

Тема: Data Science.

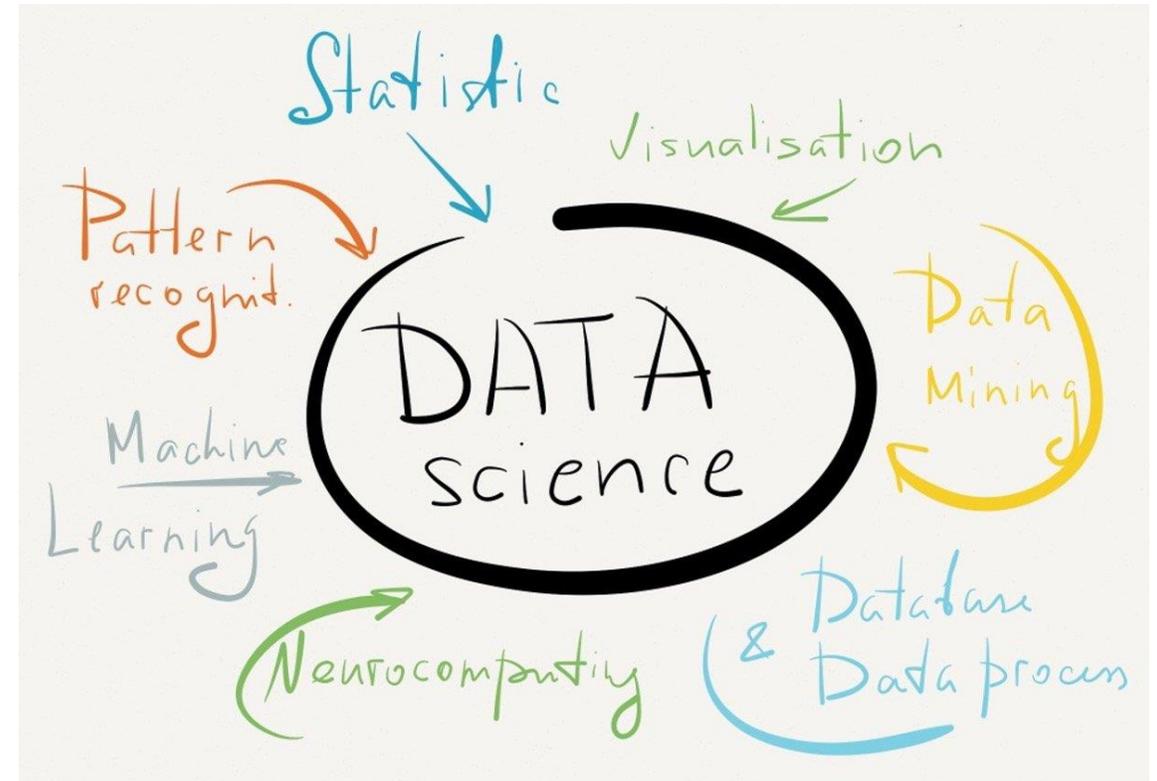
Исследовательский анализ данных о том, как ученики онлайн-школы «Сириус.Курсы» учились на нескольких курсах.

Что такое Data Science?

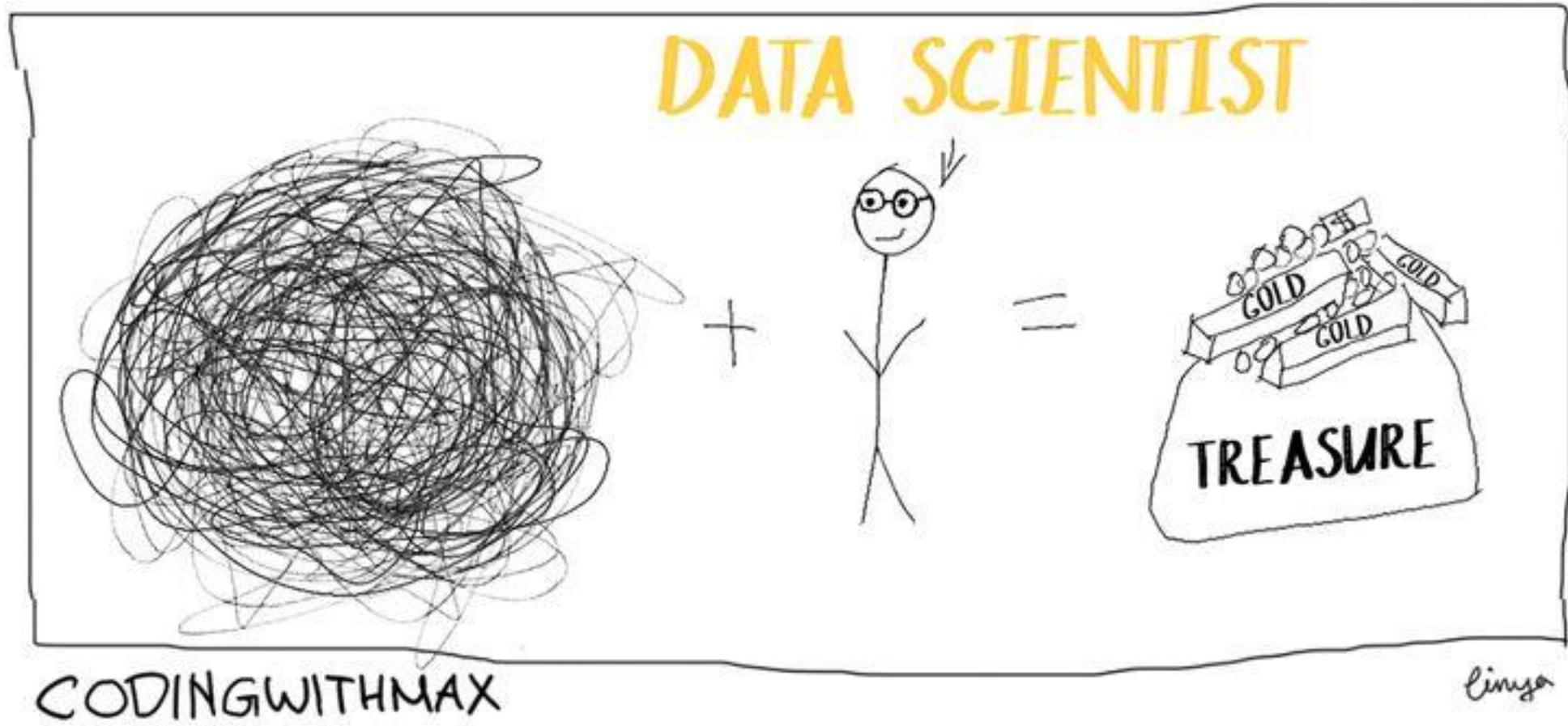
Наука о данных (англ. data science) — раздел информатики, изучающий проблемы анализа, обработки и представления данных в цифровой форме.

В рамках науки о данных рассматриваются:

- методы обработки данных в условиях больших объёмов и высокого уровня параллелизма;
- статистические методы;
- методы интеллектуального анализа данных и приложения искусственного интеллекта для работы с данными;
- методы проектирования и разработки баз данных.



Чем занимается Data Scientist?



Что такое Pandas?

- — это библиотека Python для обработки и анализа структурированных данных, её название происходит от «panel data» («панельные данные»). Панельными данными называют информацию, полученную в результате исследований и структурированную в виде таблиц. Для работы с такими массивами данных и создан Pandas.



Для чего используется Pandas?

Pandas может использоваться во всех процессах анализа данных. С помощью этой библиотеки можно:

- Импортировать наборы данных из баз данных, электронных таблиц, CSV-файлов и т.д.
- Очищать наборы данных, например, устраняя пропущенные значения.
- Упорядочивать наборы данных путем преобразования их структуры в формат, пригодный для анализа.
- Агрегировать данные, вычисляя сводную статистику, например, среднее значение столбцов, корреляцию между ними и т.д.
- Визуализировать наборы данных и открывать новые возможности.
- Pandas также имеет функционал для анализа временных рядов и текстовых данных.

Создание инструкции по Pandas

- Были изучены и подробно описаны в отдельном текстовом документе многие функции библиотеки Pandas

Библиотека `Pandas` в `Python` – это мощный инструмент для анализа данных и их обработки. Вот подробная инструкция по основным функциям и возможностям `Pandas` с примерами.

1) Установка `Pandas`

Для начала убедитесь, что у вас установлен `pip`. Если нет, установите его через `pip`:

```
python  
pip install pandas
```

2) Импорт `Pandas`

Импортируйте `Pandas` в ваш скрипт:

```
python  
import pandas as pd
```

3) Основные структуры данных: `Series` и `DataFrame`

`Series` – это одномерный массив, способный хранить данные любого типа (целые числа, строки, числа с плавающей точкой, объекты `Python` и т.д.).

Пример создания `Series`:

```
python  
data = pd.Series([1, 3, 5, 7, 9])
```

`DataFrame`

`DataFrame` – это двумерная структура данных, похожая на таблицу с рядами и столбцами.

Пример создания `DataFrame`:

```
python  
data = {  
    'Country': ['Russia', 'Colombia', 'Chile', 'Ecuador', 'Nigeria'],  
    'Capital': ['Moscow', 'Bogota', 'Santiago', 'Quito', 'Abuja'],  
    'Population': [144, 50, 18, 17, 206]  
}  
df = pd.DataFrame(data)
```

4) Выборка данных

Выборка данных в `Pandas` может быть осуществлена разными способами.

Выборка по столбцам

```
python  
df['Country'] # Возвращает столбец 'Country'
```

Выборка по строкам

```
python  
df.iloc[0] # Возвращает первую строку  
df.iloc[0:2] # Возвращает строки с индексами от 0 до 2
```

5) Фильтрация данных

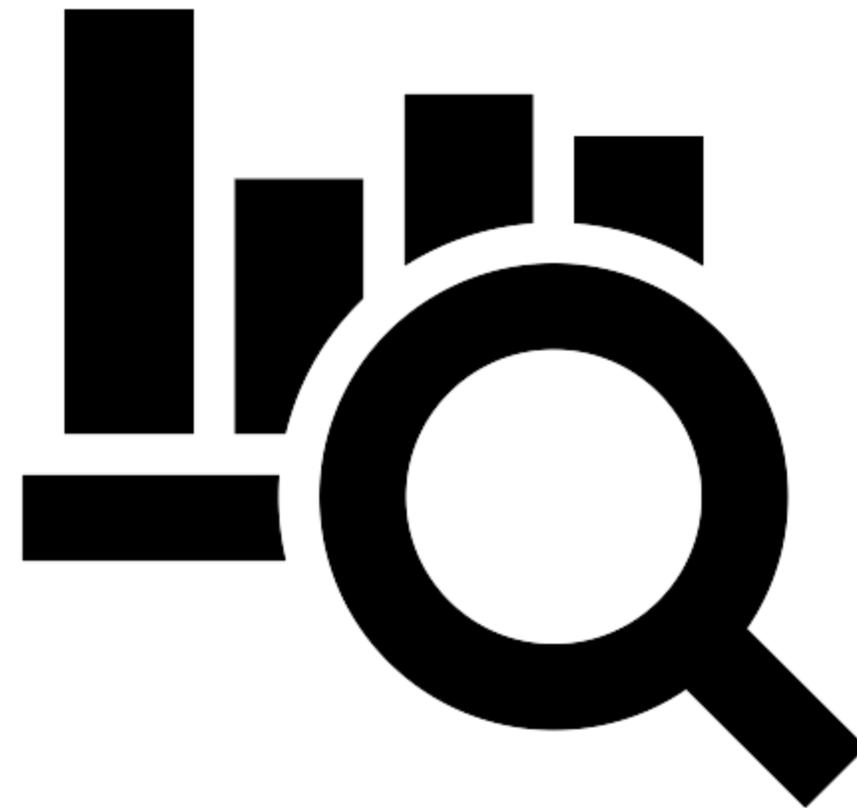
```
python
```

**Основная цель - найти
признаки по которым можно
определить будет ли продолжать
обучение ученик**

Выбор метрик.

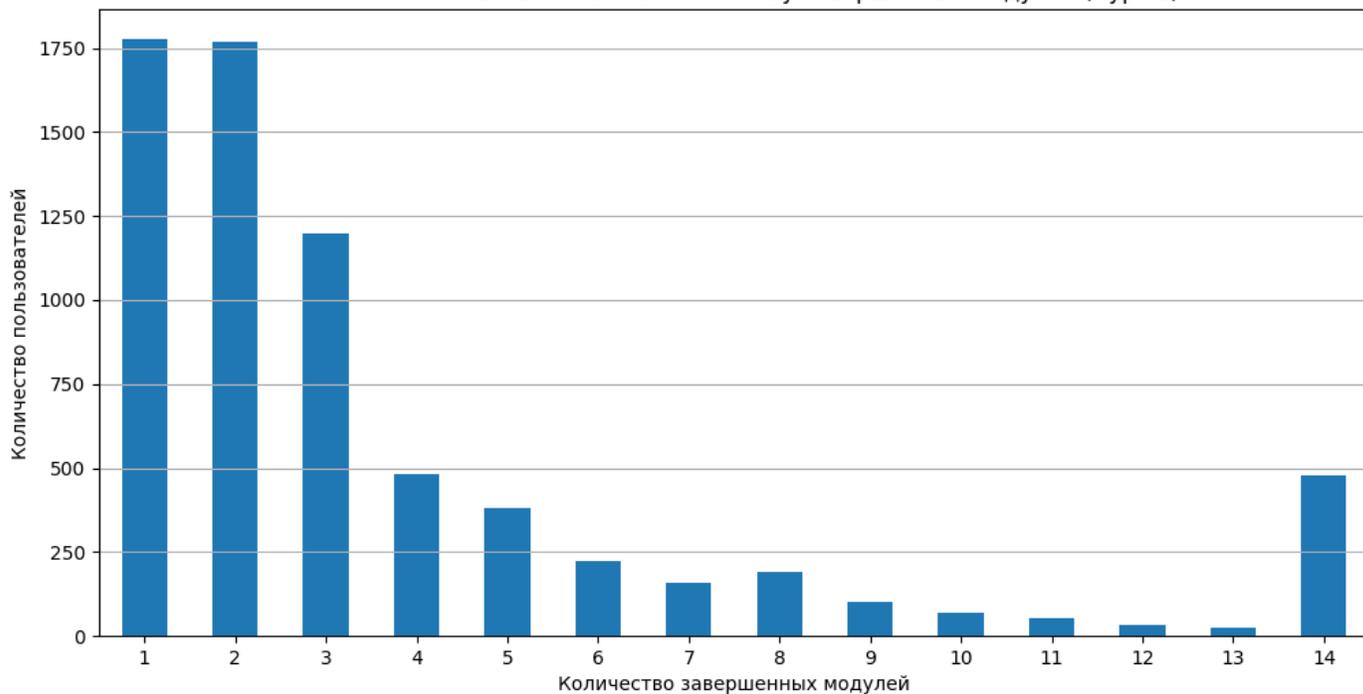
- 1) Число завершённых модулей
- 2) Среднее количество попыток
- 3) Соотношение успешных и неудачных попыток
- 4) Регулярность занятий
- 5) Временные интервалы между входами в систему
- 6) Среднее время на выполнение модуля

**Анализ
полученных
данных.**



Число завершённых модулей

Количество пользователей по количеству завершённых модулей (Курс 1)



```
1 #первый про модули
2 import pandas as pd
3 import matplotlib.pyplot as plt
4
5 data = pd.read_csv('user_module_progress.csv')
6
7 course_1_data = data[data['course_id'] == 1]
8
9 completed_modules_course_1 = course_1_data[course_1_data['is_achieved'] == True]
10
11 user_completed_modules = completed_modules_course_1.groupby('user_id').size()
12
13 user_count_by_completed_modules = user_completed_modules.value_counts().sort_index()
14
15 # Построение графика
16 plt.figure(figsize=(12, 6))
17 user_count_by_completed_modules.plot(kind='bar')
18 plt.title('Количество пользователей по количеству завершённых модулей (Курс 1)')
19 plt.xlabel('Количество завершённых модулей')
20 plt.ylabel('Количество пользователей')
21 plt.xticks(rotation=0)
22 plt.grid(axis='y')
23 plt.show()
24
```

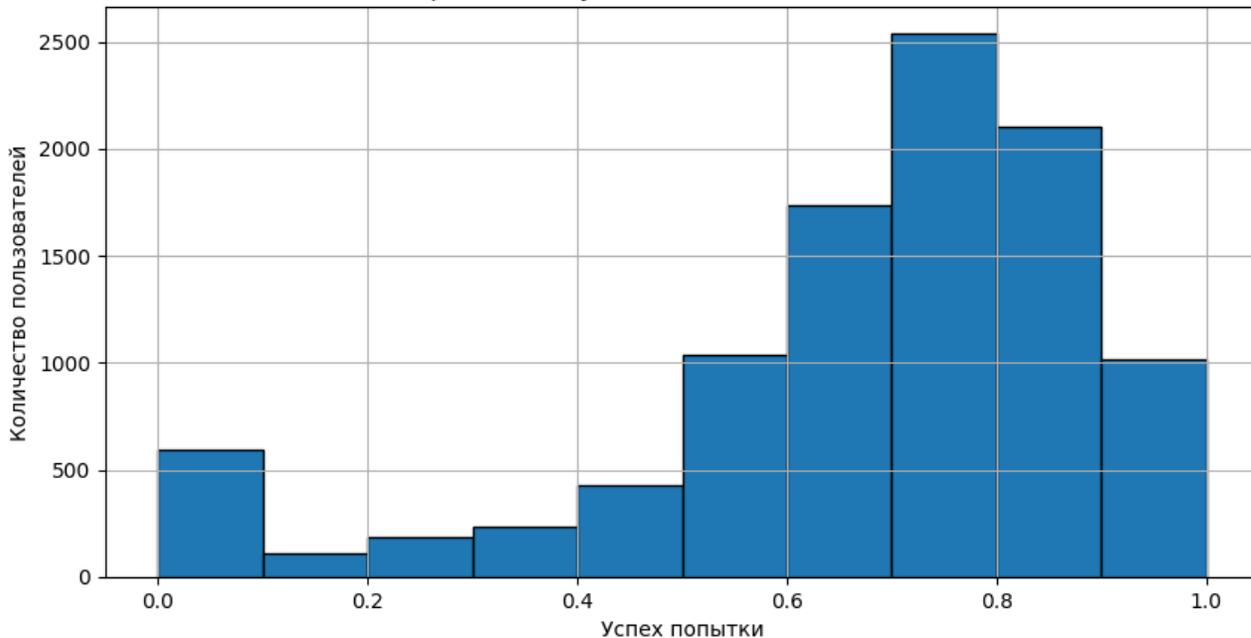
Среднее количество попыток



```
1 #про кол-во попыток
2 import pandas as pd
3 import matplotlib.pyplot as plt
4
5 user_element_progress = pd.read_csv('user_element_progress.csv')
6
7 course_1_progress = user_element_progress[user_element_progress['course_id'] == 1]
8
9 modules_completed_course_1 = course_1_progress[course_1_progress['is_achieved'] == True] \
10     .groupby('user_id')['course_module_id'] \
11     .nunique()
12
13 avg_tries_per_user_course_1 = course_1_progress.groupby('user_id')['tries_count'].mean()
14
15 user_stats_course_1 = pd.DataFrame({'modules_completed': modules_completed_course_1, 'avg_tries': avg_tries_per_user_course_1})
16
17 avg_tries_by_module_count_course_1 = user_stats_course_1.groupby('modules_completed')['avg_tries'].mean()
18
19 #график
20 plt.figure(figsize=(12, 6))
21 avg_tries_by_module_count_course_1.plot(kind='bar')
22 plt.title('Среднее количество попыток')
23 plt.xlabel('Количество выполненных модулей у ученика')
24 plt.ylabel('Количество попыток')
25 plt.xticks(rotation=0)
26 plt.grid(axis='y')
27 plt.show()
```

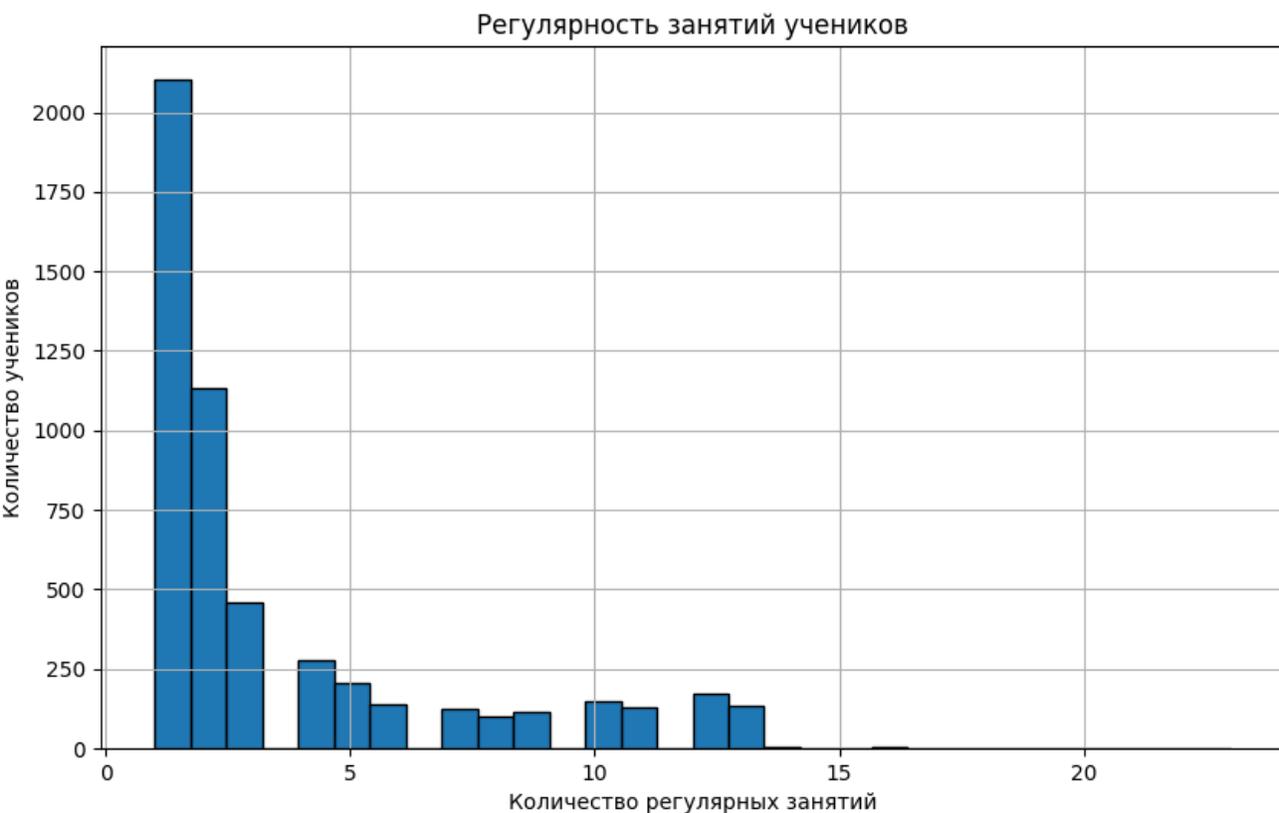
Соотношение успешных и неудачных попыток

Распределение успеха попыток пользователей



```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 element_solution = pd.read_csv('user_element_solution.csv')
5 element_progress = pd.read_csv('user_element_progress.csv')
6
7 # Оценка 'verdict' как 1 для 'ok' и 0 для других значений
8 element_solution['verdict'] = element_solution['verdict'].apply(lambda x: 1 if x == 'ok' else 0)
9
10 solution_stats = element_solution.groupby('element_progress_id')['verdict'].mean().reset_index()
11
12 progress_with_solutions = pd.merge(element_progress, solution_stats, how='left', left_on='id', right_on='element_progress_id')
13
14 user_success_rate = progress_with_solutions.groupby('user_id')['verdict'].mean()
15
16 plt.figure(figsize=(10, 5))
17 user_success_rate.hist(bins=10, edgecolor='black')
18 plt.title('Распределение успеха попыток пользователей')
19 plt.xlabel('Успех попытки')
20 plt.ylabel('Количество пользователей')
21 plt.show()
22
```

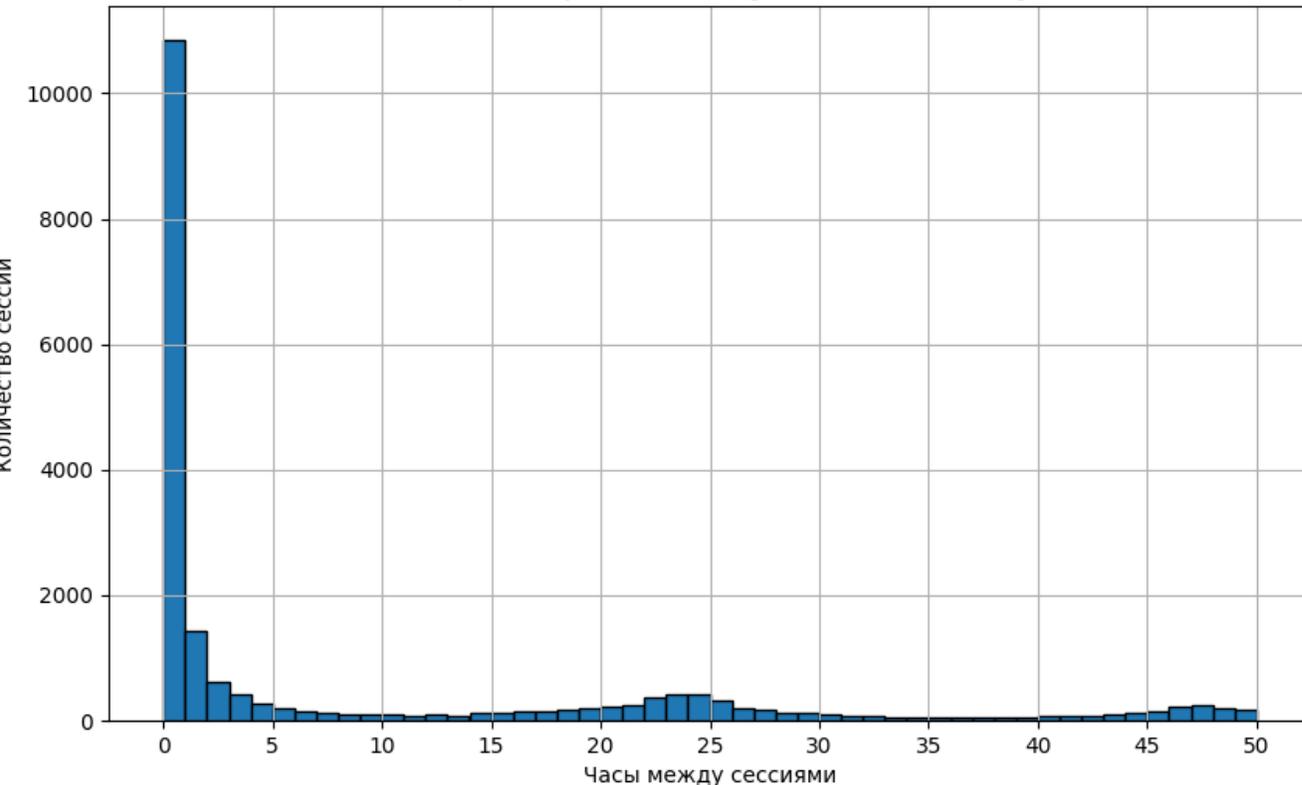
Регулярность занятий



```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 user_module_progress = pd.read_csv('user_module_progress.csv')
5
6 user_module_progress['time_achieved'] = pd.to_datetime(user_module_progress['time_achieved'], errors='coerce')
7
8 completed_modules = user_module_progress[user_module_progress['is_achieved'] == True].copy()
9
10 completed_modules.sort_values(by=['user_id', 'time_achieved'], inplace=True)
11 completed_modules['time_diff'] = completed_modules.groupby('user_id')['time_achieved'].diff()
12
13 completed_modules['days_between_achievements'] = completed_modules['time_diff'].dt.days
14
15 regular_achievements = completed_modules[completed_modules['days_between_achievements'] <= 7]
16
17 regularity_counts = regular_achievements.groupby('user_id').size()
18
19 plt.figure(figsize=(10, 6))
20 regularity_counts.hist(bins=30, edgecolor='black')
21 plt.title('Регулярность занятий учеников')
22 plt.xlabel('Количество регулярных занятий')
23 plt.ylabel('Количество учеников')
24 plt.show()
```

Временные интервалы между входами в систему

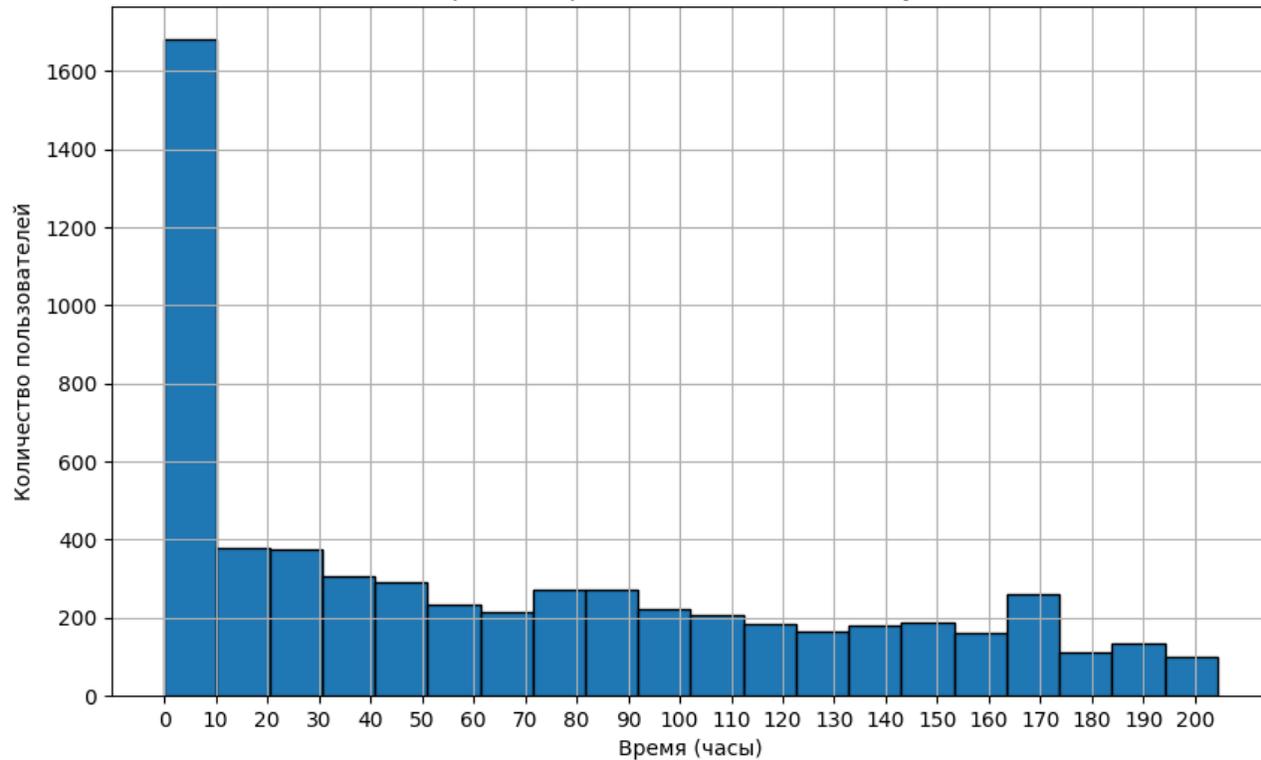
Интервалы времени между входами в систему



```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 user_module_progress = pd.read_csv('user_module_progress.csv')
5
6 user_module_progress['time_unlocked'] = pd.to_datetime(user_module_progress['time_unlocked'], errors='coerce')
7
8 sorted_progress = user_module_progress.sort_values(by=['user_id', 'time_unlocked'])
9 sorted_progress['time_diff'] = sorted_progress.groupby('user_id')['time_unlocked'].diff()
10 sorted_progress['hours_between_sessions'] = sorted_progress['time_diff'].dt.total_seconds() / 3600
11
12 interval_data = sorted_progress[sorted_progress['hours_between_sessions'].notna()]
13 plt.figure(figsize=(10, 6))
14 interval_data['hours_between_sessions'].hist(bins=50, range=(0, 50), edgecolor='black')
15 plt.title('Интервалы времени между входами в систему')
16 plt.xlabel('Часы между сессиями')
17 plt.ylabel('Количество сессий')
18 plt.xticks(range(0, 51, 5)) # Установка меток с шагом в 5 часов
19 plt.grid(True)
20 plt.show()
```

Среднее время на выполнение модуля

Среднее время на достижение модуля



```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 # Загрузка данных
5 user_module_progress = pd.read_csv('user_module_progress.csv')
6
7 # Преобразование дат
8 user_module_progress['time_unlocked'] = pd.to_datetime(user_module_progress['time_unlocked'], errors='coerce')
9 user_module_progress['time_achieved'] = pd.to_datetime(user_module_progress['time_achieved'], errors='coerce')
10
11 # Расчет времени на достижение модуля
12 user_module_progress['time_to_achieve'] = (user_module_progress['time_achieved'] - user_module_progress['time_unlocked']).dt.total_seconds() / 3600 # в часах
13 completed_modules_time = user_module_progress[user_module_progress['is_achieved'] == True]
14 average_time_per_module = completed_modules_time.groupby('user_id')['time_to_achieve'].mean()
15
16 # Визуализация с уточненными данными
17 plt.figure(figsize=(10, 6))
18 upper_limit = average_time_per_module.quantile(0.75)
19 num_bins = 20
20 average_time_per_module.hist(bins=num_bins, range=(0, upper_limit), edgecolor='black')
21 plt.title('Среднее время на достижение модуля')
22 plt.xlabel('Время (часы)')
23 plt.ylabel('Количество пользователей')
24 plt.xticks(range(0, int(upper_limit) + 1, int(upper_limit / num_bins)))
25 plt.grid(True)
26 plt.show()
```