

# ОБЗОР БИБЛИОТЕКИ МАШИННОГО ОБУЧЕНИЯ SCIKIT-LEARN

ВЯЧЕСЛАВ МУРАШКИН



# ВЯЧЕСЛАВ МУРАШКИН

Разработчик-  
исследователь, Яндекс



[mvjacheslav@gmail.com](mailto:mvjacheslav@gmail.com)



[a4tunado\\_](https://t.me/a4tunado_)

---

# ЦЕЛИ ЗАНЯТИЯ

---

## В КОНЦЕ ЗАНЯТИЯ ВЫ:

- изучите API библиотеки scikit-learn
- познакомитесь с основными модулями библиотеки scikit-learn
- научитесь применять библиотеку на реальных данных

---

О ЧЁМ ПОГОВОРИМ И ЧТО  
СДЕЛАЕМ

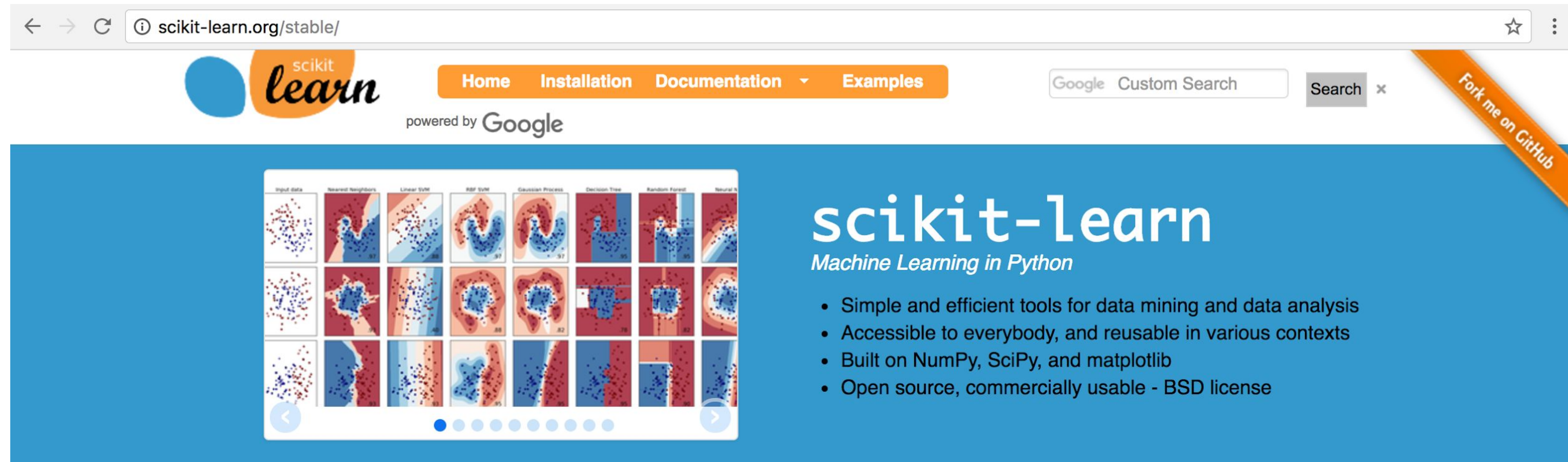
- 
1. Возможности и ограничения scikit-learn
  2. Модули библиотеки scikit-learn
  3. Модель линейной регрессии
  4. Кросс-валидация и подбор гиперпараметров модели
  5. Пример применения библиотеки на реальных данных

---

# 1. БИБЛИОТЕКА SCIKIT-LEARN



# БИБЛИОТЕКА SCIKIT-LEARN



## Classification

Identifying to which category an object belongs to.

**Applications:** Spam detection, Image recognition.

**Algorithms:** SVM, nearest neighbors, random forest, ... — Examples

## Regression

Predicting a continuous-valued attribute associated with an object.

**Applications:** Drug response, Stock prices.  
**Algorithms:** SVR, ridge regression, Lasso, ... — Examples

## Clustering

Automatic grouping of similar objects into sets.

**Applications:** Customer segmentation, Grouping experiment outcomes  
**Algorithms:** k-Means, spectral clustering, mean-shift, ... — Examples

## Dimensionality reduction

Reducing the number of random variables to consider.

**Applications:** Visualization, Increased efficiency

## Model selection

Comparing, validating and choosing parameters and models.

**Goal:** Improved accuracy via parameter tuning

## Preprocessing

Feature extraction and normalization.

**Application:** Transforming input data such as text for use with machine learning algorithms.  
**Modules:** preprocessing, feature extraction



# ВОЗМОЖНОСТИ SCIKIT-LEARN

- **Построение моделей** машинного обучения
- Реализованы модели **классификации, регрессии и кластеризации**
- Для оценки качества моделей и подбора гипер-параметров реализованы алгоритмы **предобработки данных, алгоритмы кросс-валидации**
- Предоставляет API интерфейс для **Python**

# ОСОБЕННОСТИ SCIKIT-LEARN

- Использует **методы оптимизации** из библиотеки [SciPy](#)
- Для большинства алгоритмов реализована возможность **параллельного обучения** с использованием библиотеки [Joblib](#)
- Для загрузки данных необходимо использовать **внешние библиотеки**, например [NumPy](#) или [Pandas](#)
- Данные для обучения модели должны быть загружены в **оперативную память**

---

## 2. МОДУЛИ SCIKIT-LEARN

# МОДЕЛИ МАШИННОГО ОБУЧЕНИЯ

- `sklearn.linear_model` - линейные модели классификации и регрессии
  - [LinearRegression](#)
  - [Ridge](#)
  - [LogisticRegression](#)
- `sklearn.tree` - дерево решений
  - [DecisionTreeClassifier](#)

# МОДЕЛИ МАШИННОГО ОБУЧЕНИЯ

- **sklearn.ensemble** - ансамбли решений: бустинг, лес
  - [RandomForestClassifier](#)
  - [AdaBoostClassifier](#)
  - [GradientBoostingClassifier](#)
- **sklearn.cluster** - обучение без учителя
  - [KMeans](#)
  - [DBSCAN](#)

# ОЦЕНКА КАЧЕСТВА

- **sklearn.metrics** - метрики качества алгоритмов
  - [classification\\_report](#)
  - [mean\\_squared\\_error](#)
- **sklearn.feature\_selection** - оценка важности признаков
  - [RFE](#)
- **sklearn.model\_selection** - оценка качества и подбор гипер-параметров
  - [cross\\_val\\_score](#)
  - [GridSearchCV](#)



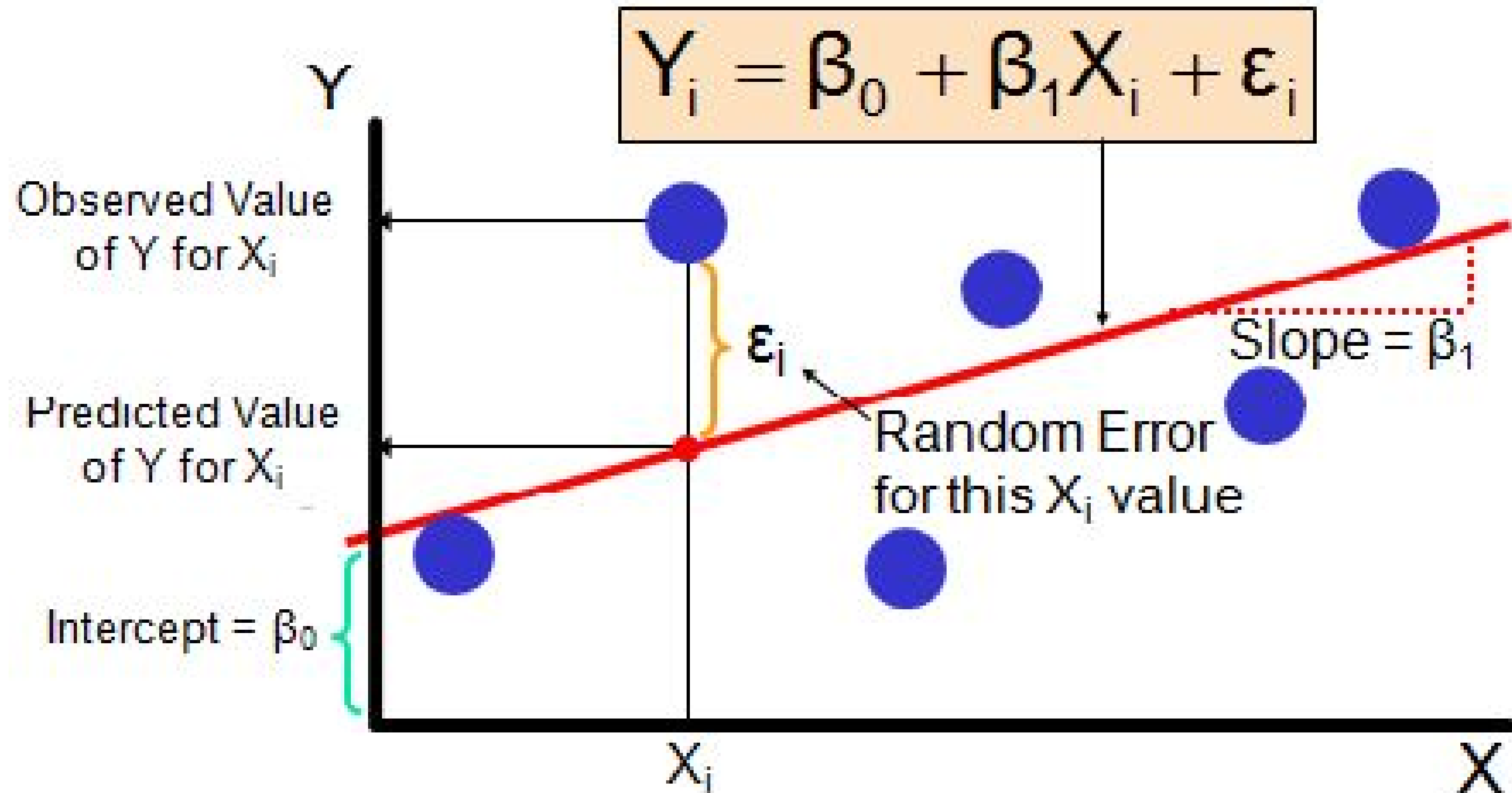
# ПРЕДОБРАБОТКА ДАННЫХ

- **sklearn.preprocessing** - нормализация, центрирование, бинаризация
  - [StandardScaler](#)
- **sklearn.feature\_extraction** - предобработка сырых данных
  - [HashingVectorizer](#)
  - [TfidfTransformer](#)
- **sklearn.decomposition** - разложение матриц и снижение размерности
  - [PCA](#)
  - [TruncatedSVD](#)

---

## 3. ЛИНЕЙНАЯ РЕГРЕССИЯ

# ЗАДАЧА



## ОЦЕНКА ПАРАМЕТРОВ

$$Y_i = \beta_0 + \beta_1 X_i \qquad SSE = \sum_i^n (Y_i - \hat{Y}_i)^2$$

$$\hat{\beta}_1 = \frac{\sum_i^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_i^n (x_i - \bar{x})^2} \qquad \hat{\beta}_0 = \bar{y} - \beta_i \bar{x}$$

## РЕГУЛЯРИЗАЦИЯ (RIDGE)

$$SSE + \lambda \sum_{j=1}^k \beta_j^2$$

$$\sum_i^n (Y_i - \beta_0 - \beta_1 X_i)^2 + \lambda \sum_{j=1}^k \beta_j^2$$

# МЕТРИКИ КАЧЕСТВА

1. Среднеквадратичная  
ошибка

$$MSE = \frac{1}{n} \sum_i^n (y_i - \hat{y}_i)^2$$

2. Коэффициент  
детерминации

$$R^2 = 1 - \frac{\sum_i^n (y_i - \hat{y}_i)^2}{\sum_i^n (y_i - \bar{y})^2}$$



---

## 4. КРОСС-ВАЛИДАЦИЯ

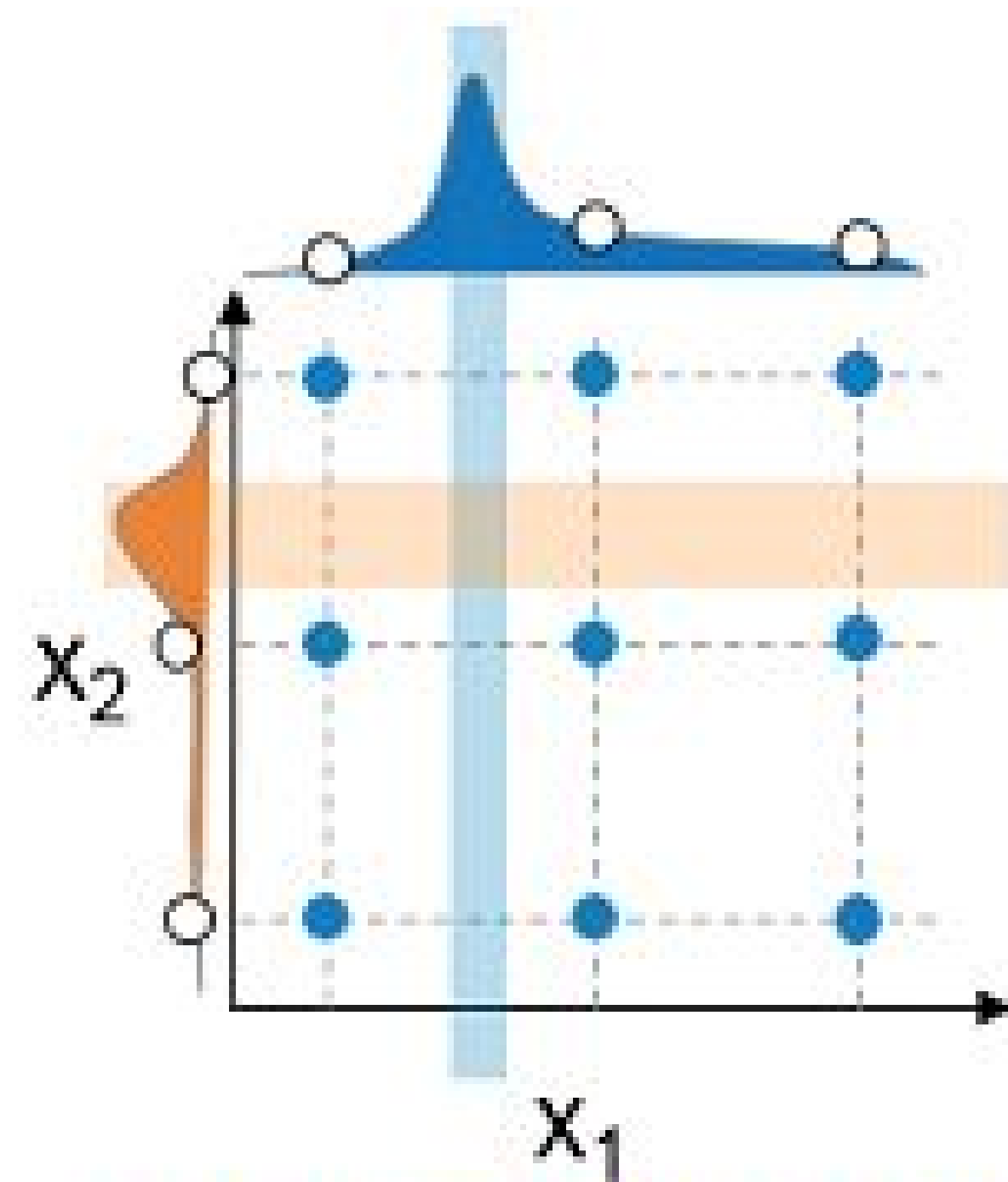
# ОБОБЩАЮЩАЯ СПОСОБНОСТЬ

- Алгоритм обучения обладает **обобщающей способностью**, если вероятность ошибки на тестовой выборке не сильно отличается от ошибки на обучающей выборке
- **Переобучение** - нежелательное явление, при котором вероятность ошибки обученного алгоритма на объектах тестовой выборки оказывается существенно выше, чем средняя ошибка на обучающей выборке
- Переобучение связано с **избыточной сложностью** используемой модели

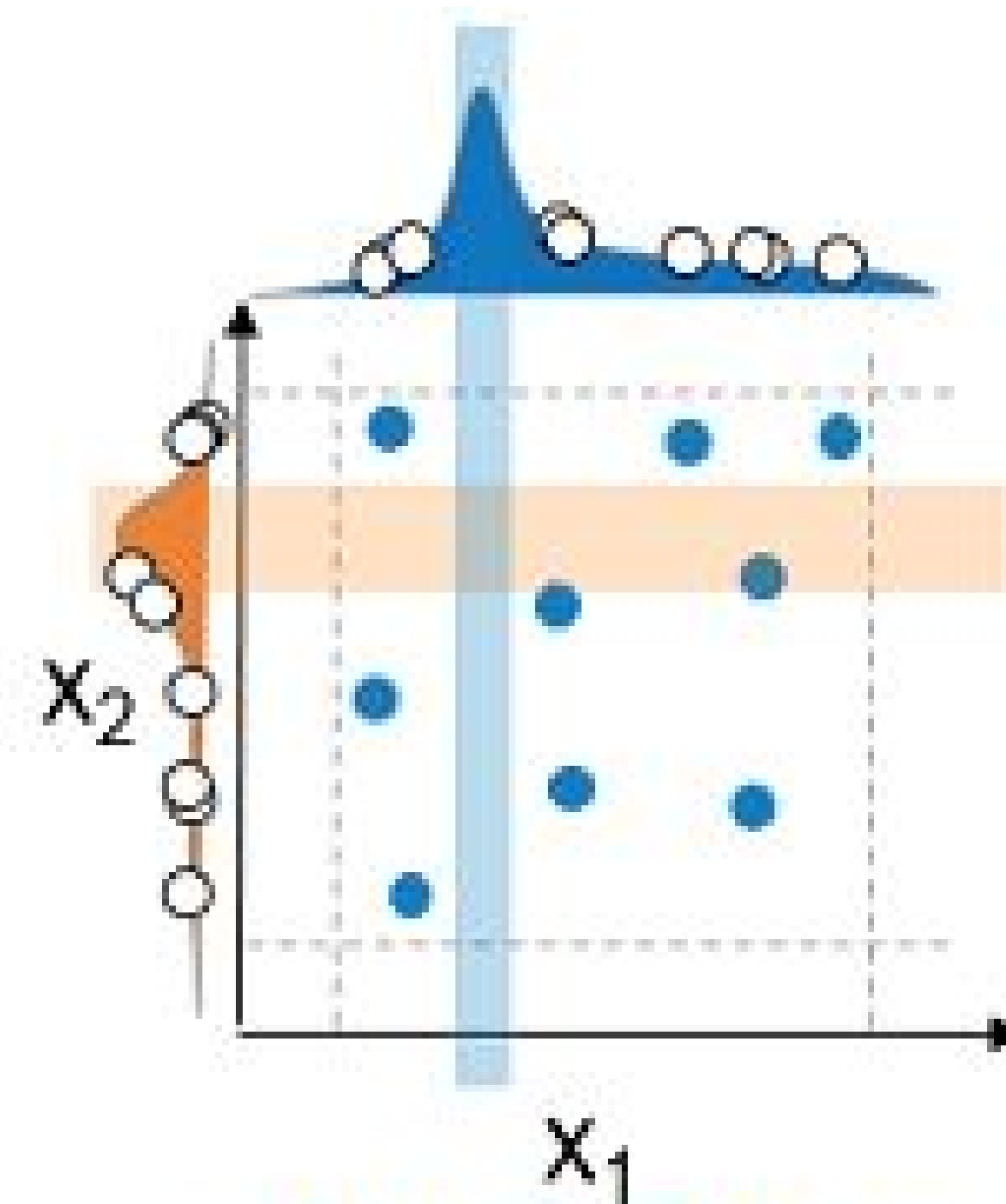
# К-FOLD КРОСС-ВАЛИДАЦИЯ



# ПОДБОР ГИПЕР-ПАРАМЕТРОВ



Standard Grid Search



Random Search

---

# ПРАКТИЧЕСКОЕ ЗАДАНИЕ

# ЗАДАНИЕ

1. Построить регрессионную модель на данных Boston Housing Dataset
2. Оценить качество модели на кросс-валидации
3. Подобрать гипер-параметры модели
4. Сравнить различные алгоритмы регрессии для данной задачи



---

ЧТО МЫ СЕГОДНЯ УЗНАЛИ

## ЧТО МЫ СЕГОДНЯ УЗНАЛИ

---

1. Познакомились с API библиотеки scikit-learn
2. Изучили возможности и особенности библиотеки
3. Научились применять scikit-learn для решения практических задач

---

ПОЛЕЗНЫЕ МАТЕРИАЛЫ

## ПОЛЕЗНЫЕ МАТЕРИАЛЫ

---

1. [An introduction to machine learning with scikit-learn](#)
2. [Scikit\\_Learn\\_Cheat\\_Sheet\\_Python.pdf](#)
3. [MachineLearning.ru](#)



НЕТОЛОГИЯ  
групп

# Спасибо за внимание!

МКРАШКИН ВЯЧЕСЛАВ



[mvjacheslav@gmail.com](mailto:mvjacheslav@gmail.com)



[a4tunado\\_](https://t.me/a4tunado_)