

ТЕХНОЛОГИЧЕСКАЯ КАРТА ЗАНЯТИЯ

Тема занятия: Практика решающего леса

Аннотация к занятию: на данном занятии обучающиеся познакомятся с алгоритмом машинного обучения — случайным лесом. Рассмотрят на примере реальных данных качество алгоритма в зависимости от количества деревьев.

Цель занятия: сформировать у обучающихся представление об алгоритме машинного обучения — случайном лесе.

Задачи занятия:

- познакомить обучающихся с алгоритмом машинного обучения — случайным лесом;
- рассмотреть, как выглядят разделяющие поверхности в случае леса и одного дерева;
- применить полученные знания на практике.

Ход занятия

Этап занятия	Время	Деятельность педагога	Комментарии, рекомендации для педагогов
Организационный этап	2 мин.	Добрый день! Добро пожаловать на урок.	Приветствие. Создание в классе атмосферы психологического комфорта
Постановка цели и задач занятия. Мотивация учебной деятельности обучающихся	10 мин.	<p>Вопрос для обсуждения Как вы думаете, если в процессе прогнозирования участвует множество деревьев решений, данный метод будет точнее?</p> <p>Ответы обучающихся На этом занятии мы рассмотрим конвейер машинного обучения для алгоритма случайного леса. Мы загрузим и исследуем данные. На синтетическом примере визуализируем, как выглядят разделяющие поверхности в дереве и случайном лесу.</p>	Способствовать обсуждению мотивационных вопросов
Изучение нового материала	50 мин.	Импортируем библиотеки, которые будут использованы в ходе анализа	<p>Для справки: Файл для работы: https://drive.google.com/file/d/10VYvu1lkYtCuEFLtw</p>


```

1 import pandas as pd
2 from matplotlib.colors import ListedColormap
3 import numpy as np
4 from sklearn import datasets
5 import matplotlib.pyplot as plt
6 from sklearn.tree import DecisionTreeClassifier
7 from sklearn.ensemble import RandomForestClassifier
8
9 cmap_light = ListedColormap(['#FFAAAA', '#AAFFAA', '#AAAAFF'])

```

Сравнение дерева и случайного леса

Сравним, как происходит разделение классов деревом и решающим лесом. Для этого генерируем данные, которые использовались для теоретического занятия. Не будем в деталях останавливаться на механике реализации, её суть — создание точек 3 классов по 3 спиралям, закручивающимся к центру.

 Код генерации датасета с занятия

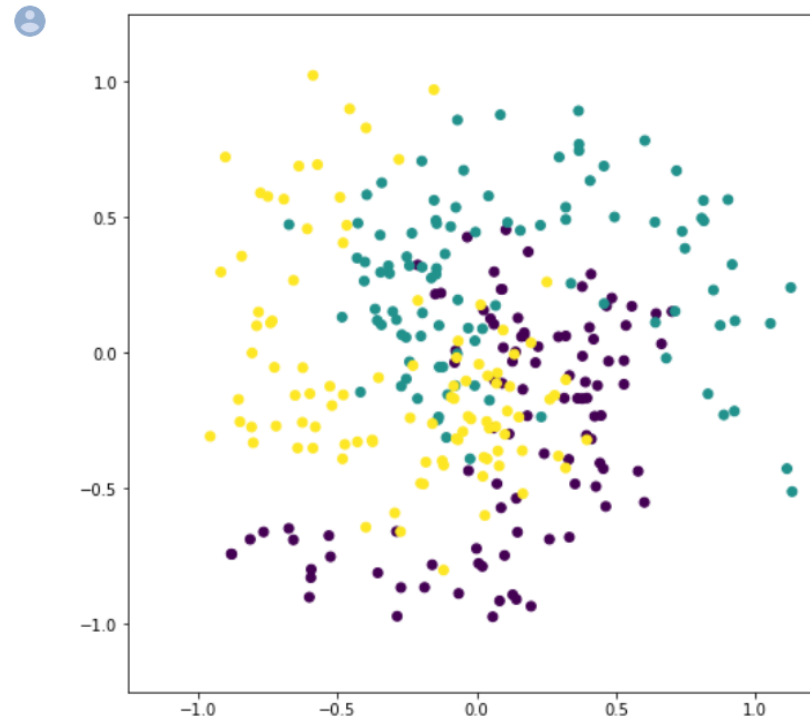
[Показать код](#)

```
[ ] 1 X, y = devil(100)
```

Визуализируем эти данные, чтобы убедиться, что всё отработало корректно:

7G0NO2hT4b39oET/view?usp=sharing

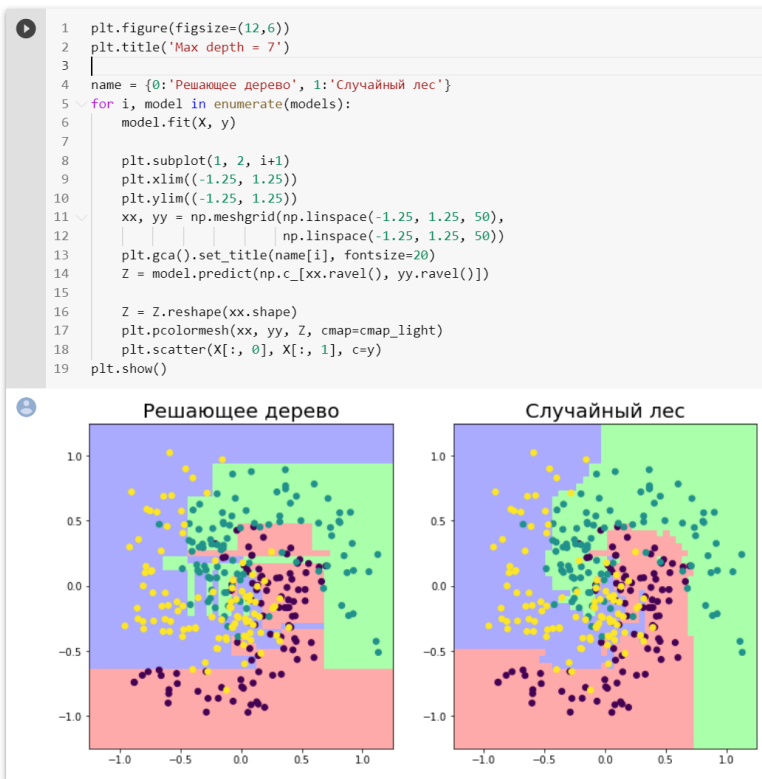
```
1 plt.figure(figsize=(8, 8))
2
3 plt.scatter(X[:,0], X[:,1], c=y)
4
5 plt.xlim([-1.25, 1.25])
6 plt.ylim([-1.25, 1.25])
7
8 plt.show()
```



Инициализируем две модели — решающего дерева и леса — с одинаковой глубиной.

```
1 models = [DecisionTreeClassifier(max_depth=10),
2           RandomForestClassifier(max_depth=10)]
```

Создадим сетку, покрывающую диапазон данных по x , y .
Обучим модели на созданных выше данных. Подскажем отношение класса к каждому элементу сетки.
Визуализируем полученный прогноз для дерева и для леса.
Видно, что границы разделяющего леса более плавные, что будет позитивно сказываться на метриках качества в дальнейшем.



Данные

В ходе анализа данных будет использован датасет мобильных телефонов:
<https://www.kaggle.com/datasets/iabhishekoofficial/mobile-price-classification?resource=download&select=train.csv>

Подобные данные могут использоваться для создания автоматических подборок телефонов в зависимости от их характеристик ещё до анонса цены.

Наши данные были скачаны с платформы Kaggle и загружены на облачный сервер. Загрузим и проверим корректность данных. Видно, что все данные были загружены в свои колонки и явных сбоев с кодировкой нет.

```
1 df = pd.read_csv('https://dl.uploadgram.me/62def1f38f9cbh?raw')
```

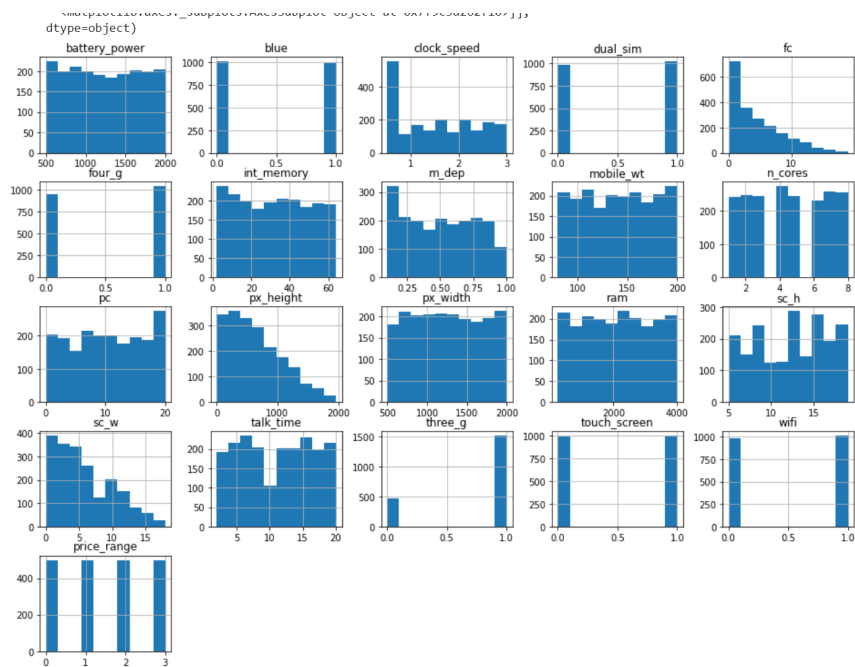
```
1 df.head()
```

	battery_power	blue	clock_speed	dual_sim	fc	four_g	int_memory	m_dep	m
0	842	0	2.2	0	1	0	7	0.6	
1	1021	1	0.5	1	0	1	53	0.7	
2	563	1	0.5	1	2	1	41	0.9	
3	615	1	2.5	0	0	0	10	0.8	
4	1821	1	1.2	0	13	1	44	0.6	

5 rows x 10 columns

Исследование

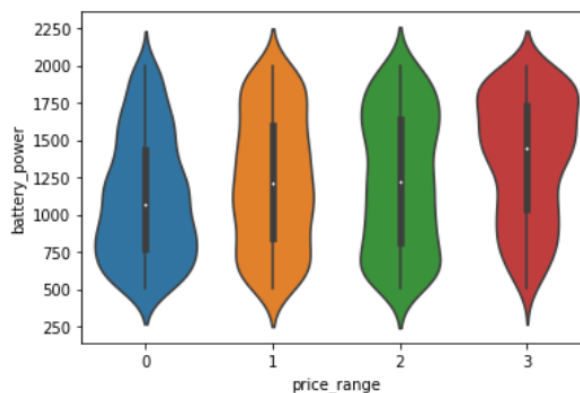
С помощью встроенного метода .hist визуализируем гистограммы по характеристикам телефонов. Можно видеть, что особенно аномальных значений нет.



Помимо гистограмм, можно визуализировать связь характеристик с ценой. Например, медианный объем батареи растёт с ростом класса цены.

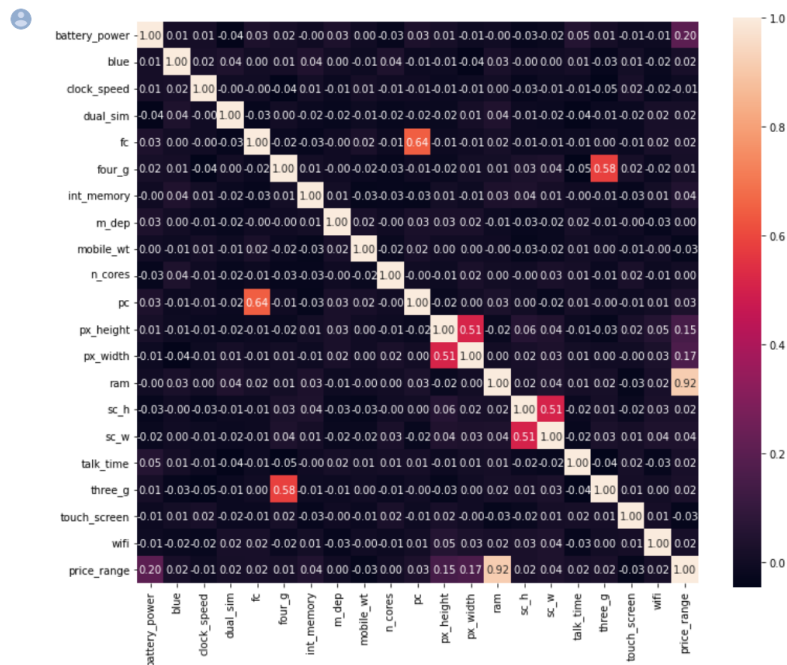

```
1 import seaborn as sns
2
3 sns.violinplot(x="price_range", y='battery_power', data=df)
```

<matplotlib.axes._subplots.AxesSubplot at 0x7f9c3c609fd0>



Поскольку категориальных признаков в датасете нет, этапа кодирования тоже нет. Визуализируем гистограмму корреляций по Пирсону (мера линейной связи). Видно, например, что оперативная память линейно коррелирует с ценовой категорией.

```
1 import matplotlib.pyplot as plt
2
3 plt.subplots(figsize=(12, 10))
4 sns.heatmap(df.corr(),fmt='.2f', square = True, annot=True)
5 plt.show()
```



Алгоритм

Разделим данные на train/test в соотношении 80 на 20.
Инициализируем алгоритм леса и обучим его на части train.
Проверим качество на test: усреднённый f1 в 0.86

показывает неплохие результаты алгоритма даже без перебора его параметров.

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.ensemble import RandomForestClassifier
3
4 X_train, X_test, y_train, y_test = train_test_split(df.drop(columns = ['price_range']), df['price_range'], test_size=0.2)
5
6 # создание и обучение алгоритма
7 rf = RandomForestClassifier()
8 rf.fit(X_train, y_train)
```

```
RandomForestClassifier()
```

```
[ ] 1 from sklearn.metrics import classification_report
2
3 preds = rf.predict(X_test)
4 print(classification_report(y_test, preds))
```

	precision	recall	f1-score	support
0	0.92	0.93	0.92	105
1	0.83	0.81	0.82	106
2	0.78	0.81	0.80	96
3	0.92	0.88	0.90	93
accuracy			0.86	400
macro avg	0.86	0.86	0.86	400
weighted avg	0.86	0.86	0.86	400

Попробуем увеличить количество решающих деревьев и критерий разделения. Видим, что это даёт прирост в два процентных пункта.

```
1 rf = RandomForestClassifier(n_estimators=1000, criterion='entropy')
2 rf.fit(X_train, y_train)
```

```
RandomForestClassifier(criterion='entropy', n_estimators=1000)
```

```
1 preds = rf.predict(X_test)
2 print(classification_report(y_test, preds))
```

```

              precision    recall  f1-score   support

    0         0.93        0.95        0.94         105
    1         0.86        0.87        0.86         106
    2         0.82        0.81        0.82          96
    3         0.91        0.89        0.90          93

 accuracy          0.88
 macro avg          0.88
weighted avg          0.88
```

Для наглядности рассмотрим несколько примеров из тестовой части и спрогнозируем их. Видно, что прогнозы полностью совпадают с разметкой.

```
[ ] 1 X_test.iloc[:5]
```


	battery_power	blue	clock_speed	dual_sim
1972	1191	0	0.8	(
859	623	0	2.0	(
384	625	1	1.9	(
1053	623	1	0.9	(
297	1329	1	1.0	(

```
1 print('Предсказания:')
2 print(rf.predict(X_test.iloc[0:5,].values))
3 print('Реальные ответ:')
4 y_test[:5].values
```

Предсказания:
[0 2 3 3 1]
Реальные ответ:
/usr/local/lib/python3.7/dist-packages/sklearn/l
"X does not have valid feature names, but"
array([0, 2, 2, 3, 1])

Сравним с алгоритмом логистической регрессии и решающего дерева. В обоих случаях результаты макро-f1 хуже.

```
1 from sklearn.linear_model import LogisticRegression
2
3 tree = DecisionTreeClassifier()
4 tree.fit(X_train, y_train)
5
6 lr = LogisticRegression()
7 lr.fit(X_train, y_train)
```

 /usr/local/lib/python3.7/dist-packages/sklearn/linear_model
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the d
<https://scikit-learn.org/stable/modules/preprocessing.h>
Please also refer to the documentation for alternative solv
https://scikit-learn.org/stable/modules/linear_model.ht
extra_warning_msg=_LOGISTIC_SOLVER_CONVERGENCE_MSG,
LogisticRegression()

```
1 preds = tree.predict(X_test)
2 print(classification_report(y_test, preds))
```

	precision	recall	f1-score	support
0	0.90	0.90	0.90	105
1	0.80	0.85	0.83	106
2	0.80	0.78	0.79	96
3	0.90	0.86	0.88	93
accuracy			0.85	400
macro avg	0.85	0.85	0.85	400
weighted avg	0.85	0.85	0.85	400

		<div><div><div>▶</div><div><pre>1 preds = lr.predict(X_test) 2 print(classification_report(y_test, preds))</pre></div></div><div><div><div>👤</div><table><thead><tr><th></th><th>precision</th><th>recall</th><th>f1-score</th><th>support</th></tr></thead><tbody><tr><td>0</td><td>0.84</td><td>0.74</td><td>0.79</td><td>105</td></tr><tr><td>1</td><td>0.57</td><td>0.62</td><td>0.60</td><td>106</td></tr><tr><td>2</td><td>0.52</td><td>0.46</td><td>0.49</td><td>96</td></tr><tr><td>3</td><td>0.68</td><td>0.78</td><td>0.73</td><td>93</td></tr><tr><td>accuracy</td><td></td><td></td><td>0.65</td><td>400</td></tr><tr><td>macro avg</td><td>0.65</td><td>0.65</td><td>0.65</td><td>400</td></tr><tr><td>weighted avg</td><td>0.66</td><td>0.65</td><td>0.65</td><td>400</td></tr></tbody></table></div></div></div>		precision	recall	f1-score	support	0	0.84	0.74	0.79	105	1	0.57	0.62	0.60	106	2	0.52	0.46	0.49	96	3	0.68	0.78	0.73	93	accuracy			0.65	400	macro avg	0.65	0.65	0.65	400	weighted avg	0.66	0.65	0.65	400	
	precision	recall	f1-score	support																																							
0	0.84	0.74	0.79	105																																							
1	0.57	0.62	0.60	106																																							
2	0.52	0.46	0.49	96																																							
3	0.68	0.78	0.73	93																																							
accuracy			0.65	400																																							
macro avg	0.65	0.65	0.65	400																																							
weighted avg	0.66	0.65	0.65	400																																							
Закрепление изученного материала	15 мин.	<div><div><div>Вопросы для обсуждения</div><div><ul style="list-style-type: none">В чём заключается смысл алгоритма случайного леса?При каком алгоритме получается более точный результат?</div></div></div>	Педагог организует беседу по вопросам																																								
Этап подведения итогов занятия (рефлексия)	8 мин.	<div><div><div>Вопросы для обсуждения</div><div><ul style="list-style-type: none">Чему я научился?С какими трудностями я столкнулся?Какие вопросы остались? Что осталось непонятным?</div></div></div>	Педагог способствует размышлению обучающихся над вопросами																																								

Информация о домашнем задании, инструктаж по его применению	5 мин.	-	
---	--------	---	--

Рекомендуемые ресурсы для дополнительного изучения:

1. Введение в машинное обучение с помощью Scikit-learn. [Электронный ресурс] – Режим доступа: <https://habr.com/ru/post/264241/>
2. Библиотека Scikit-learn в Python. [Электронный ресурс] – Режим доступа: <https://pythonim.ru/libraries/biblioteka-scikit-learn-v-python>
3. Решающие деревья. [Электронный ресурс] – Режим доступа: https://ml-handbook.ru/chapters/decision_tree/intro
4. Метод случайного леса. [Электронный ресурс] – Режим доступа: <https://habr.com/ru/company/ods/blog/324402/>