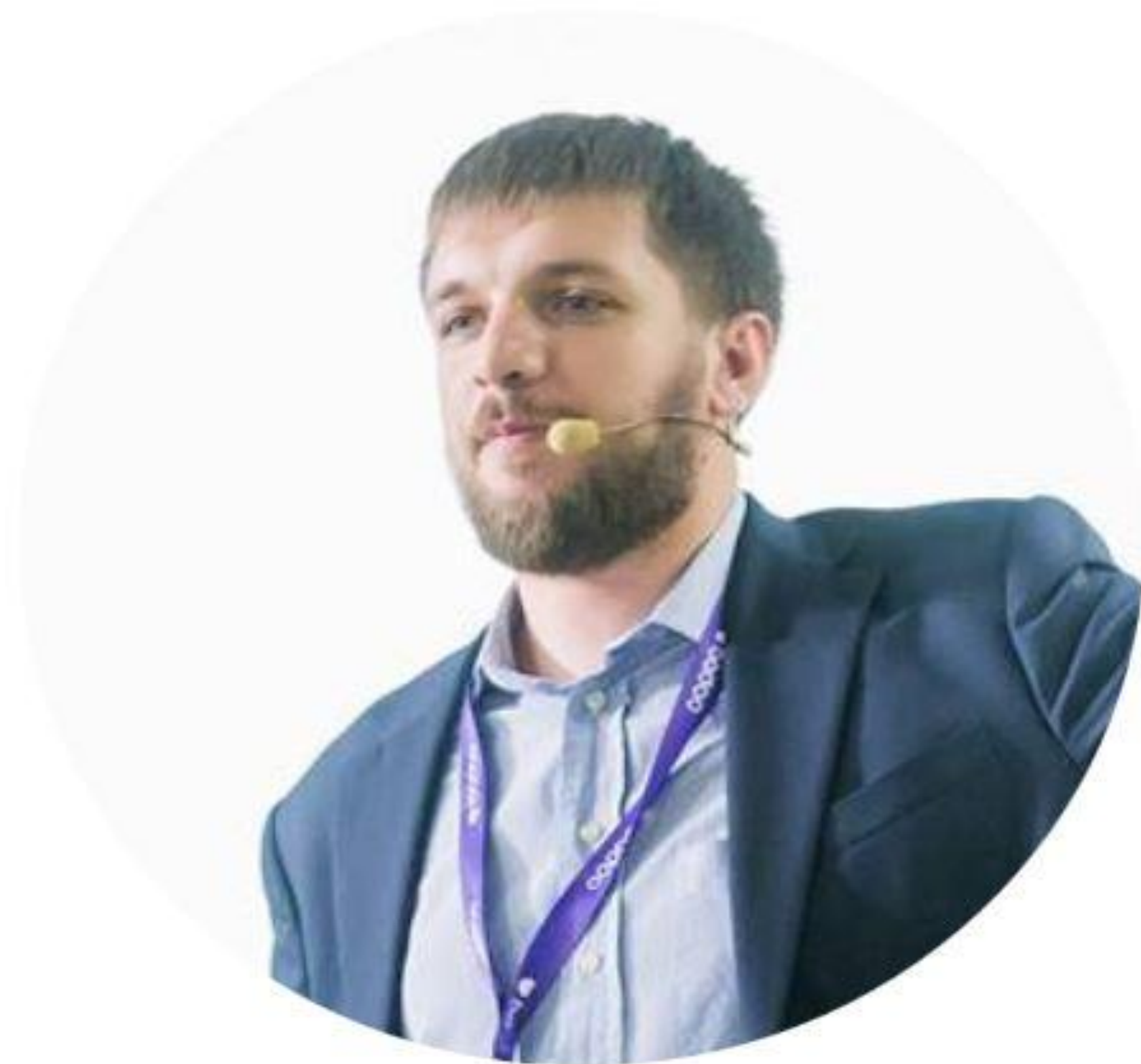


ЗАНЯТИЕ 1.8

МЕТРИКИ КАЧЕСТВА МОДЕЛИ И ПЕРЕОБУЧЕНИЕ



Алексей Кузьмин

Директор разработки; Data Scientist

ДомКлик.ру



aleksej.kyzmin@gmail.com

ЦЕЛИ ЗАНЯТИЯ

В КОНЦЕ ЗАНЯТИЯ ВЫ:

- будете знать как проводить кросс-валидацию модели;
- сможете оценить качество разных версий модели по AUC;
- подберете параметры модели для борьбы с переобучением.

О ЧЁМ ПОГОВОРИМ И ЧТО
СДЕЛАЕМ

-
1. Обучающая и тестовая выборка, кросс-валидация: немного теории;
 2. Метрики качества: accuracy, precision, recall: определения и практическое задание;
 3. Смещение и разброс (bias-variance tradeoff): немного теории;
 4. Признаки переобучения и регуляризация: основы и практическое задание.

ОБУЧАЮЩАЯ, ТЕСТОВАЯ
ВЫБОРКА И
ПЕРЕОБУЧЕНИЕ

Обучающая выборка

Содержит значения признаков и целевой переменной.

На обучающей выборке строим модель.

ТЕСТОВАЯ ВЫБОРКА

Содержит значения признаков, по которым необходимо предсказать значение целевой переменной.

Оцениваем качество различных вариантов модели.

ПРОБЛЕМЫ

Модель может хорошо работать на обучающей выборке, однако сильно терять в качестве на тестовой (один из вариантов - переобучение).

Преобразования данных на обучающей выборке должны быть повторены и иметь смысл для тестовой.

РАЗБИВАЕМ ОБУЧАЮЩУЮ ВЫБОРКУ

Разбиваем обучающую выборку на 2 части. На одной будем тренировать модель, на другой – проверять (т. е. использовать в качестве тестовой, только с известной целевой переменной)

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split( X, y, test_size = 0.3, random_state = 0 )
```

ОБУЧАЮЩАЯ ВЫБОРКА

```
graph TD; A[ОБУЧАЮЩАЯ ВЫБОРКА] --> B[TRAINING]; A --> C[TEST];
```

TRAINING

TEST

К БЛИЖАЙШИХ СОСЕДЕЙ

НЕМНОГО ПРАКТИКИ

LOGRES_AFFAIR.ІРҮНВ

ОЦЕНКА КАЧЕСТВА МОДЕЛИ

PRECISION RECALL
ТОЧНОСТЬ И ПОЛНОТА

ПОРОГ ДЛЯ ТЕСТОВОЙ ВЫБОРКИ

```
model = LogisticRegression()
```

```
model.fit(X_train, y_train)  
predictions = model.predict_proba(X_test)
```

```
zip(predictions[:, 1], y_test)
```

```
[(0.64583193796528038, 0),  
(0.075906148028446599, 0),  
(0.2704606033743272, 0),  
(0.26938542699540474, 0),  
(0.26433391263337475, 1),  
(0.1443590034736055, 0),  
(0.17840859560894495, 0),  
(0.21871761029690232, 0),  
(0.75293068528621931, 1),  
(0.2694630112685994, 0),  
(0.11209927315788928, 0),  
(0.18717054508217956, 0),  
(0.081787486664569364, 0)].
```

Выберем порог, выше которого будем считать полученное значение принадлежащим 1. А ниже – нулю

Это определит долю угаданных моделью значений

МАТРИЦА ОШИБОК ДЛЯ ПОРОГА

	Actual positive	Actual negative
Predicted positive	True positive	False positive
Predicted negative	False negative	True negative

False positive – ошибка I рода
(ложная тревога)

False negative – ошибка II рода
(пропуск цели)

ТОЧНОСТЬ

	Actual positive	Actual negative
Predicted positive	True positive	False positive
Predicted negative	False negative	True negative

Ассурасу – доля правильно предсказанных от всех вариантов

$$Accuracy = \frac{TP + TN}{TP + FP + FN + TN}$$

К БЛИЖАЙШИХ СОСЕДЕЙ

НЕМНОГО ПОСЧИТАЕМ

LOGRES_AFFAIR.ИПУНВ

ПОЧЕМУ ТОЧНОСТИ НЕДОСТАТОЧНО

100 обычных писем

True negative 90	FP 10
---------------------	----------

10 спам-писем

false negative 5	True positive
---------------------	------------------

На почту пришло 100 обычных писем. И 10 писем спама.

Наша модель из 100 обычных 10 классифицировала как спам. Из 10 спам-писем – 5 как спам

ПОЧЕМУ ТОЧНОСТИ НЕДОСТАТОЧНО

	Actual positive	Actual negative
Predicted positive	5	5
Predicted negative	10	90

Ассурасу – доля правильно предсказанных от всех вариантов

$$Accuracy = \frac{5 + 90}{5 + 90 + 10 + 5} = 86\%$$

ПОЧЕМУ ТОЧНОСТИ НЕДОСТАТОЧНО

100 обычных писем

True negative
100

10 спам-писем

False negative
10

Возьмем модель, которая
считает все письма обычными

ПОЧЕМУ ТОЧНОСТИ НЕДОСТАТОЧНО

	Actual positive	Actual negative
Predicted positive	0	10
Predicted negative	0	100

Возьмем модель, которая считает все письма обычными

$$Accuracy = \frac{0 + 100}{0 + 100 + 0 + 10} = 91\%$$

PRECISION

	Actual positive	Actual negative
Predicted positive	True positive	False positive
Predicted negative	False negative	True negative

Precision – доля правильно предсказанных среди причисленных моделью к категории 1

$$Precision = \frac{TP}{TP + FP}$$

Способность алгоритма отличать данный класс от других классов

RECALL

	Actual positive	Actual negative
Predicted positive	True positive	False positive
Predicted negative	False negative	True negative

Recall – доля правильно предсказанных среди категории 1

$$Recall = \frac{TP}{TP + FN}$$

Синоним – True Positive Rate
(sensitivity)

Способность алгоритма обнаруживать данный класс вообще

PRECISION И RECALL ДЛЯ СПАМА

100 обычных писем

True negative
100

10 спам-писем

False negative
10

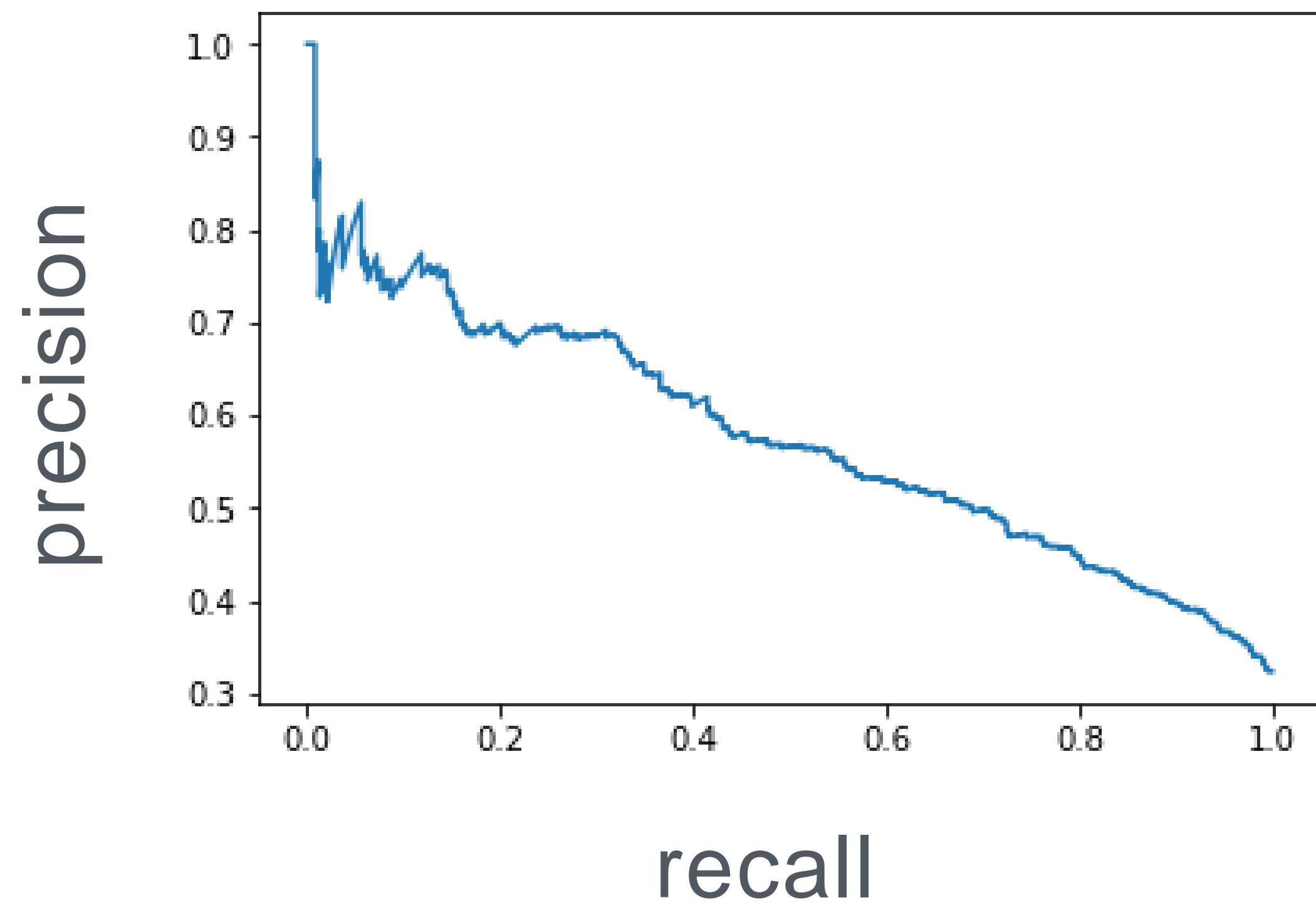
	Actual positive	Actual negative
Predicted positive	0	0
Predicted negative	10	100

К БЛИЖАЙШИХ СОСЕДЕЙ

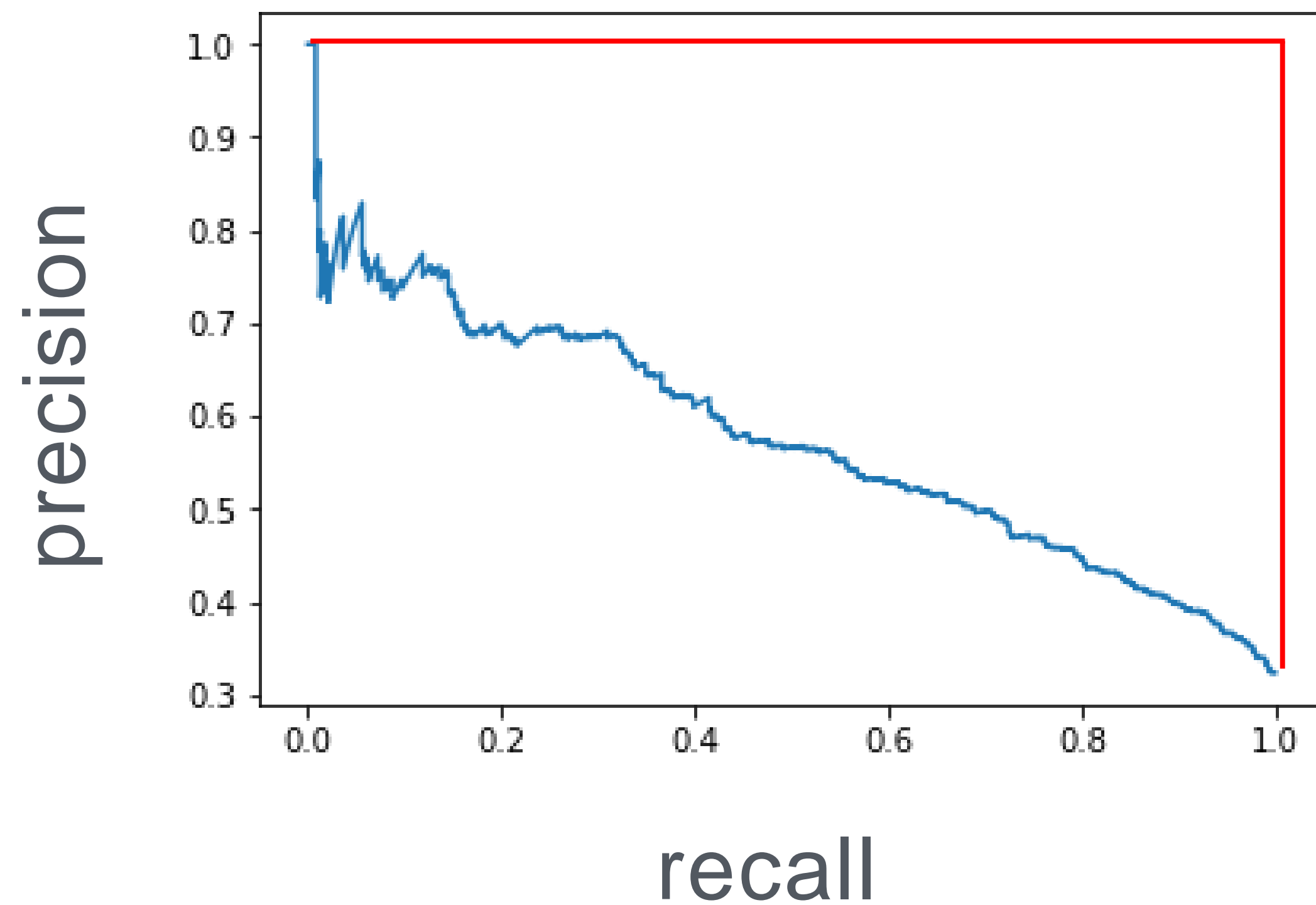
СНОВА ТОТ ЖЕ ФАЙЛ

LOGRES_AFFAIR.IPYNB

КРИВАЯ PRECISION-RECALL



КРИВАЯ PRECISION-RECALL



Модель тем лучше, чем
выше площадь под
кривой

—

AREA UNDER CURVE

TRUE POSITIVE RATE

	Actual positive	Actual negative
Predicted positive	True positive	False positive
Predicted negative	False negative	True negative

True Positive Rate – доля правильно предсказанных среди категории 1

$$TPR = \frac{TP}{TP + FN}$$

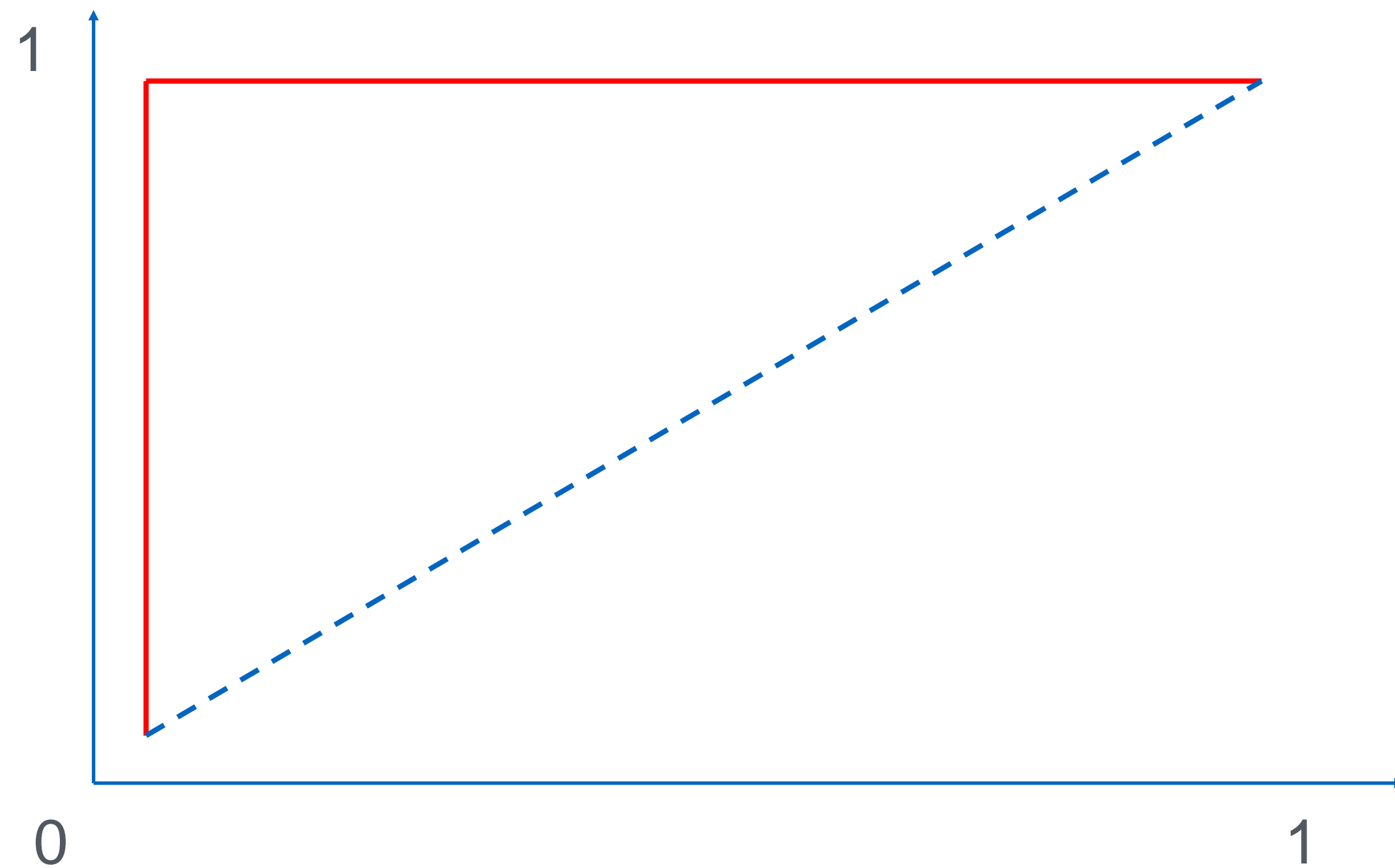
FALSE POSITIVE RATE

	Actual positive	Actual negative
Predicted positive	True positive	False positive
Predicted negative	False negative	True negative

False Positive Rate – доля неправильно предсказанных среди относящихся к категории 0

$$FPR = \frac{FP}{FP + TN}$$

ИДЕАЛЬНЫЙ СЛУЧАЙ



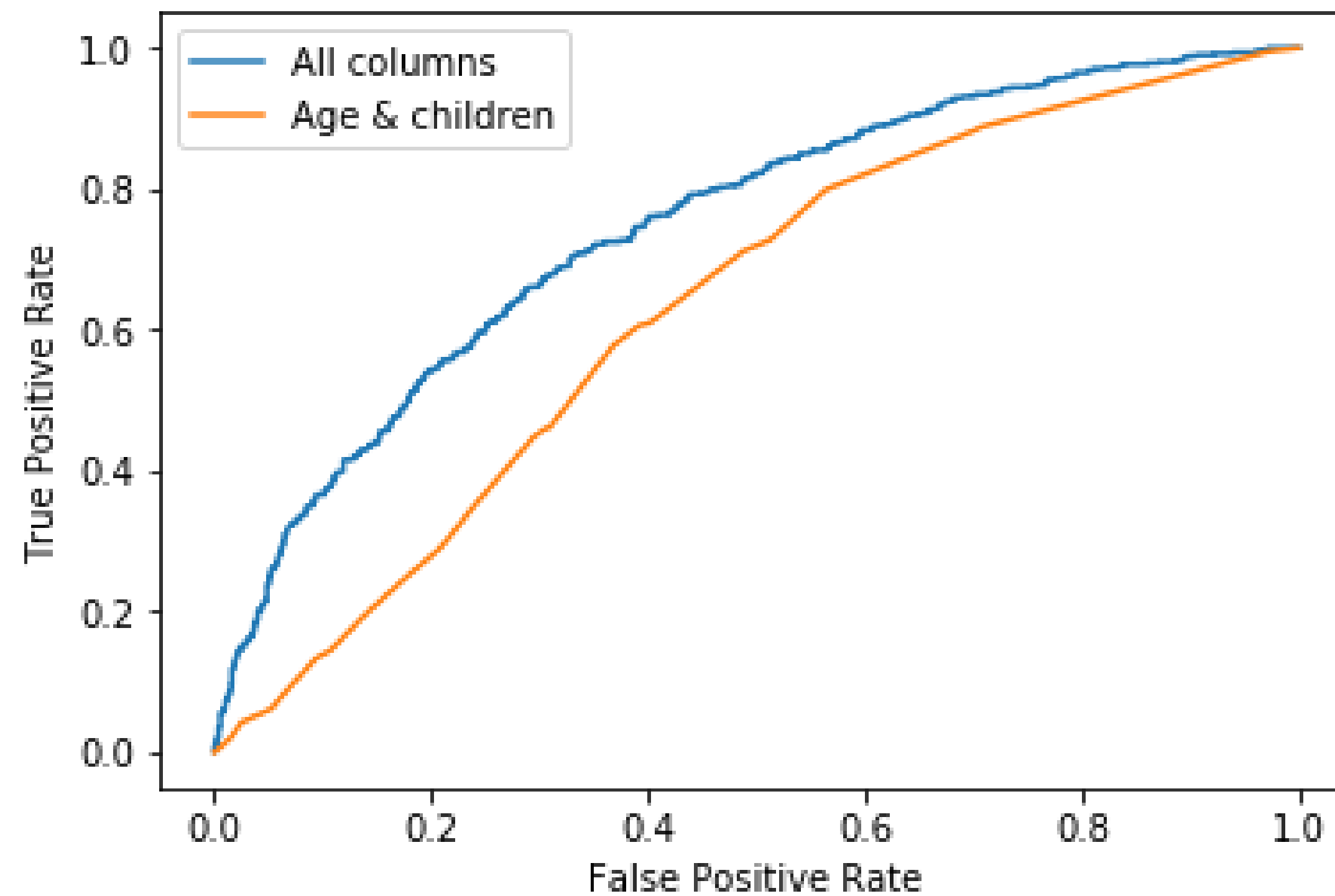
Модель предсказывает
абсолютно верно

$$\text{TPR} = 1$$

$$\text{FPR} = 0$$

----- случайные
предсказания

СРАВНЕНИЕ ДВУХ МОДЕЛЕЙ



ПРАКТИЧЕСКОЕ ЗАДАНИЕ 1

КЛАССИФИКАЦИЯ СПОРТСМЕНОВ

ATHLETES_CLASSIFIER.IPYNB

Дана статистика спортсменов ОИ 2016. Необходимо построить модель, предсказывающая пол спортсмена по имеющимся признакам (кроме столбца sex).

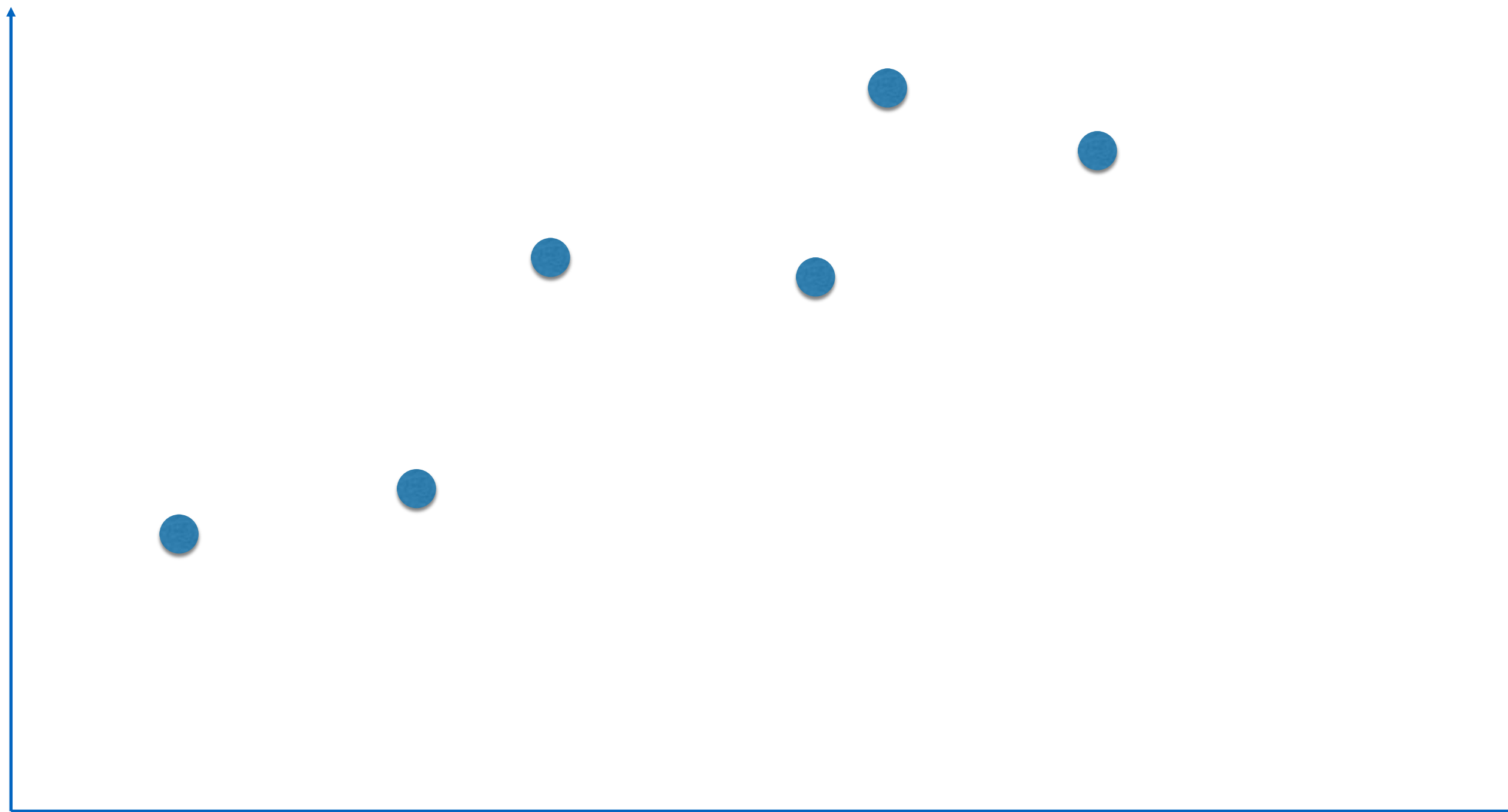
Построить графики Precision-Recall и FPR-TPR, посчитать AUC

Время на задание 20 минут

БОРЬБА С ПЕРЕОБУЧЕНИЕМ

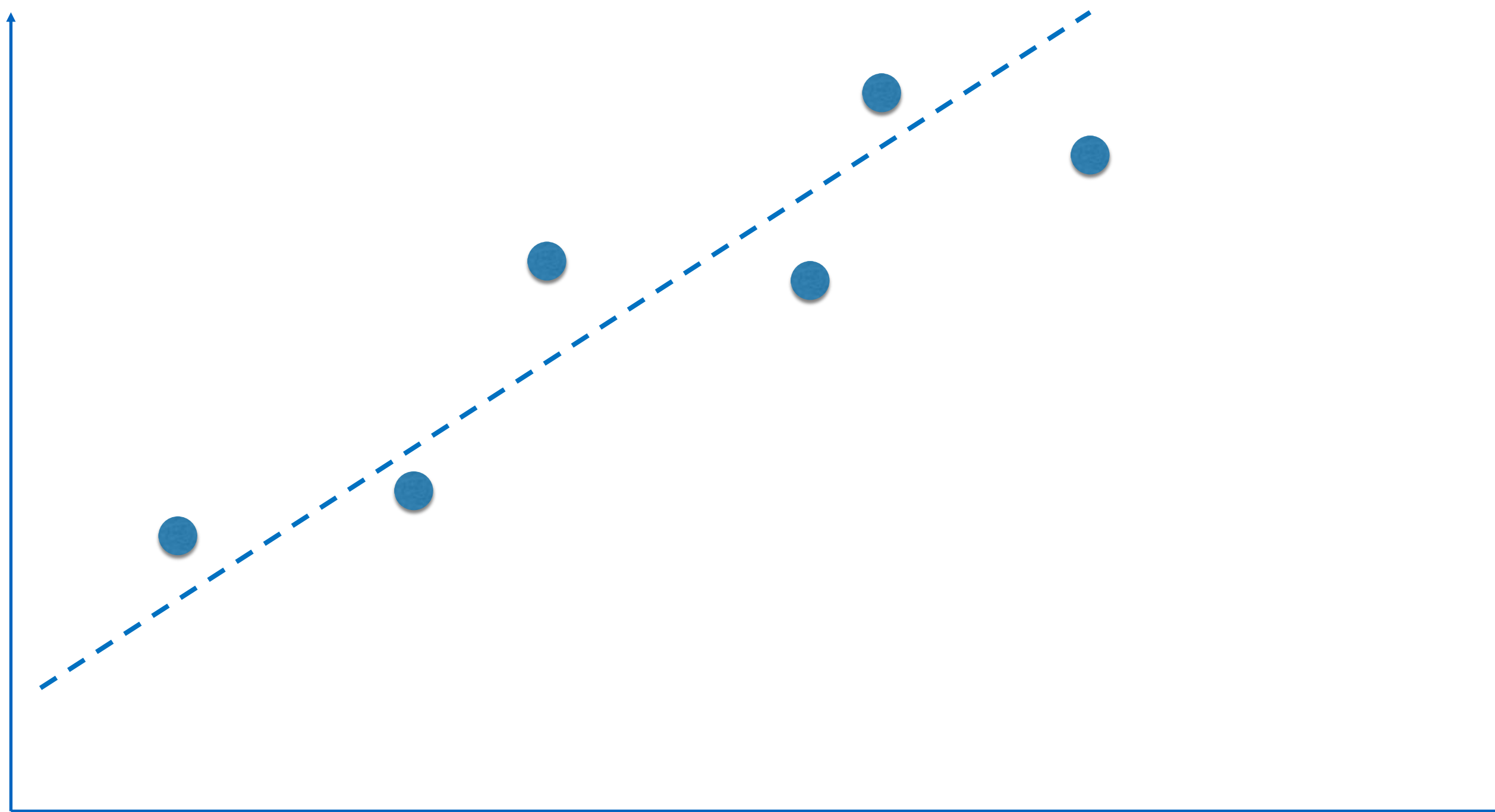
ПРИМЕР ПЕРЕОБУЧЕНИЯ

Имеются данные из 6
точек



ПРИМЕР ПЕРЕОБУЧЕНИЯ

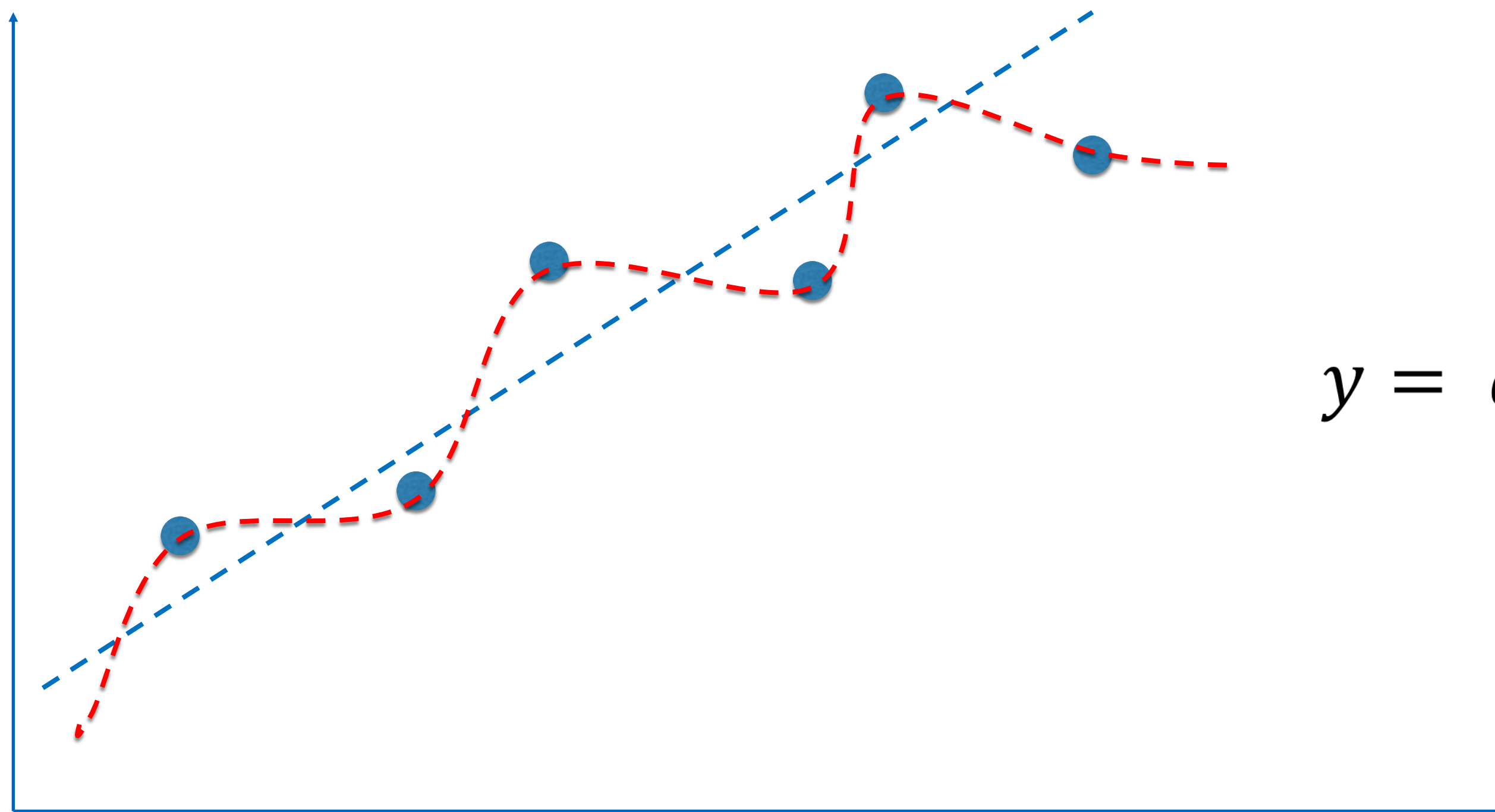
Имеются данные из 6 точек



----- $y = kx + b$; есть
ошибка > 0

ПРИМЕР ПЕРЕОБУЧЕНИЯ

Имеются данные из 6 точек



----- $y = kx + b$; есть

ошибка > 0

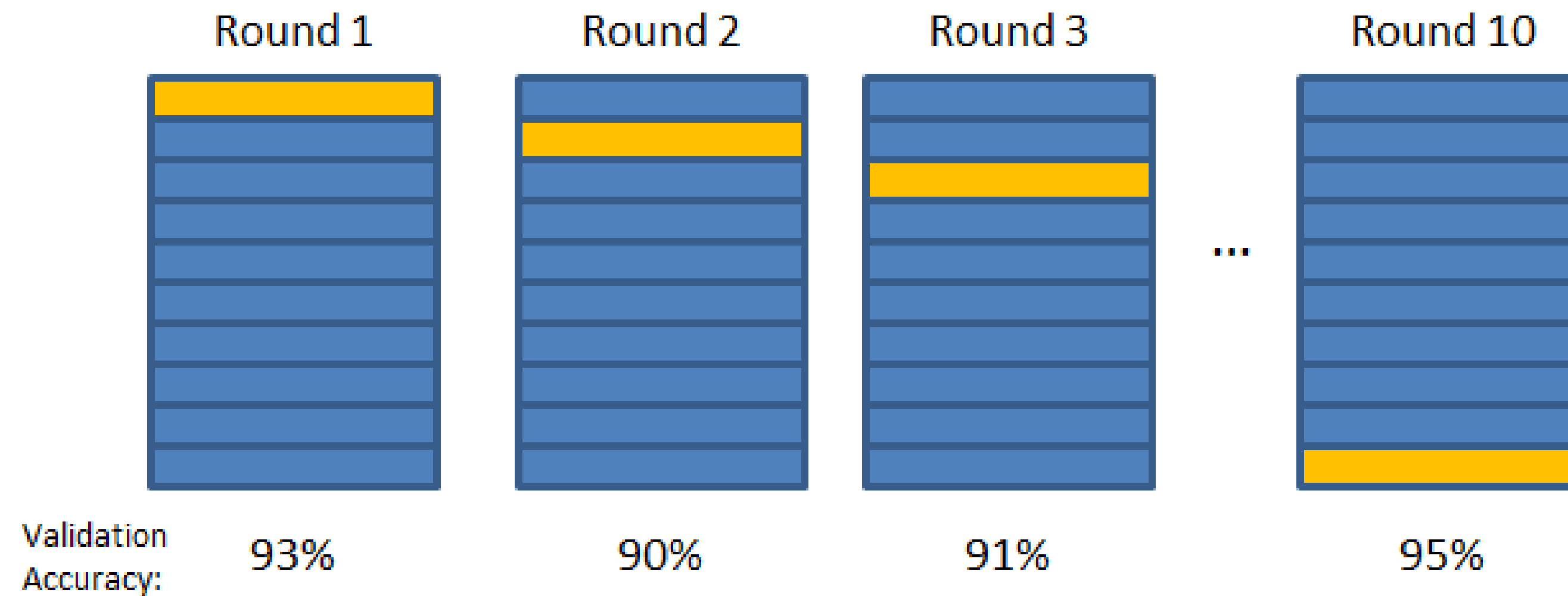
----- ошибка $= 0$. Круто?

$$y = a_0 + a_1x + a_2x^2 + a_3x^3 + a_4x^4 + a_5x^5$$

КРОСС-ВАЛИДАЦИЯ

k-fold cross validation

Validation Set
Training Set



Final Accuracy = Average(Round 1, Round 2, ...)

Лучше, чем случайная
выборка

К БЛИЖАЙШИХ СОСЕДЕЙ

CROSS_VAL_SCORE.IPYNB

ПРАКТИЧЕСКОЕ ЗАДАНИЕ 2

РАСПОЗНАВАНИЕ ЦИФР

Дана статистика картинок цифр, каждая из которых описывается набором из 64 признаков.

Используя модель `DecisionTreeClassifier`, необходимо подобрать значение параметра модели `max_depth` (от 1 до 20), при котором точность модели (accuracy) максимальна

Время на задание 20 минут

СМЕЩЕНИЕ И РАЗБРОС

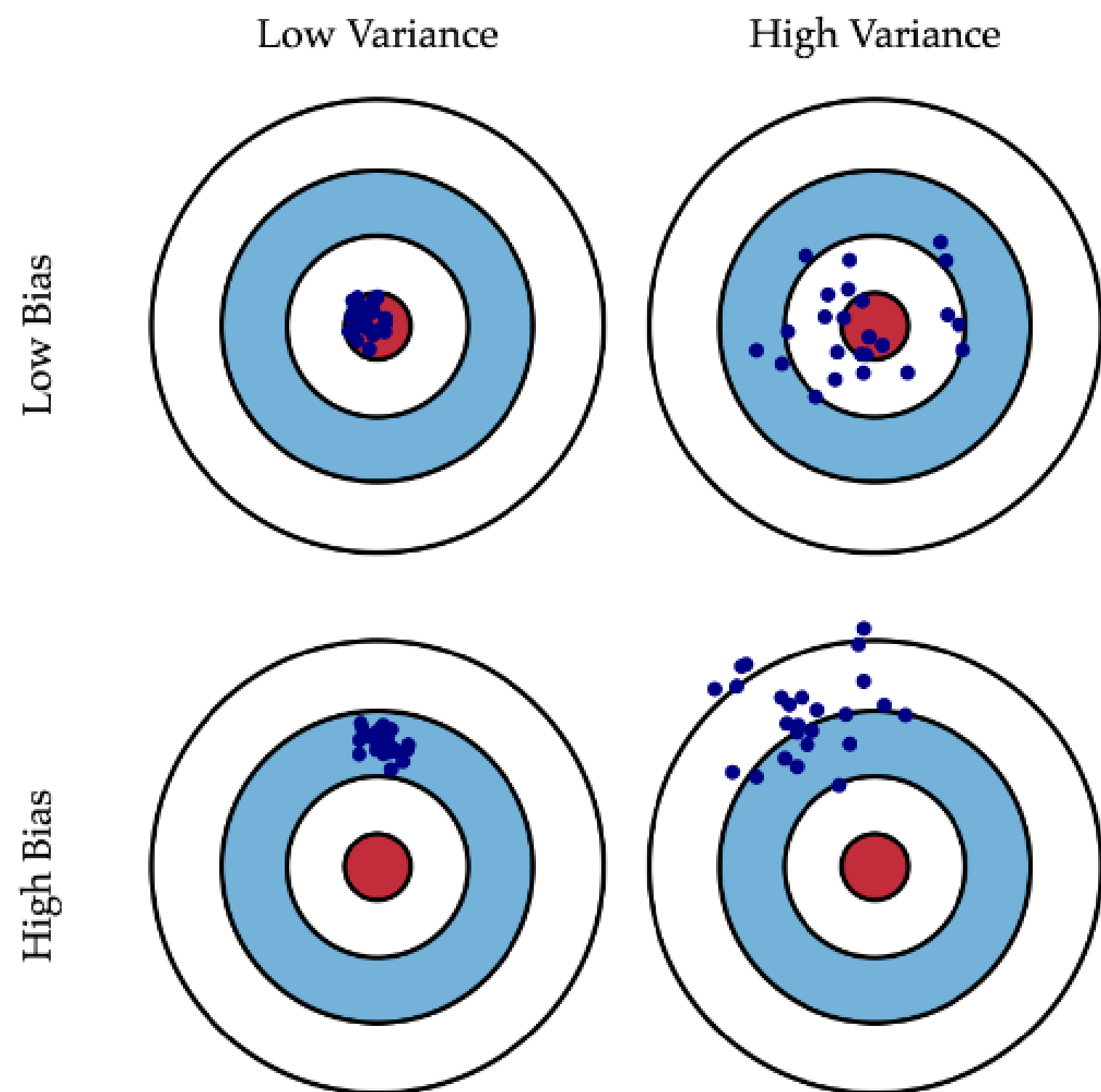
ОШИБКА ПРОГНОЗА

[HTTPS://HABRAHABR.RU/COMPANY/ODS/BLOG/323890/#RAZLOZHENIE-OSHIBKI-NA-SMESCHENIE-I-RAZBROS-BIAS-VARIANCE-DECOMPOSITION](https://habrahabr.ru/company/ods/blog/323890/#razlozhenie-oshibki-na-smeschenie-i-razbros-bias-variance-decomposition)

Можем разложить на слагаемые:

- Bias – средняя ошибка прогноза
- Variance – изменение ошибки при обучении на разных наборах данных
- Неустраняемая ошибка

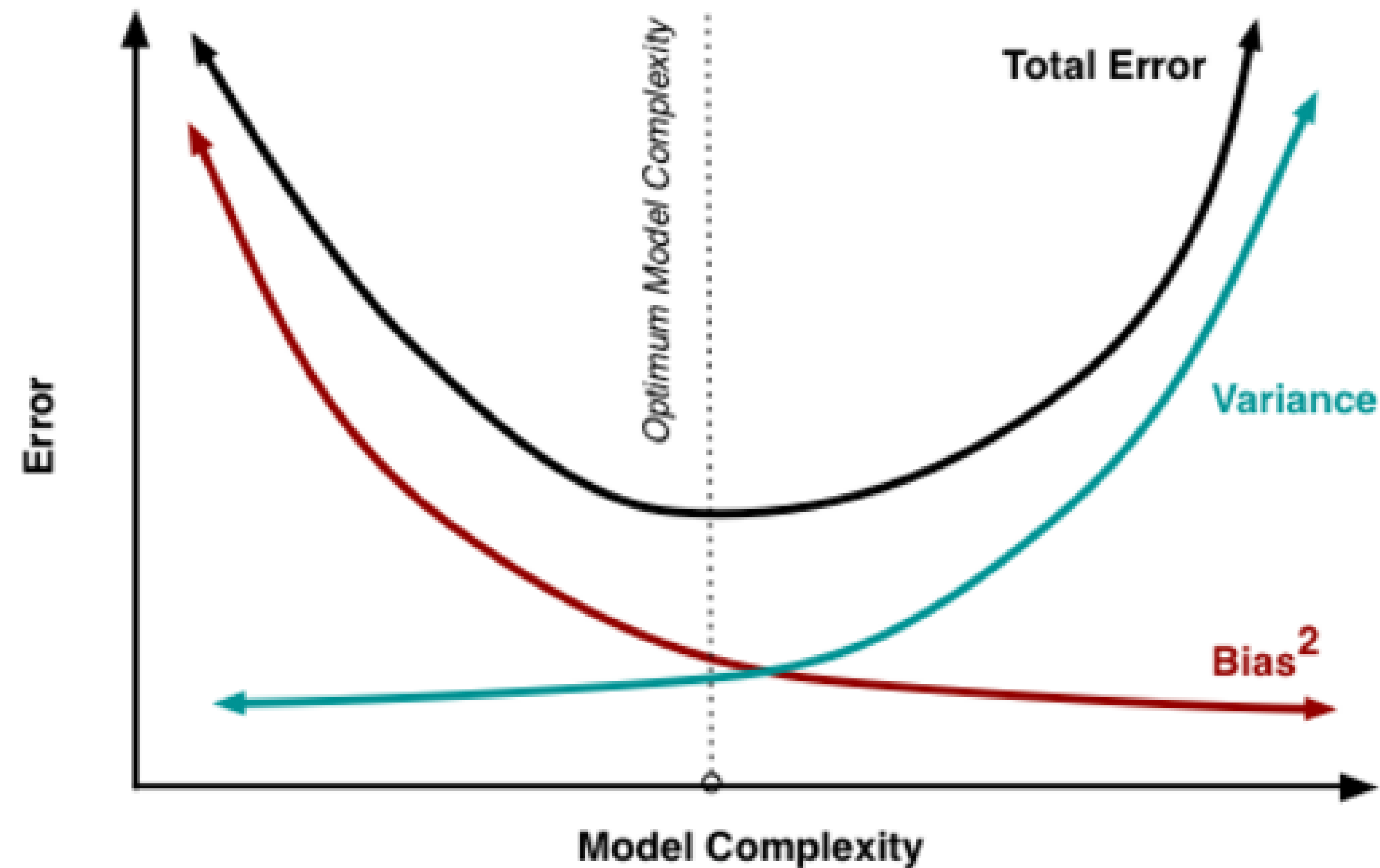
ОШИБКА ПРОГНОЗА



Сложная модель (учитывает много признаков) – увеличивает разброс ошибки

Слишком простая модель (мало признаков) – вызывает смещение в пользу одного признака

ОПТИМАЛЬНЫЙ ВАРИАНТ

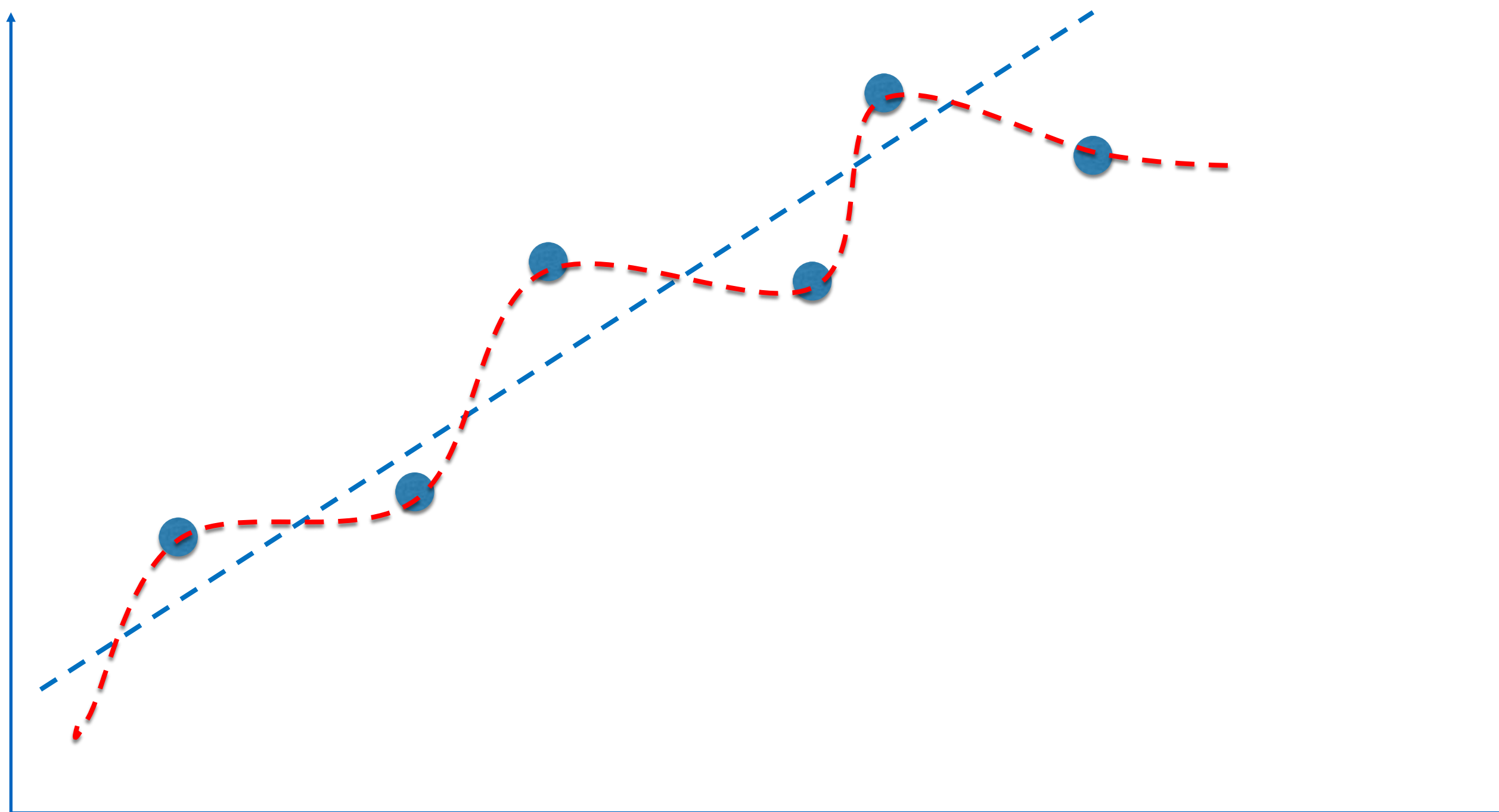


Можно ли повлиять на стабильность модели, т. е. уменьшить Variance?

—

L1 И L2 РЕГУЛЯРИЗАЦИЯ

ПРОШЛЫЙ ПРИМЕР ПЕРЕОБУЧЕНИЯ



Переберем модели,
увеличивая степень
функции

$$y = a_0 + a_1x$$

$$y = a_0 + a_1x + a_2x^2$$

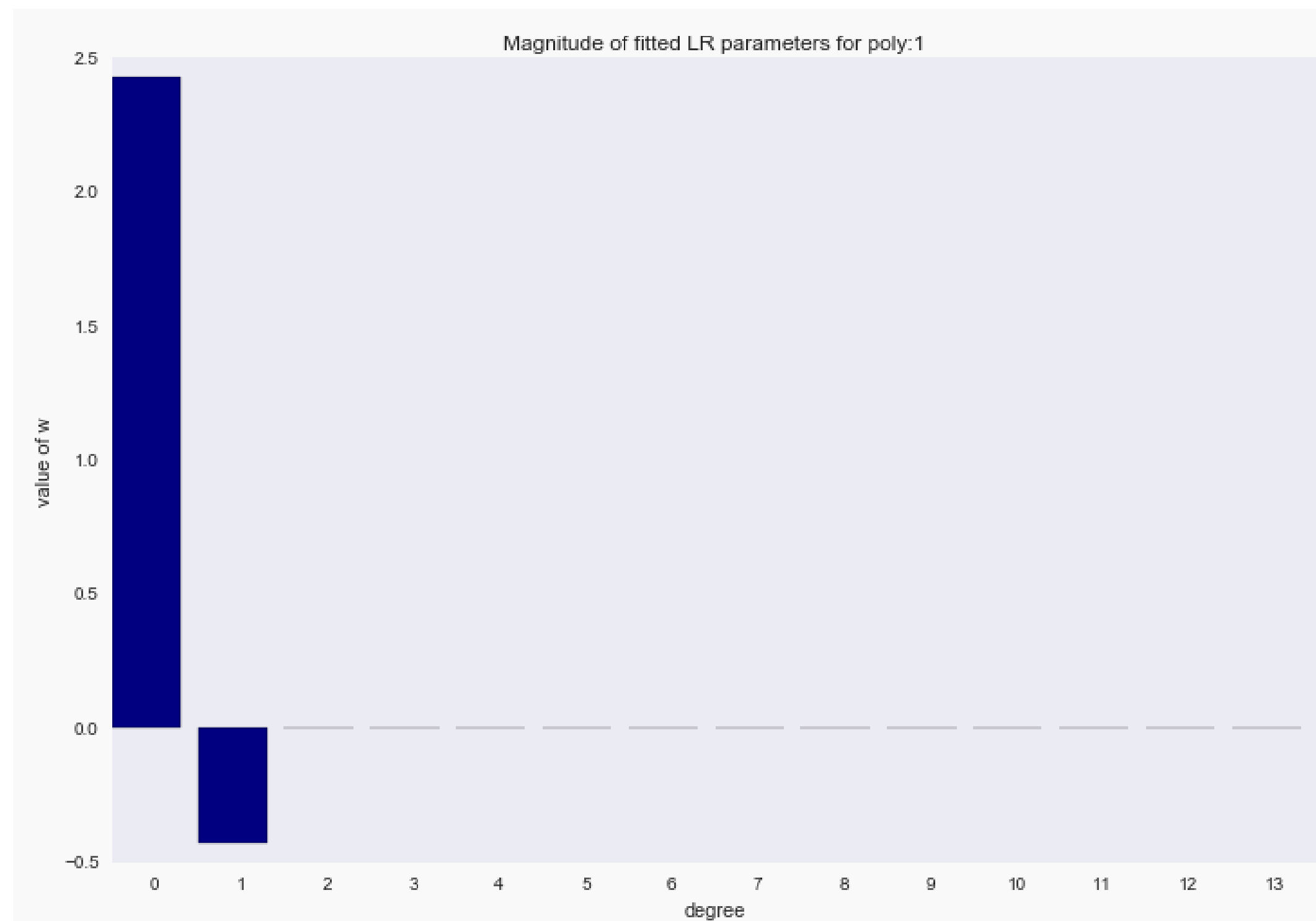
$$y = a_0 + a_1x + a_2x^2 + a_3x^3$$

...

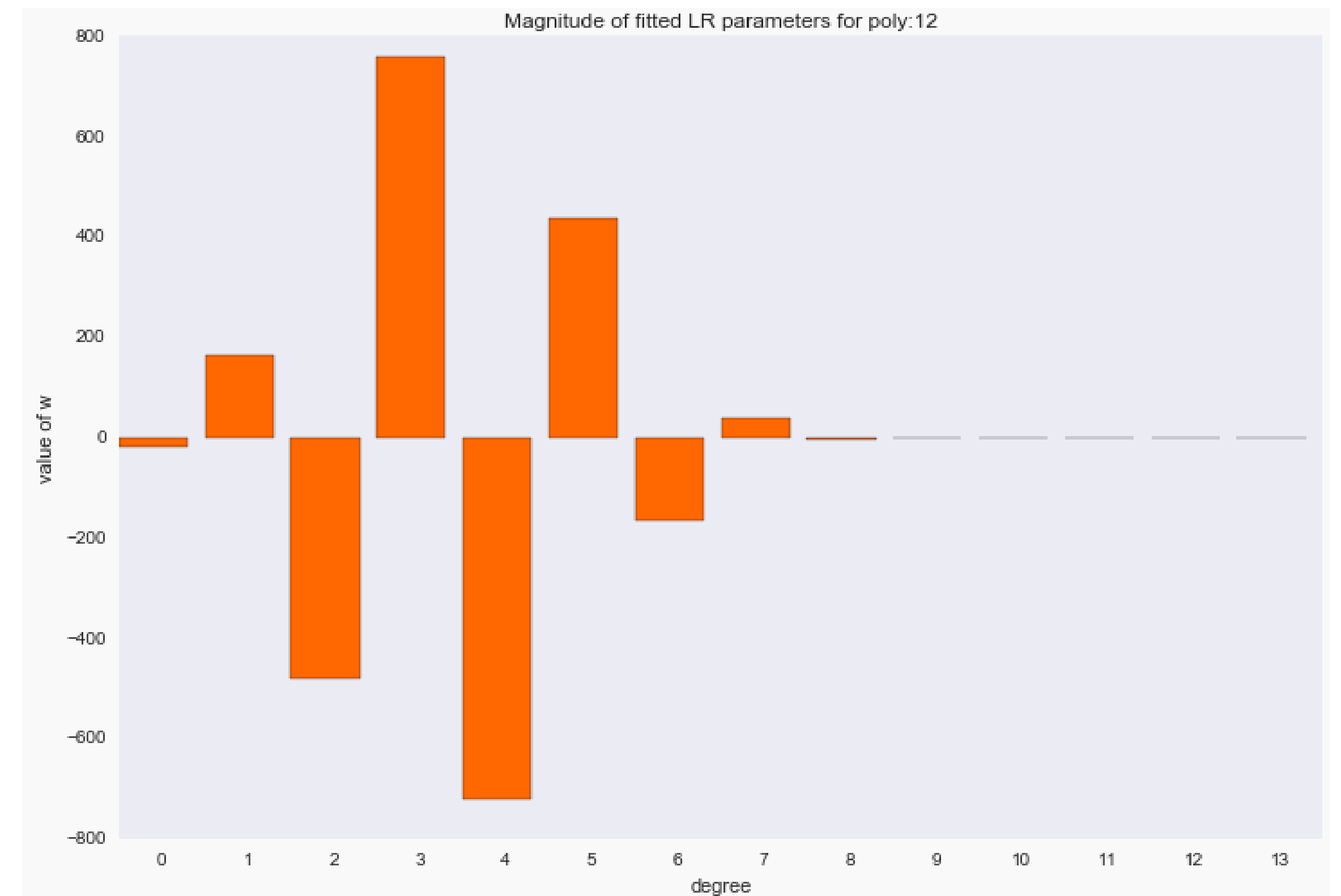
$$y = a_0 + a_1x + a_2x^2 + \dots + a_5x^5$$

КАК БУДУТ ВАРИИРОВАТЬСЯ a ?

При увеличении степени полинома вариация коэффициентов быстро растет



Картинка из блога ODS



НАДО УМЕНЬШИТЬ РАЗБРОС КОЭФФИЦИЕНТОВ

Имеем модель целевой переменной y и
коэффициентами a

$$\text{Целевая функция} = \sum_i (y_{\text{факт}} - Xa)^2$$

ШТРАФ ЗА СЛОЖНОСТЬ

Основные варианты регуляризации

$$L_1 = \sum_i (y_{\text{факт}} - Xa)^2 + \lambda \sum_i |a_i|$$

$$L_2 = \sum_i (y_{\text{факт}} - Xa)^2 + \lambda \sum_i a_i^2$$

ПРАКТИЧЕСКОЕ ЗАДАНИЕ 3

ПРЕДСКАЗАНИЕ УРОВНЯ ДОХОДА

Дана статистика пользователей `adult.csv`.

Получите значения AUC для различных моделей и их параметров

ЧТО МЫ СЕГОДНЯ УЗНАЛИ

1. Изучили метрики оценки качества моделей.
2. На практике потренировались в проведении кросс-валидации моделей.
3. Изучили признаки и способы борьбы с переобучением на примере L1 и L2 регуляризации.

ПОЛЕЗНЫЕ МАТЕРИАЛЫ

1. Наглядные примеры переобучения модели и теоретические выкладки регуляризации
<https://habrahabr.ru/company/ods/blog/322076/>
2. О разнице между L1 и L2 регуляризацией
<http://www.chioka.in/differences-between-l1-and-l2-as-loss-function-and-regularization/>
3. Более сложный пример регуляризации
<https://habrahabr.ru/company/ods/blog/323890/#3-naglyadnyy-primer-regulyarizacii-logisticheskoy-regressii>



НЕТОЛОГИЯ
групп

Спасибо за внимание!

Алексей Кузьмин



aleksej.kyzmin@gmail.com