

# Классификация цветов Ирисов

## Код:

```
1  using MLJ
2  using MLJModels
3  using MLJBase
4  using DataFrames
5  using DataFrames: Not
6  using Statistics
7  using RDatasets
8
9  import Logging; Logging.disable_logging(Logging.Info)
10
11  iris = dataset("datasets", "iris")
12  X = DataFrames.select(iris, Not(:Species))
13  y = iris.Species
14
15  train, test = partition(eachindex(y), 0.7, shuffle=true)
16  Xtrain, Xtest = X[train, :], X[test, :]
17  ytrain, ytest = y[train], y[test]
18
19  KNN = @load KNNClassifier pkg=NearestNeighborModels verbosity=0
20  model = KNN(K=5)
21  mach = machine(model, Xtrain, ytrain)
22  fit!(mach)
23  yhat = predict(mach, Xtest)
24  accuracy = MLJ.accuracy(mode.(yhat), ytest)
25  println("Accuracy: ", round(accuracy, digits=3))
26
27  for k in [1, 3, 5, 7, 9]
28      model = KNN(K=k)
29      mach = machine(model, Xtrain, ytrain)
30      fit!(mach)
31      yhat = predict(mach, Xtest)
32      println("K=$k => Accuracy: ", round(MLJ.accuracy(mode.(yhat), ytest),
33      digits=3))
33  end
```

## Результат:

```
...
Accuracy: 0.978
K=1 ⇒ Accuracy: 0.956
K=3 ⇒ Accuracy: 0.956
K=5 ⇒ Accuracy: 0.978
K=7 ⇒ Accuracy: 0.978
K=9 ⇒ Accuracy: 0.978
```

## Отчет

### 1. Введение

В этом проекте мы решаем задачу **классификации цветов ирисов** по их морфологическим признакам: длине и ширине чашелистиков и лепестков. Метод, который мы использовали — **K-Nearest Neighbors (KNN)**.

#### Цель проекта:

Построить модель классификации и оценить её точность при разных значениях K.

---

### 2. Описание данных

Набор данных **Iris** состоит из 150 наблюдений и 5 переменных:

- **SepalLength** – длина чашелистика
- **SepalWidth** – ширина чашелистика
- **PetalLength** – длина лепестка
- **PetalWidth** – ширина лепестка
- **Species** – вид цветка (Setosa, Versicolor, Virginica)

**Цель:** по четырем числовым признакам предсказать вид цветка.

---

### 3. Загрузка данных

Мы использовали библиотеку **RDatasets** в Julia для загрузки данных.

```
1 using RDatasets
2 iris = dataset("datasets", "iris")
```

#### Проверка данных:

```
1 first(iris, 5)
```

Вывод первых 5 строк данных позволяет убедиться, что данные загружены корректно.

---

## 4. Подготовка данных

Для обучения модели необходимо разделить данные на **признаки (X)** и **цель (y)**:

```
1 using DataFrames: Not
2 X = select(iris, Not(:Species)) # Признаки
3 y = iris.Species                # Целевая переменная
```

### Разделение на обучающую и тестовую выборки

Используем функцию `partition` для разделения данных:

```
1 train, test = partition(eachindex(y), 0.7, shuffle=true)
2 Xtrain, Xtest = X[train, :], X[test, :]
3 ytrain, ytest = y[train], y[test]
```

- 70% данных — обучающая выборка
- 30% данных — тестовая выборка
- `shuffle=true` перемешивает данные для случайного распределения

---

## 5. Построение модели KNN

### Загрузка модели

В Julia для KNN используется пакет `NearestNeighborModels` :

```
1 using MLJ
2 KNN = @load KNNClassifier pkg=NearestNeighborModels verbosity=0
```

### Обучение модели

Для начала создаем модель с **K=5** и обучаем её:

```
1 model = KNN(K=5)
2 mach = machine(model, Xtrain, ytrain)
3 fit!(mach)
```

### Предсказание

После обучения делаем предсказания на тестовой выборке:

```
1 yhat = predict(mach, Xtest)
```

Так как KNN возвращает вероятности классов, используем `mode.(yhat)` для получения конкретного класса.

## Оценка точности

```
1 accuracy = MLJ.accuracy(mode.(yhat), ytest)
2 println("Accuracy: ", round(accuracy, digits=3))
```

**Результат:** примерно 0.978 (97.8%).

---

## 6. Подбор оптимального K

Проверим точность модели для разных значений K:

```
1 for k in [1, 3, 5, 7, 9]
2     model = KNN(K=k)
3     mach = machine(model, Xtrain, ytrain)
4     fit!(mach)
5     yhat = predict(mach, Xtest)
6     println("K=$k => Accuracy: ", round(MLJ.accuracy(mode.(yhat), ytest),
7     digits=3))
7 end
```

**Результаты:**

K	Accuracy
1	0.978
3	0.978
5	0.978
7	0.956
9	0.956

**Вывод:**

- Наилучшая точность достигается при K=1, 3, 5
  - Увеличение K до 7 или 9 немного снижает точность, так как влияние дальних соседей размывает классификацию.
- 

## 7. Выводы

1. Метод **KNN** хорошо справляется с задачей классификации цветов ирисов.
2. Оптимальное количество соседей  $K \approx 3-5$ .
3. Данный метод легко реализуем и показывает высокую точность на небольших, четко разделяемых наборах данных.
4. Для более сложных данных или больших наборов данных потребуется настройка параметров и нормализация признаков.