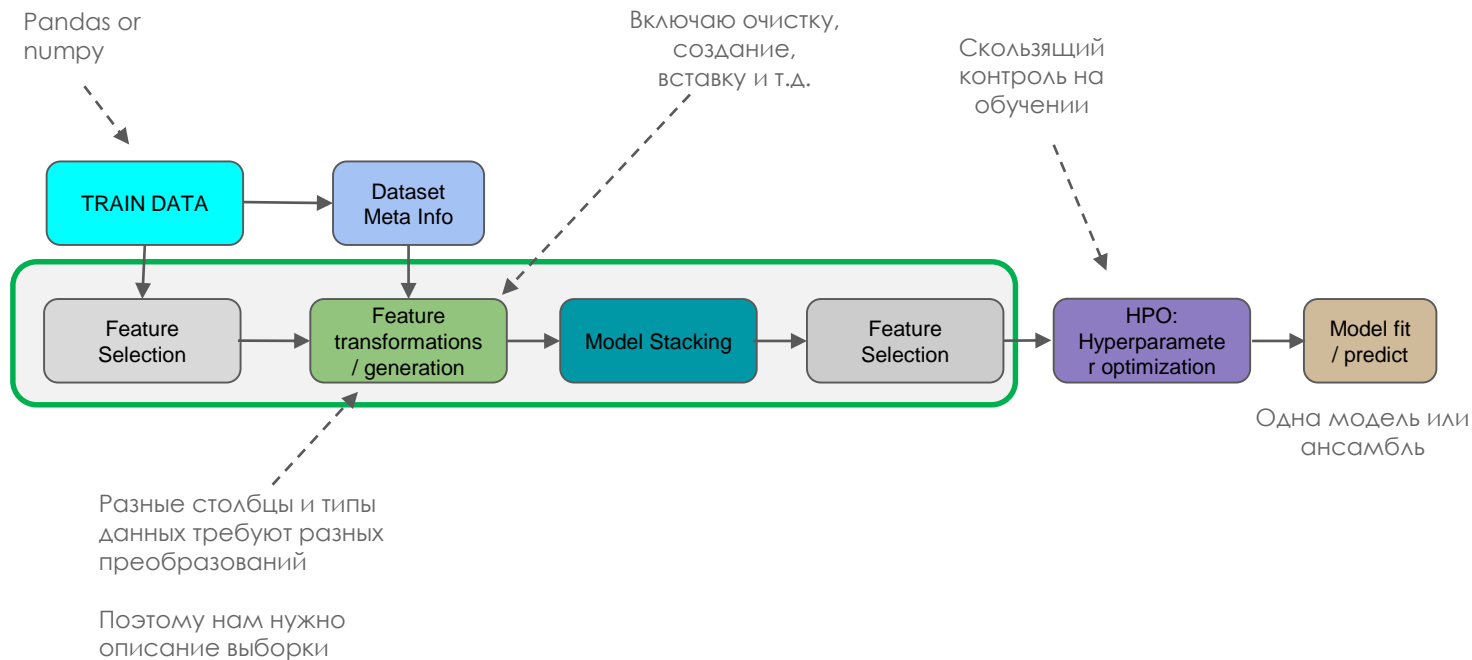
O P T U N A

# Optuna - доступные методы оптимизации

- Grid Search
- Random Search
- Tree-structured Parzen Estimator algorithm
- CMA-ES based algorithm
- Algorithm to enable partial fixed parameters
- Nondominated Sorting Genetic Algorithm II
- A Quasi Monte Carlo sampling algorithm



# Схема пайплайна для обработки данных



# AutoGluon by Amazon

Основные свойства:

- Применима для табличных данных, изображений и текстовых данных
- Для запуска необходимо несколько строк кода
- Встроены подходы для оптимизации гиперпараметров (HPO)
- Кастомизируема



# LightAutoML by Sber

Совокупность инструментов для автоматизированного решения задач построения эффективных моделей на основе данных.

AutoML позволяет быстро получать решение высокого качества, учитывает особенности входных данных.

Базовые компоненты библиотеки AutoML:

- Подбор гиперпараметров
- Кодирование признаков
- Стекинг моделей
- Аналитика входных данных





# LightAutoML by Sber - реализованные блоки

Model adapters	Data transformation	Data Meta	Main Module	Feature selection
<p>a wrapper for various models with unified:</p> <ul style="list-style-type: none"> <li>• fit</li> <li>• predict</li> <li>• save</li> <li>• load</li> <li>• get_params</li> <li>• set_params</li> </ul>	<p>a wrapper for data processing methods with unified:</p> <ul style="list-style-type: none"> <li>• fit</li> <li>• transform</li> <li>• k-fold split during training</li> </ul>	<p>data transformations history logging, utils for data slicing and keeping feature types labeled</p>	<p>Includes right now</p> <ul style="list-style-type: none"> <li>• fit</li> <li>• predict with internal data processing and feature generation</li> <li>• cross validation</li> </ul>	<ul style="list-style-type: none"> <li>• permutation selection</li> <li>• shap tree selector based on LGBM</li> </ul>
<p>Stacking adapter</p> <p>Stacking as a part of feature generation</p> <ul style="list-style-type: none"> <li>• one model stack</li> <li>• multi model stacker</li> <li>• hyperparameters for stackers</li> <li>• different stacking strategies</li> </ul>	<p>Cross validation with stratification by a number of groups for cross validation</p> <p>Splits data into k-folds and keeps distribution in each fold for specified columns roughly the same. Not deterministic but works pretty well.</p>	<p>HPO - hyperparameters optimization</p> <ul style="list-style-type: none"> <li>• HPO for individual models</li> <li>• logging</li> <li>• HPO history for one model</li> </ul>	<p>SAVING</p> <ul style="list-style-type: none"> <li>• saving separate models</li> </ul>	

# **Примеры использования из коробки**

---

# GridSearch

**Идея:** использовать полный перебор возможных параметров, выбрать параметры лучшего результата (например, лучшая валидация)

Плюсы и минусы подхода:

- + Перебор по всем возможным параметрам, решение будет гарантировано лучшее на рассматриваемом множестве
- + Простой в реализации подход
- Время работы растет экспоненциально с числом новых параметров
- В каждом случае нужно ждать конца обучения модели

```
>>> from sklearn import svm, datasets
>>> from sklearn.model_selection import GridSearchCV
>>> iris = datasets.load_iris()
>>> parameters = {'kernel':('linear', 'rbf'), 'C':[1, 10]}
>>> svc = svm.SVC()
>>> clf = GridSearchCV(svc, parameters)
>>> clf.fit(iris.data, iris.target)
GridSearchCV(estimator=SVC(),
              param_grid={'C': [1, 10], 'kernel': ('linear', 'rbf')})
>>> sorted(clf.cv_results_.keys())
['mean_fit_time', 'mean_score_time', 'mean_test_score',...
 'param_C', 'param_kernel', 'params',...
 'rank_test_score', 'split0_test_score',...
 'split2_test_score', ...
 'std_fit_time', 'std_score_time', 'std_test_score']
```

# Optuna

**Идея:** по умолчанию используется Tree-structured Parzen Estimator algorithm - подход, основанный на гауссовских процессах.

Плюсы и минусы подхода:

- + Для получения оптимальных параметров не надо перебирать всевозможные значения
- + За малое количество шагов оптимизации ищет решение, близкое к оптимальному
- Нужно полностью ждать обучения модели (однако с помощью инструмента pruner можно это нивелировать)
- При сложном пространстве параметров не гарантирует нахождение оптимальных

# Optuna - разберем пример

Здесь присутствуют три основных сущности подхода:

- objective - целевая функция, которую максимизируем (минимизируем)
- trial - объект, содержит в себе информацию об оптимизируемых гиперпараметрах, также обозначает шаг оптимизатора
- study - объект обучения

```
import optuna

def objective(trial):
    x = trial.suggest_float('x', -10, 10)
    return (x - 2) ** 2

study = optuna.create_study()
study.optimize(objective, n_trials=100)

study.best_params # E.g. {'x': 2.002108042}
```

# AutoGluon

Пример работы AutoGluon на табличных данных

```
from autogluon.tabular import TabularDataset, TabularPredictor

data_root = 'https://autogluon.s3.amazonaws.com/datasets/Inc/'
train_data = TabularDataset(data_root + 'train.csv')
test_data = TabularDataset(data_root + 'test.csv')

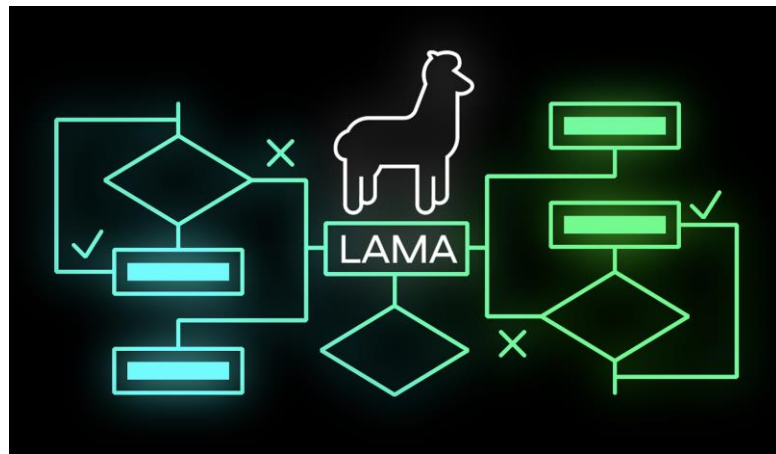
predictor = TabularPredictor(label='class').fit(train_data=train_data)
predictions = predictor.predict(test_data)
```

# AutoML

- Новые блоки легко добавлять с помощью общих классов адаптеров и файла конфигурации
- Порядок блоков и выполнения операций задается пользователем и под каждую задачу можно быстро собрать свой уникальный пайплайн.

## Пример работы AutoML

```
aml_light = AMLib(  
    features_processing=['MEstimateEncoder', 'LeaveOneOutEncoder', 'FeSelectorPer'],  
  
    models_names=['LGBM'],  
    problem_type="classification",  
    default_fe_processing_sample = 1,  
    hpo_n_trials=25,  
    cv_groups = ['X5', 'X6'],  
    metric="roc_auc",  
    user_params = {'crossval': {'n_splits': 5},  
  
                   'FeSelectorPer': {'num_buckets': 7, 'choose_k_first': 3},  
  
                   "cat_encoders": {"n_splits_train": 3, "n_splits_bagging": 30}},  
    dense_transform="RankGauss",)  
  
aml_light.fit(train,  
              target_col='TARGET',  
              column_types=column_types)  
  
preds = aml_light.predict(val)  
print('LGBM', aml_light.metric(val['TARGET'], preds['LGBM'][:,1]))
```



# Как валидировать модели?

---



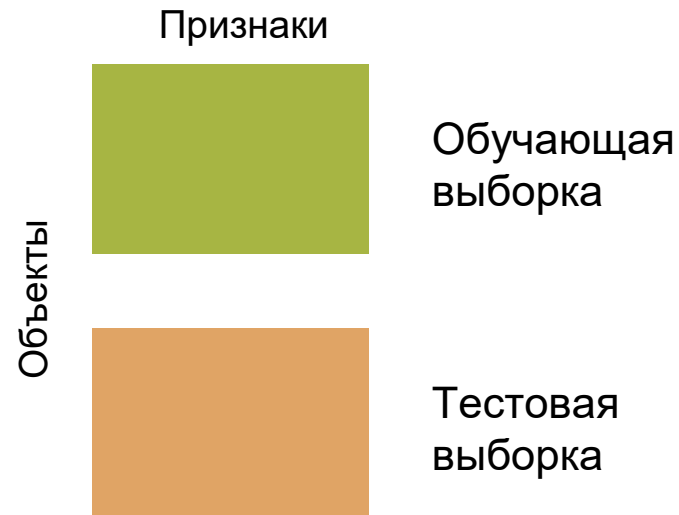


# Тестирование моделей машинного обучения

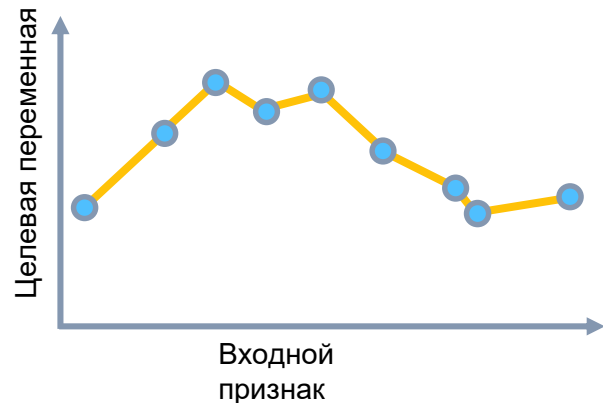
Идеальный случай: есть независимая тестовая выборка, которая похожа на то, как все будет работать в действительности

Решение: обучаем модель на обучающей выборке, тестируем на тестовой (или другое название – валидационной)

Пример задачи: обучение моделей распознавания изображений на ImageNet

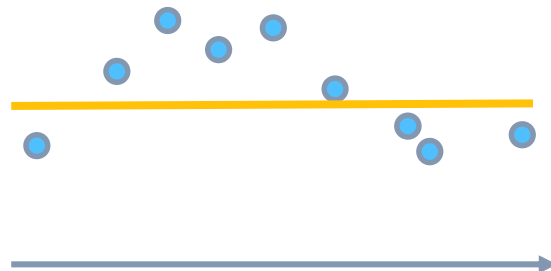


# Какие типы поведения бывают?



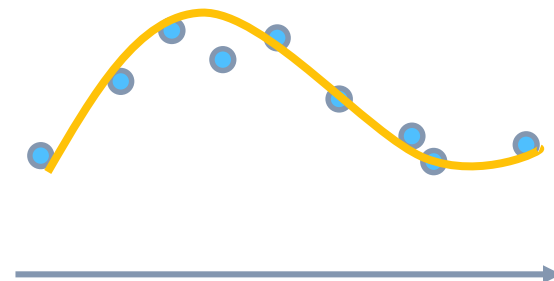
## Переобучение

Модель подстроилась под обучающие данные, но будет плохо работать на тестовой выборке



## Недообучение

Модель недостаточно подстроилась под обучающие данные, на тестовой выборке тоже все будет плохо



## Нормальная модель

Модель ошибается, но не слишком сильно

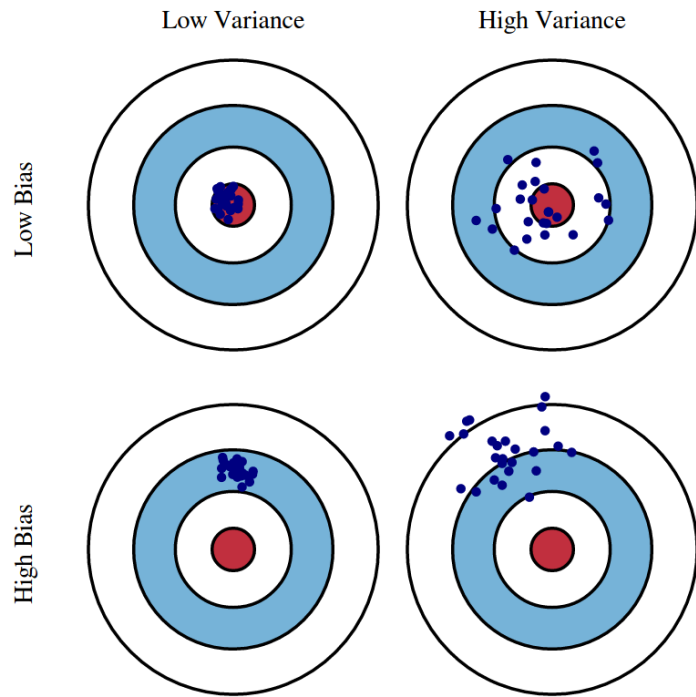
# Оценка качества моделей машинного обучения

Ошибка на обучающей выборке	Ошибка на тестовой выборке	Анализ
Низкая	Низкая	Все ОК
Низкая	Высокая	Переобучение
Высокая	Высокая	Недообучение
Высокая	Низкая	???

# Выбор между смещением и разбросом

*Смещение.* Если возьмем слишком мало данных для обучения – то качество модели, обученной на подвыборке будет меньше, чем качество модели, обученной на всех доступных данных

*Разброс.* Мы усредняем ошибку по тестовым данным. Если их мало – то оценка ошибки будет не очень точной.

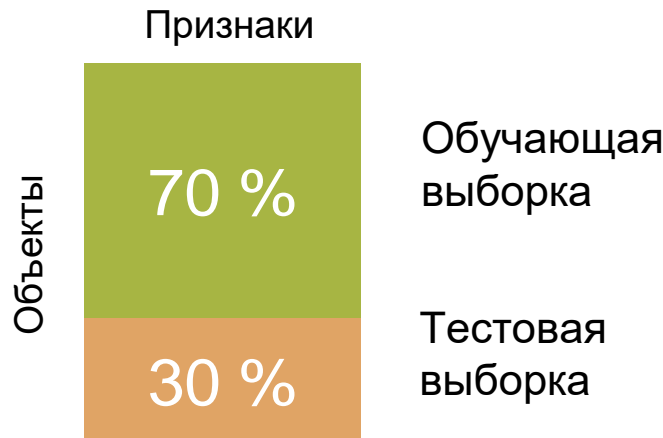


# Выделение тестовой выборки

Обычный случай: есть одна выборка данных, по которой нужно оценивать качество модели

Решение: положить в тестовую выборку то, что больше всего похоже на использование модели в реальности

Пример: есть данные за 2017-2020, для обучения используем 2017-2019, для теста - 2020

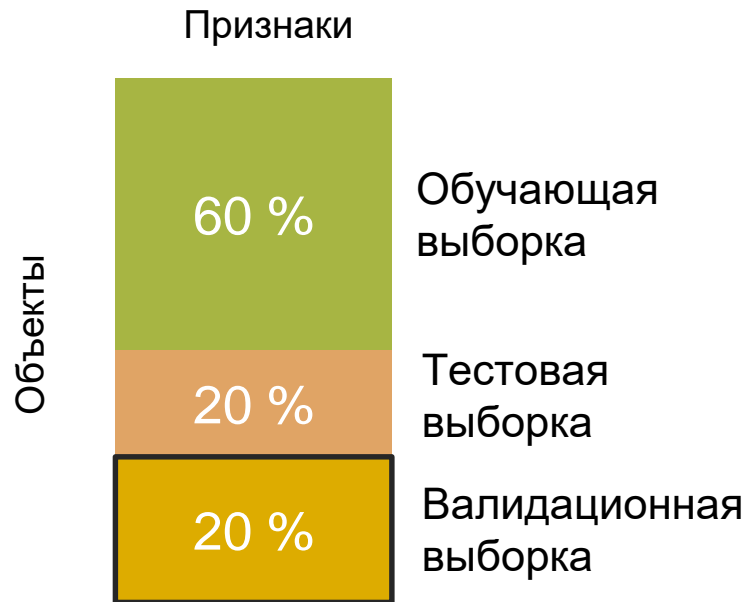


# Выделение *тестовой* и *валидационной* выборки

Обычный случай: есть одна выборка данных, по которой нужно оценивать качество модели

Решение: положить в тестовую выборку то, что больше всего похоже на использование модели в реальности, выделить дополнительную валидационную выборку для итогового тестирования

Пример: есть данные за 2017-2020, для обучения используем 2017-2018, для теста – первые два квартала 2020, для валидации – последние два квартала.

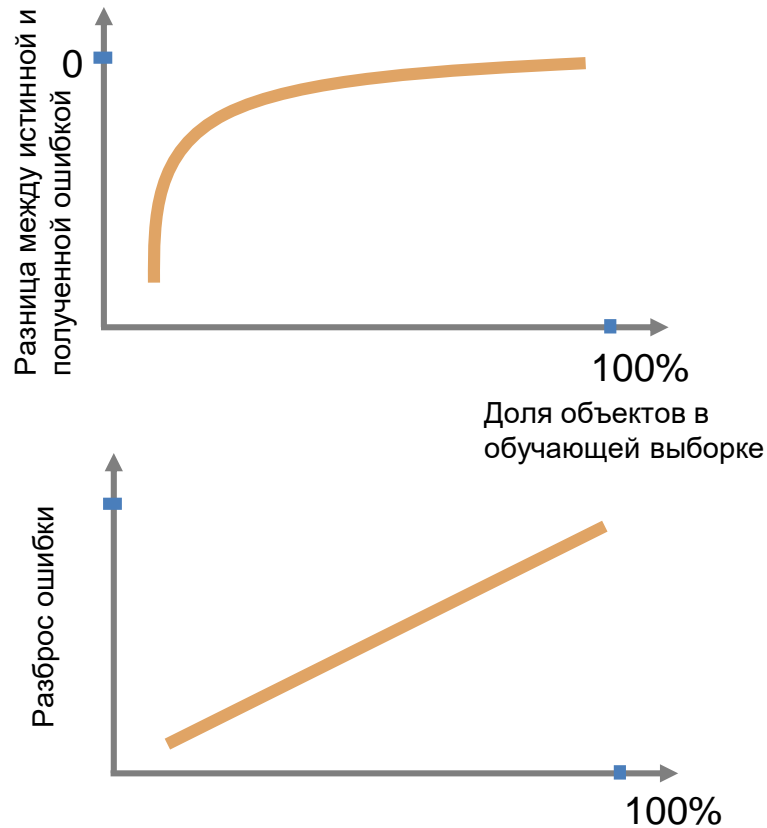


# Выбор размера обучающей, тестовой и валидационной выборки

*Смещение.* Если возьмем слишком мало данных для обучения – то качество модели, обученной на подвыборке будет меньше, чем качество модели, обученной на всех доступных данных

*Разброс.* Мы усредняем ошибку по тестовым данным. Если их мало – то оценка ошибки будет не очень точной.

Получается, что в таком подходе мы используем не все данные для обучения и не все для тестирования. Можно ли делать иначе?

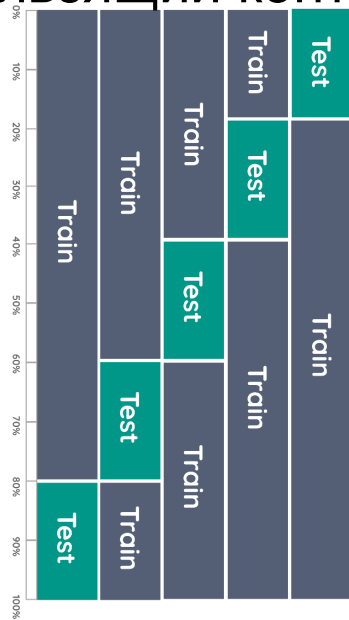


# Скользящий контроль / Cross-validation

## Обычное разбиение



## Скользящий контроль





# Преимущества и недостатки скользящего контроля

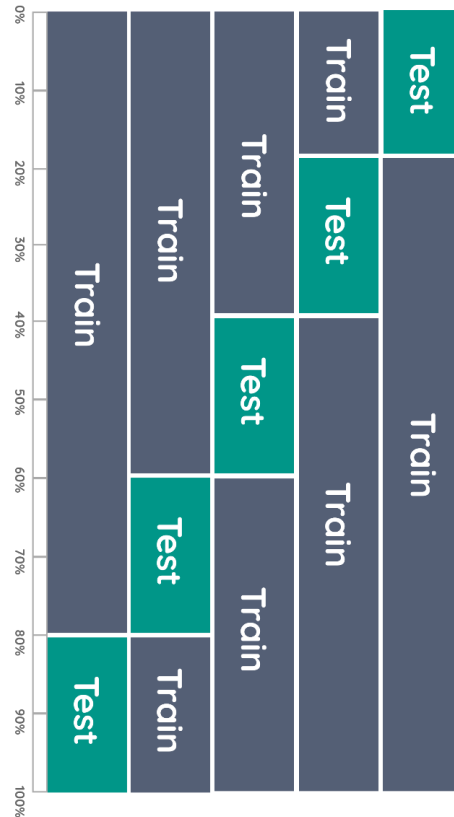
Обычно в процессе скользящего контроля делают 5 разбиений

Можно оценить устойчивость модели: если ошибки сильно отличаются, то модель неустойчива.

Нужно учить много моделей!

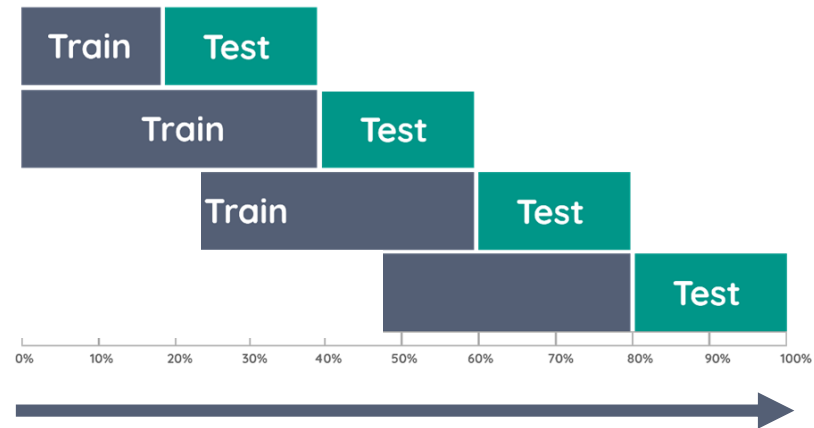
Leave-one-out cross-validation:

Размер тестовой выборки в процессе скользящего контроля – 1

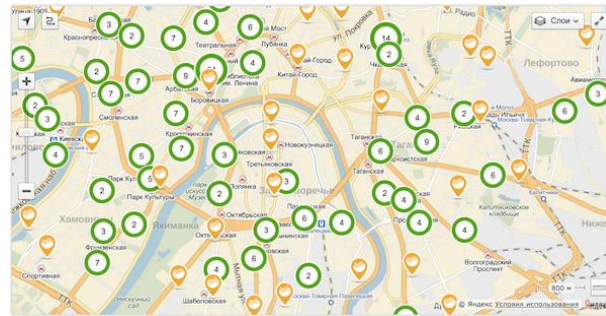


# Специальные случаи оценки качества моделей

Данные удобно разбивать по времени: это похоже на реальный сценарий использования модели, где нам доступны только данные из прошлого, а предсказывать мы должны будущее.

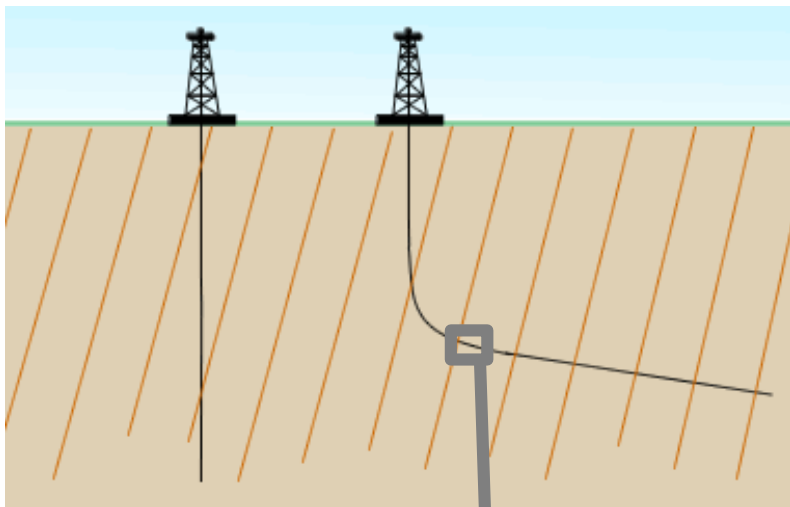


Данные удобно разбивать по географии: если данные для разных регионов – это тоже нужно учесть при разбиении выборки на обучение и контроль!



Время

# Кейс, прогноз типа породы



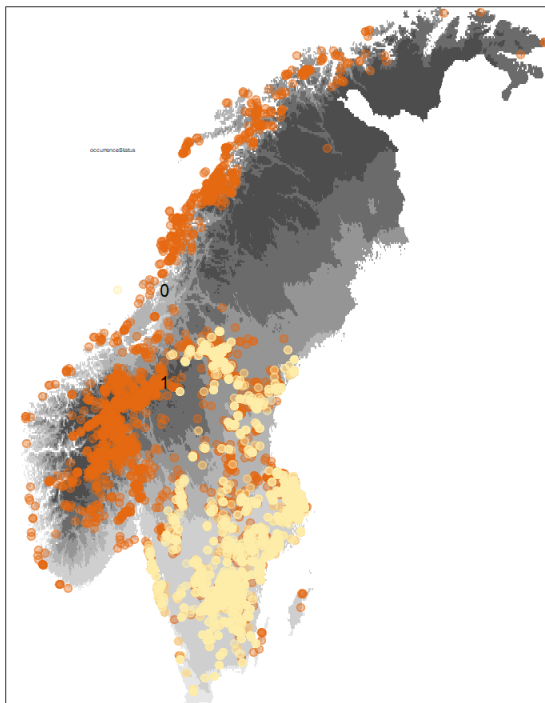
Мы обучали модель на одних скважинах, а тестировали на других

Если разбивать данные без учета принадлежности к скважинам – получается оптимистичная оценка качества модели.

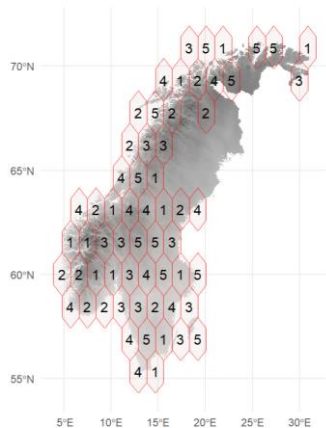
Для проекта про нефтяные месторождения  
объект – часть интервала бурения

Целевая переменная – тип породы в интервале

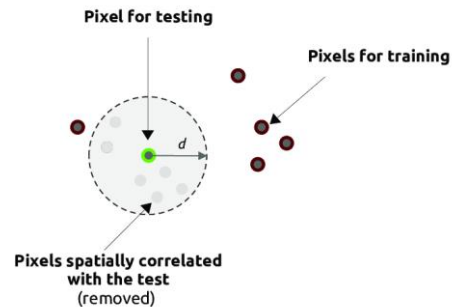
# Кейс, прогноз наличия вида



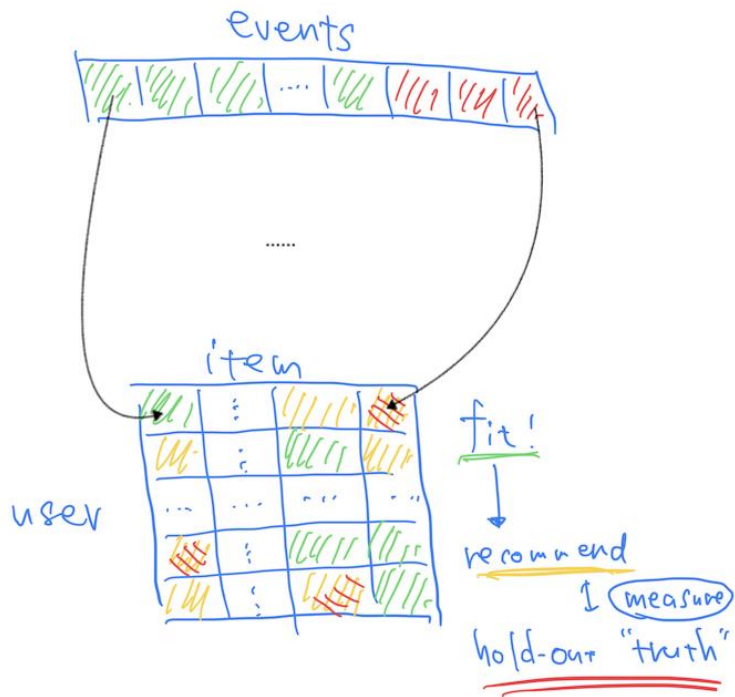
- - наличие
- - отсутствие



Разбиение на фолды



# Валидация для рекомендательных систем

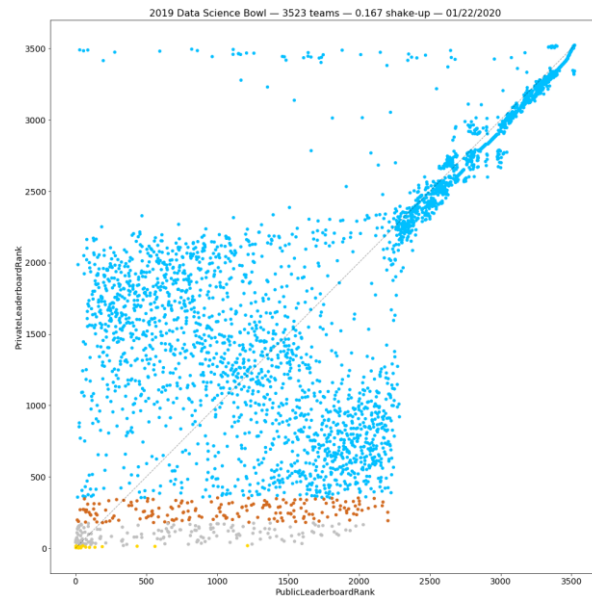


# Shake-up на kaggle

Есть три выборки:

- Train
  - Public test
  - Private test
- Метки участникам доступны только на Train.
  - Они с помощью модели ставят метки на Test,
  - По части Public test считается качество модели, которое видно участникам соревнования
  - По части Private test считается качество модели, которое видно участникам соревнования

Всегда происходит shake-up



Public Private

The private leaderboard is calculated with approximately 75% of the test data.

■ Prize Contenders

#	△	Team	Members	Score	Entries	Last	Code
1	▲ 149	QNS		0.760	52	2d	
2	▲ 17	元宵快乐		0.737	294	2d	< >
3	▲ 1	Team Hydrogen		0.735	382	2d	
4	▲ 3	outrunner		0.728	195	2d	
5	▲ 21	bestfitting		0.728	165	2d	

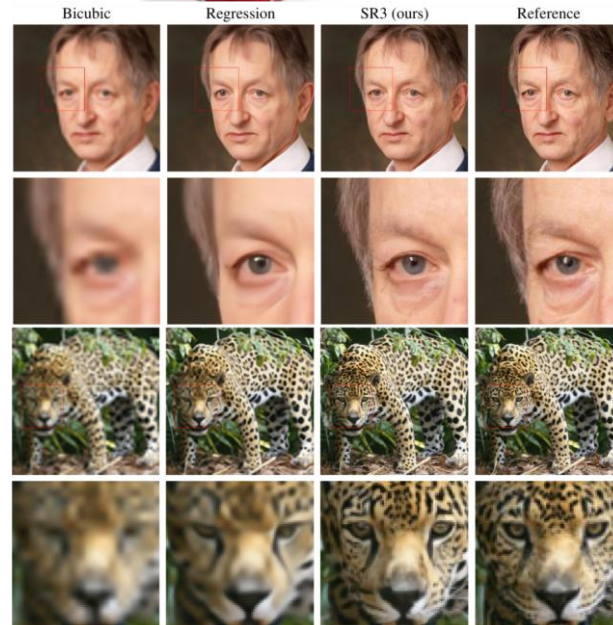
Данные по соревнованиям: <https://www.kaggle.com/datasets/daniboy370/competitions-shakeup>

# Как делать скользящий контроль правильно?

1. Отдельная валидационная выборка, не трогаем ее до конца проекта
2. Скользящий контроль с большим количеством разбиений, чем больше – тем лучше
3. Out-of-time валидация, имитируем реальные условия
4. Принимаем во внимание специфику задачи, когда строит процедуру скользящего контроля – **думаем во время построения процедуры валидации**

# Еще про метрики

- Выбор метрики: качество бизнес-процесса не всегда коррелирует с выбранной метрикой качества.  
Пример: Netflix prize
- Может не быть единой метрики  
Пример: NDCG@5, NDCG@20, NDCG@50, surprise
- Нормальных метрик может не существовать  
Пример: PSNR для super-resolution, генеративные модели





# Анализ внутренних данных по проведенным валидациям в крупном банке



# Внутренние данные по проведенным валидациям

В исходных данных 1014 пар валидаций моделей:  
первичная и периодическая

**Целевая переменная:** перешла ли модель в красную зону  
при периодической валидации

Признаков всего 26:

- дней с прошлой валидации, результат прошлой валидации
- блок-заказчик модели,
- целевая метрика (GINI, R2)
- тип алгоритма и др.

Целевая переменная  
для анализа



Результат периодической  
валидации

# Прогноз исхода валидации моделей

1. Верхнеуровневый прогноз перехода моделей в красную зону в разрезе определенных классов моделей **с точностью 75-85% и полнотой 60-70%**
2. Качество прогноза может быть повышено при использовании низкоуровневых данных по динамике качества каждой модели и изменению факторов



Выборка валидаций: 662 объекта, из них для 154, переход в красную зону

Результаты скользящего контроля в процентах, усреднение по десяти повторениям пяти разбиений

	CatBoost	XGBoost	LightGBM	Logistic Regression	Decision Tree
Precision	86%	<b>81%</b>	80%	75%	77%
Recall	62%	<b>71%</b>	70%	57%	58%
Accuracy	89%	<b>90%</b>	89%	86%	86%

# Анализ эффектов

Наиболее важные признаки:

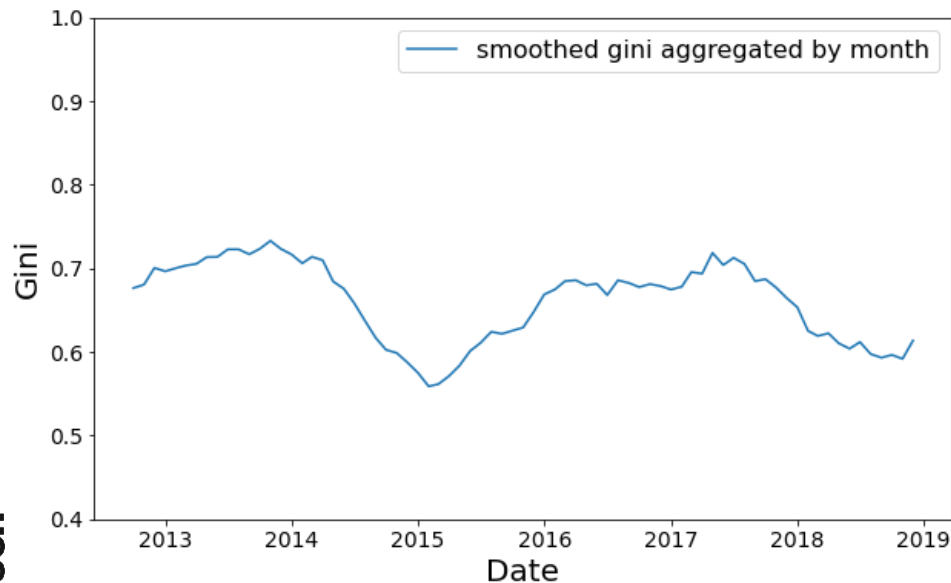
- Число дней с последней валидации
- Результат последней валидации — Желтый
- Значение метрики при последней валидации
- Длительность последней валидации

	Ранг значимости фактора (усредненный по оценкам алгоритмов)	Направление связи с вероятностью перехода модели в красную зону (при росте фактора)
Результат последней валидации — Желтый	1	
Число дней с последней валидации	2	
Задача класса PD	3	
При последней валидации значение метрики улучшилось	4	
Длительность последней валидации	5	
Задача класса отклика для розничных клиентов	6	
Модель для розничных клиентов (независимо от типа задачи)	7	
Тип задачи — прогноз денежных потоков	8	
Значение метрики при последней валидации	9	
Длительность последней валидации относительно среднего	10	

# Динамика метрики качества GINI

Данные для рейтинговой модели, 01/2012-12/2018

Агрегация качества модели GINI по месяцам, скользящее среднее



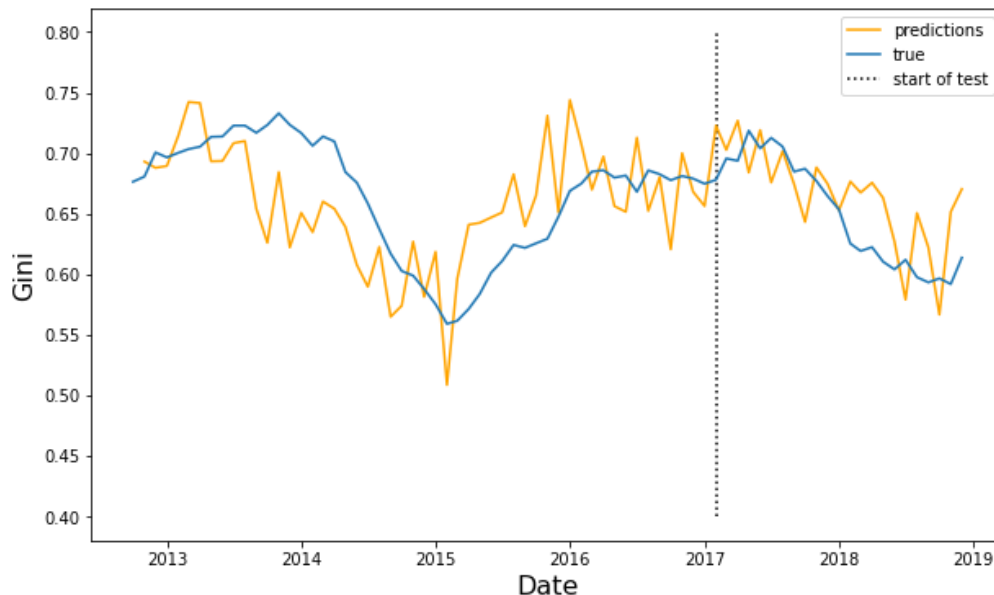
Skoltech

- У качества модели выраженная волатильность
- На качество модели влияют кризисы

# Предсказание метрик качества работы модели

3 статистически значимых фактора:

1. Среднее по кварталу для признака *qual\_factor\_7* – качественная оценка возможности оттока капитала в компании
2. Дисперсия по кварталу для признака *quant\_factor\_4* – соотношение операционной прибыли к выручке
3. Цена нефти *oil\_price*



# Выводы



# Выводы

- Подходы автоматизированного машинного обучения помогают подбирать гиперпараметры на всем этапе решения задачи
- Для более простого случая, выбора модели, достаточно использовать подходы с оптимизацией гиперпараметров
- Существует множество подходов для эффективной и быстрой оптимизации гиперпараметров. Не стоит использовать наивный поиск по сетке
- Оценка ошибки по выборке – нетривиальная задача