

Natural Language Processing Course

Autumn 2023 Stream 5

15 October 2023

Цель Проекта

mllibs проектная работа

Создания инструмента который позволит пользователю
(**без знания программирования**) выполнять проекты **машинного обучения** используя **текстовые запросы**

Цель Проекта

mllibs проектная работа

Создания инструмента который позволит пользователю
(без знания программирования) выполнять проекты **машинного обучения** используя **текстовые запросы**

Интерпретатор
(и исполнитель)

1

Минимальное или без
знания написания кода

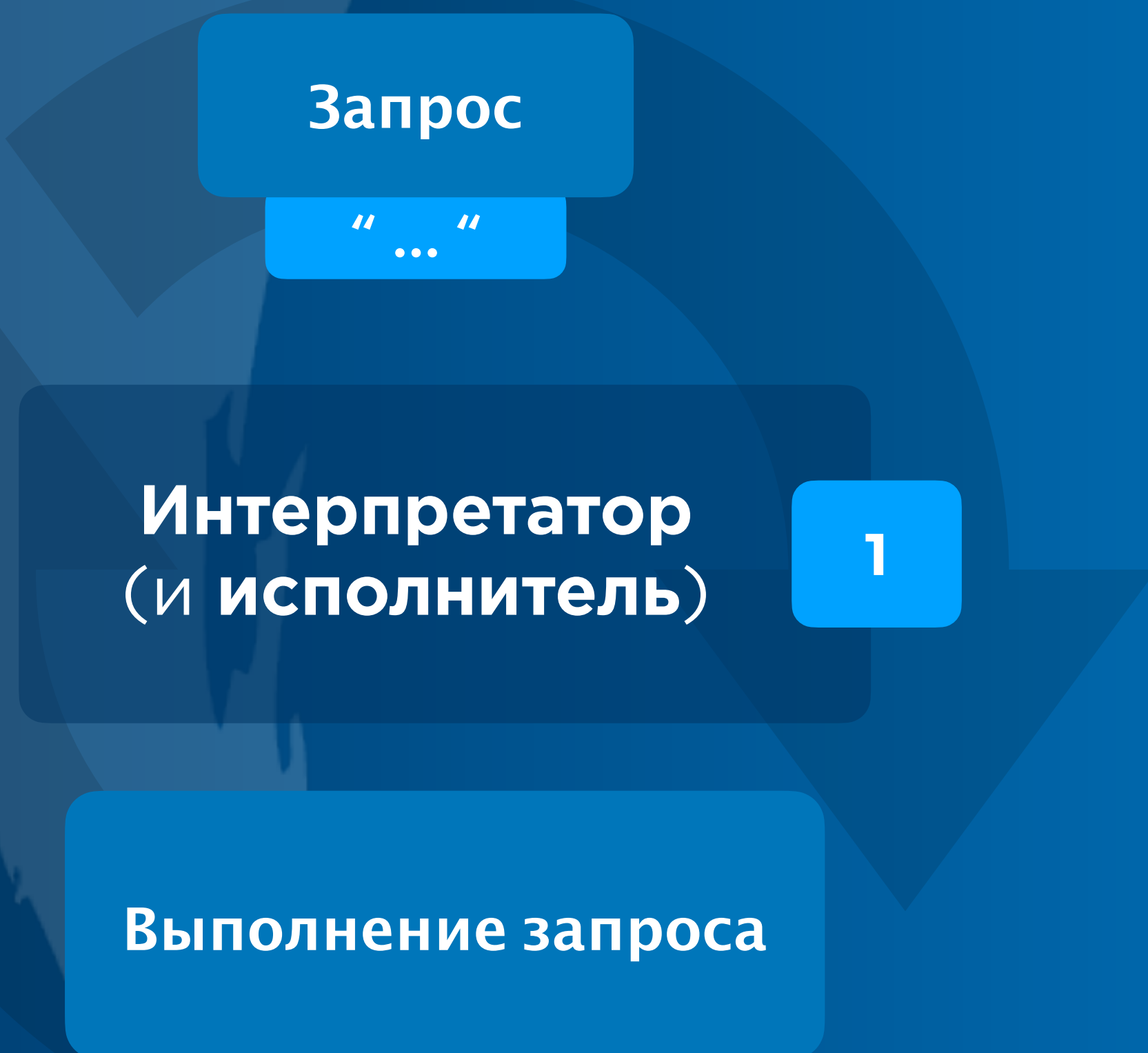
Генераторы

2

Есть знание написания кода

Цель Проекта

mllibs проектная работа



```
1 import pandas as pd
2
3 def read_csv(path:str):
4     pd.read_csv(path)
```

Запрос

" ... "

which pandas function can I use to drop missing data

Генераторы 2

You can use the `dropna()` function in pandas to drop missing data.

```
1 import pandas as pd
2
3 def drop_na(data):
4     data.dropna()
```

```
1 import seaborn as sns
2
3 def plot_scatter(data,x,y):
4     sns.scatterplot(data,x,y)
```

Этапы ML проекта

mllibs проектная работа

Выбор модели: на этом этапе выбирается модель машинного обучения, которая будет использоваться для решения задачи

1

Сбор и подготовка данных: в этом этапе происходит сбор данных, их анализ и очистка от ошибок и выбросов

2

Обучение модели: в этом этапе происходит обучение модели на тренировочных данных

3

Оценка модели: после обучения модели необходимо оценить ее качество на тестовых данных

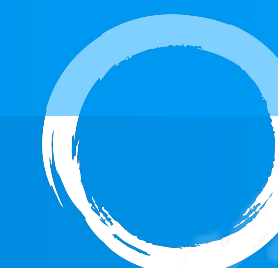
4

Тюнинг гиперпараметров: на этом этапе происходит подбор оптимальных значений гиперпараметров модели

5

Развертывание модели: после успешного обучения и оценки модели ее можно развернуть в продакшн

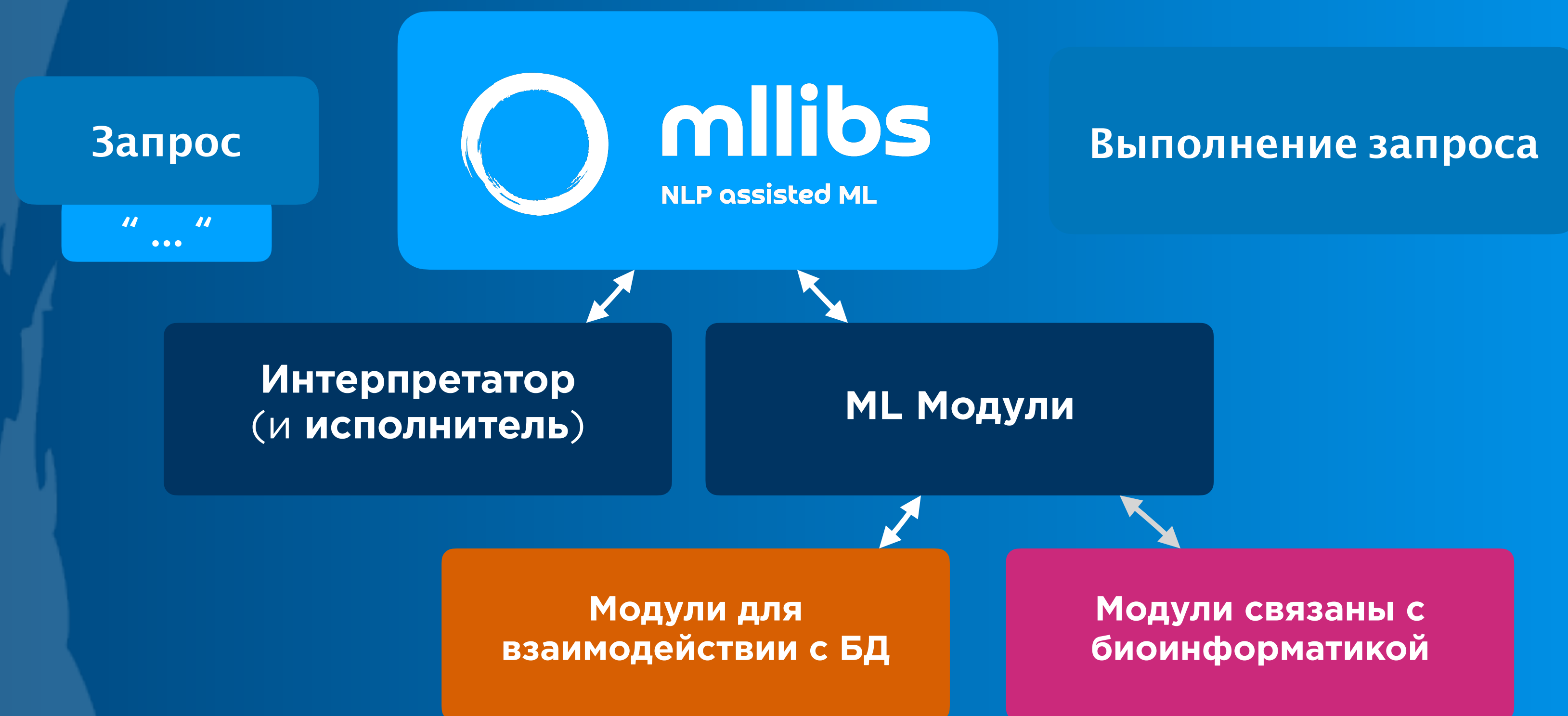
6



mllibs
NLP assisted ML

Зачем нужен такой инструмент

mlibs проектная работа



Обзор Интерпретатора

Методы Активизации Кода

mllibs проектная работа

Функция

```
1 def fib_list(n):
2     result = []
3     a,b = 0,1
4     while a<n:
5         result.append(a)
6         a,b = b, a + b
7     return result
8
9 fib_list(5)
```

Класс

```
1 class fib_list:
2
3     def __init__(self,n):
4         self.n = n
5
6     def get_list(self):
7         result = []
8         a,b = 0,1
9         while a<self.n:
10             result.append(a)
11             a,b = b, a + b
12         return result
13
14 fib = fib_list(5)
15 fib.get_list()
```

Текст

```
1 input = 'calculate the fibonacci
2         sequence for the value
3         of 5'
4 nlp_interpreter(input)
```


Как это можно реализовать

mllibs проектная работа



Как это можно реализовать

mllibs проектная работа



Функции активации

mlibs проектная работа

Сбор и подготовка данных: в этом этапе происходит сбор данных, их анализ и очистка от ошибок и выбросов

Загрузка данных

EDA

Пропуски

Выбросы

```
1 import pandas as pd
2
3 def read_csv(path:str):
4     pd.read_csv(path)
```

```
1 import seaborn as sns
2
3 def plot_scatter(data,x,y):
4     sns.scatterplot(data,x,y)
```

```
1 import pandas as pd
2
3 def drop_na(data):
4     data.dropna()
```

```
1 import pandas as pd
2 import numpy as np
3
4 def drop_na(data):
5
6     mean = np.mean(data)
7     std = np.std(data)
8
9     threshold = 3
10    outlier = []
11    for i in data:
12        z = (i-mean)/std
13        if z > threshold:
14            outlier.append(i)
15    print('outlier in dataset
is', outlier)
```

Функции активации и метки

mlibs проектная работа

Сбор и подготовка данных: в этом этапе происходит сбор данных, их анализ и очистка от ошибок и выбросов

Загрузка данных

EDA

Пропуски

Выбросы

```
1 import pandas as pd
2
3 def read_csv(path:str):
4     pd.read_csv(path)
```

load_csv

```
1 import seaborn as sns
2
3 def plot_scatter(data,x,y):
4     sns.scatterplot(data,x,y)
```

eda_scatter

```
1 import pandas as pd
2
3 def drop_na(data):
4     data.dropna()
```

pp_drop

```
1 import pandas as pd
2 import numpy as np
3
4 def drop_na(data):
5
6     mean = np.mean(data)
7     std = np.std(data)
8
9     threshold = 3
10    outlier = []
11    for i in data:
12        z = (i-mean)/std
13        if z > threshold:
14            outlier.append(i)
15    print('outlier in dataset
is', outlier)
```

pp_outlier

mlibs

NLP assisted ML

NLP подходы для вызова функции

mllibs проектная работа

NLP задачи которые могут вызвать функцию

Классификация

User пишет запрос в текстовом формате и мы предсказываем метку выполняя классификацию текста

Генерации Текста

User пишет запрос в текстовом формате и мы используем это как контекст, прогоняем через (Transformer Decoder/Seq-Seq) модель и получаем генерированный текст содержащий метку (или метки)

NLP подходы для вызова функции

mlibs проектная работа

NLP задачи которые могут вызвать функцию

Обучение

Классификация

User пишет запрос в текстовом формате и мы предсказываем метку выполняя классификацию текста

“Загрузка CSV”

“Загрузите CSV”

“Загрузите CSV используя Pandas”

load_csv

“Загрузка JSON”

“Загрузите JSON”

load_json

Предобработка текста

Токенизация + BOW

Токенизация + Word2Vec

...

RandomForestClassifier()



mlibs

NLP assisted ML

NLP подходы для вызова функции

mllibs проектная работа

NLP задачи которые могут вызвать функцию

Inference

Классификация

User пишет запрос в текстовом формате и мы предсказываем метку выполняя классификацию текста

“Загрузите мне CSV”

Предобработка текста

Токенизация + BOW
Токенизация + Word2Vec
...

RandomForestClassifier()

load_csv



mllibs
NLP assisted ML

NLP подходы для вызова функции

mllibs проектная работа

NLP задачи которые могут вызвать функцию

Генерации Текста

User пишет запрос в текстовом формате и мы используем это как контекст, прогоняем через (Transformer Decoder/Seq-Seq) модель и получаем сгенерированный текст содержащий метку (или метки)

```
1 from transformers import GPT2LMHeadModel, GPT2Tokenizer
2
3 model_name_or_path = "sberbank-ai/rugpt3large_based_on_gpt2"
4 tokenizer = GPT2Tokenizer.from_pretrained(model_name_or_path)
5 model = GPT2LMHeadModel.from_pretrained(model_name_or_path).cuda()
6
7 #вводный текст
8 text = "Александр Сергеевич Пушкин родился в "
9
10 input_ids = tokenizer.encode(text, return_tensors="pt").cuda()
11 out = model.generate(input_ids.cuda())
12 generated_text = list(map(tokenizer.decode, out))[0]
13
14 # сгенерированный текст
15 print(generated_text)
16 # Александр Сергеевич Пушкин родился в \n1799 году. Его отец был крепостным
    крестьянином, а мать – крепостной крестьянкой. Детство и юность Пушкина прошли
    в деревне Михайловское под Петербургом. В 1820-х годах семья переехала
```

NLP подходы для вызова функции

mllibs проектная работа

NLP задачи которые могут вызвать функцию

User пишет запрос в текстовом формате и мы используем это как контекст, прогоняем через (Transformer Decoder/Seq-Seq) модель и получаем генерированный текст содержащий метку (или метки)

Генерации Текста

“Загрузка CSV”

@@первый@@

“**load_csv**”

@@второй@@

Fine-tune

ruDialogPT



mllibs
NLP assisted ML

Есть связь между ФА и запросом mlibs проектная работа

Есть NLP инструмент который может привязать
вводный запрос к функции активации

Сбор библиотек/модулей

```
1 from sklearn.model_selection import
  train_test_split
2 import pandas as pd
3
4 # load the dataset
5 df = pd.read_csv('data.csv')
6
7 # split the dataset into training and
8 testing sets
9 X_train, X_test, y_train, y_test =
  train_test_split(df.drop('target', axis=1),
10 df['target'], test_size=0.2, random_state=42)
11 print(X_train.shape, y_train.shape)
12 print(X_test.shape, y_test.shape)
```

```
1 import seaborn as sns
2
3 def plot_scatter(data, x, y):
4     sns.scatterplot(data, x, y)
```

```
1 import pandas as pd
2
3 def read_csv(path: str):
4     pd.read_csv(path)
```

```
1 import pandas as pd
2
3 def drop_na(data):
4     data.dropna()
```

```
1 import pandas as pd
2 import numpy as np
3
4 def drop_na(data):
5
6     mean = np.mean(data)
7     std = np.std(data)
8
9     threshold = 3
10    outlier = []
11    for i in data:
12        z = (i - mean) / std
13        if z > threshold:
14            outlier.append(i)
15    print('outlier in dataset
  is', outlier)
```

```
1 import pandas as pd
2
3 def describe_data(data):
4     data.describe()
```

```
1 from sklearn.linear import
  LogisticRegression
2
3 def describe_data(X, y):
4     model =
  LogisticRegression()
5     model.fit(X, y)
```


Пример из 0.1.7

mllibs проектная работа

Сбор модулей для EDA

meda_scplot

Модуль для визуализации колонок в **Seaborn**)

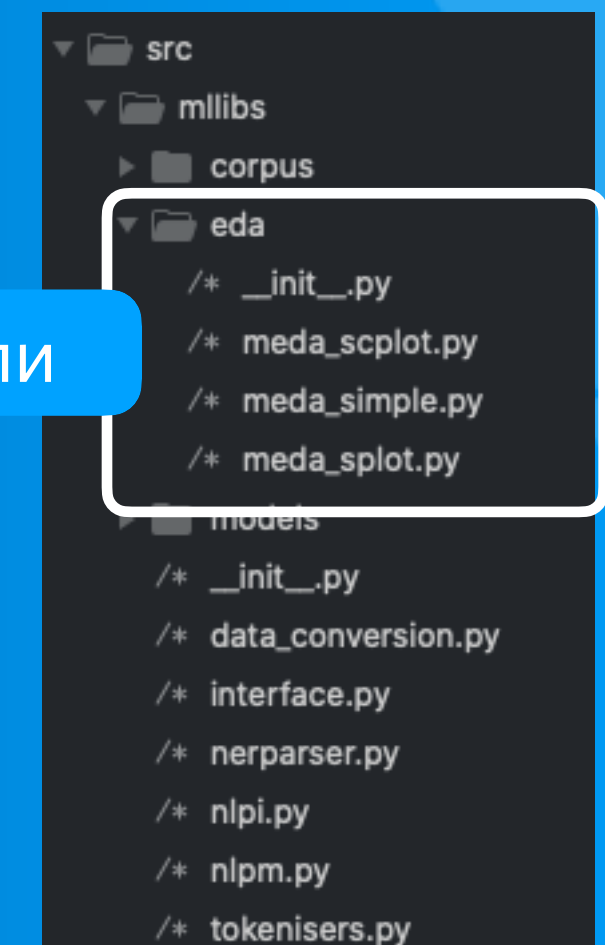
meda_simple

Модуль для простых разведывательных операции в **Pandas**)

meda_splot

Модуль для стандартных функции визуализации в **Seaborn**)

EDA модули



```
1 import seaborn as sns
2 sns.set_theme(style="ticks")
3
4 dots = sns.load_dataset("dots")
5
6 # Define the palette as a list to specify exact values
7 palette = sns.color_palette("rocket_r")
8
9 # Plot the lines on two facets
10 sns.relplot(
11     data=dots,
12     x="time", y="firing_rate",
13     hue="coherence", size="choice", col="align",
14     kind="line", size_order=["T1", "T2"], palette=palette,
15     height=5, aspect=.75, facet_kws=dict(sharex=False),
16 )
```

Нужны вводные данные (eg. **Dots**)

Как минимум нужно еще **x,y**

В **meda_splot** есть соответствующая функция
содержащее **relplot**

Пример из 0.1.7

mlibs проектная работа

Запрос

A scatter plot showing the relationship between bill length (mm) on the x-axis and bill depth (mm) on the y-axis. The x-axis ranges from approximately 32 to 60 mm, with major ticks at 35, 40, 45, 50, 55, and 60. The y-axis ranges from approximately 13 to 21 mm, with major ticks at 14, 16, 18, and 20. Data points are categorized by island: Torgersen (light green), Biscoe (teal), and Dream (dark blue). The plot shows a general positive correlation between bill length and bill depth. Torgersen points are clustered at lower bill lengths (35-45 mm) and depths (16-21 mm). Biscoe points are clustered at intermediate bill lengths (45-55 mm) and depths (14-18 mm). Dream points are clustered at higher bill lengths (50-60 mm) and depths (18-21 mm).

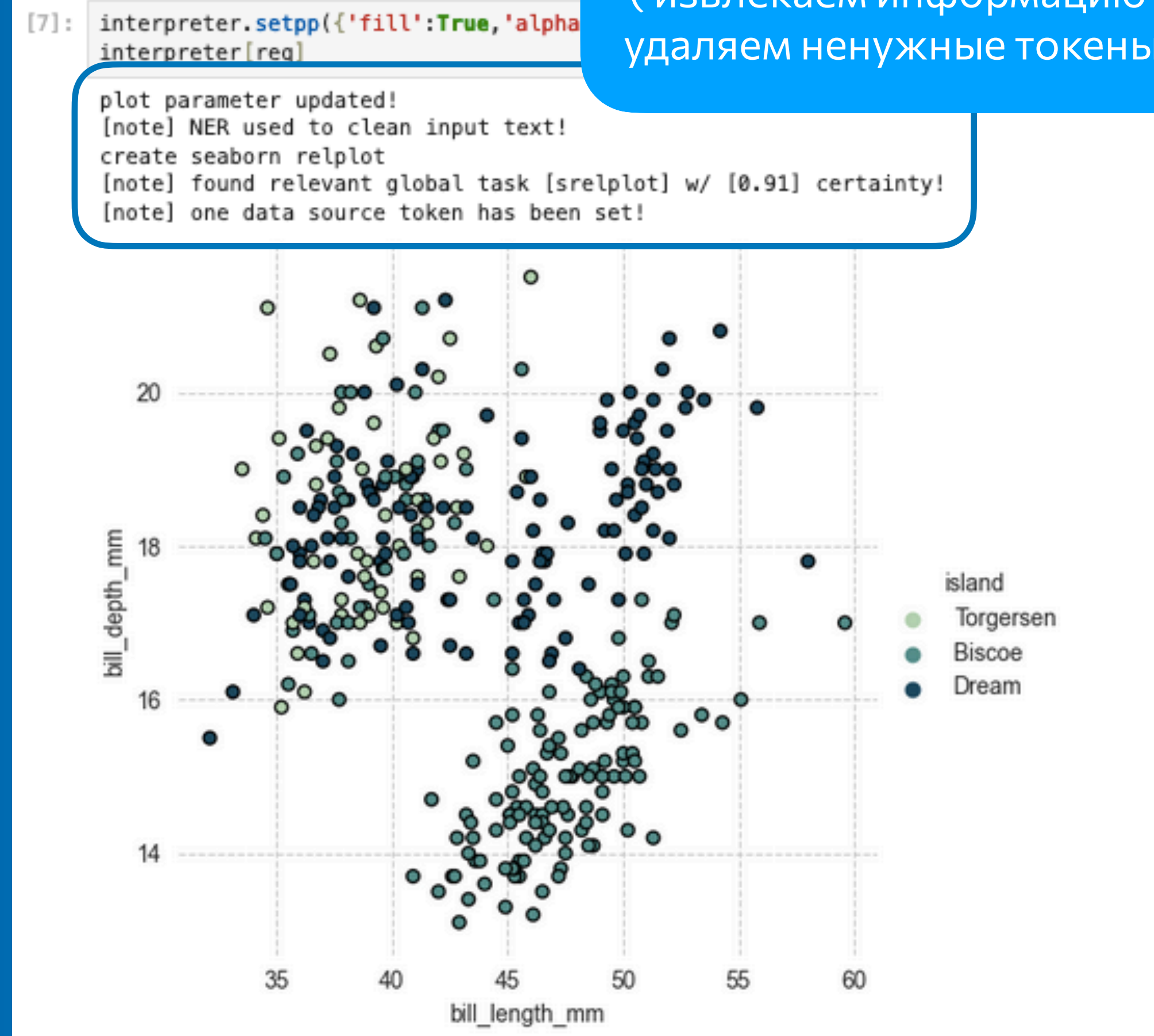


Пример из 0.1.7

mlibs проектная работа

```
1 '''  
2  
3 Seaborn Plots  
4  
5 '''  
6  
7 req = ''  
8 create seaborn relplot  
9 x: bill_length_mm  
10 y: bill_depth_mm  
11 hue: island  
12 using penguins  
13 '''  
14  
15 # interpreter['create seaborn boxplot using housing y AGE x RAD']  
16 # interpreter['create seaborn relplot x: bill_length_mm y: bill_depth_mm  
    hue island col=island alpha=1.0 s:50 mew: 1 using penguins']  
17 # interpreter['create seaborn relplot x: bill_length_mm y: bill_depth_mm  
    hue island using penguins']  
18  
19 interpreter.setpp({'fill':True,'alpha':1,'mew':1})  
20 interpreter[req]
```

Фильтрация запроса
(извлекаем информацию и
удаляем ненужные токены)



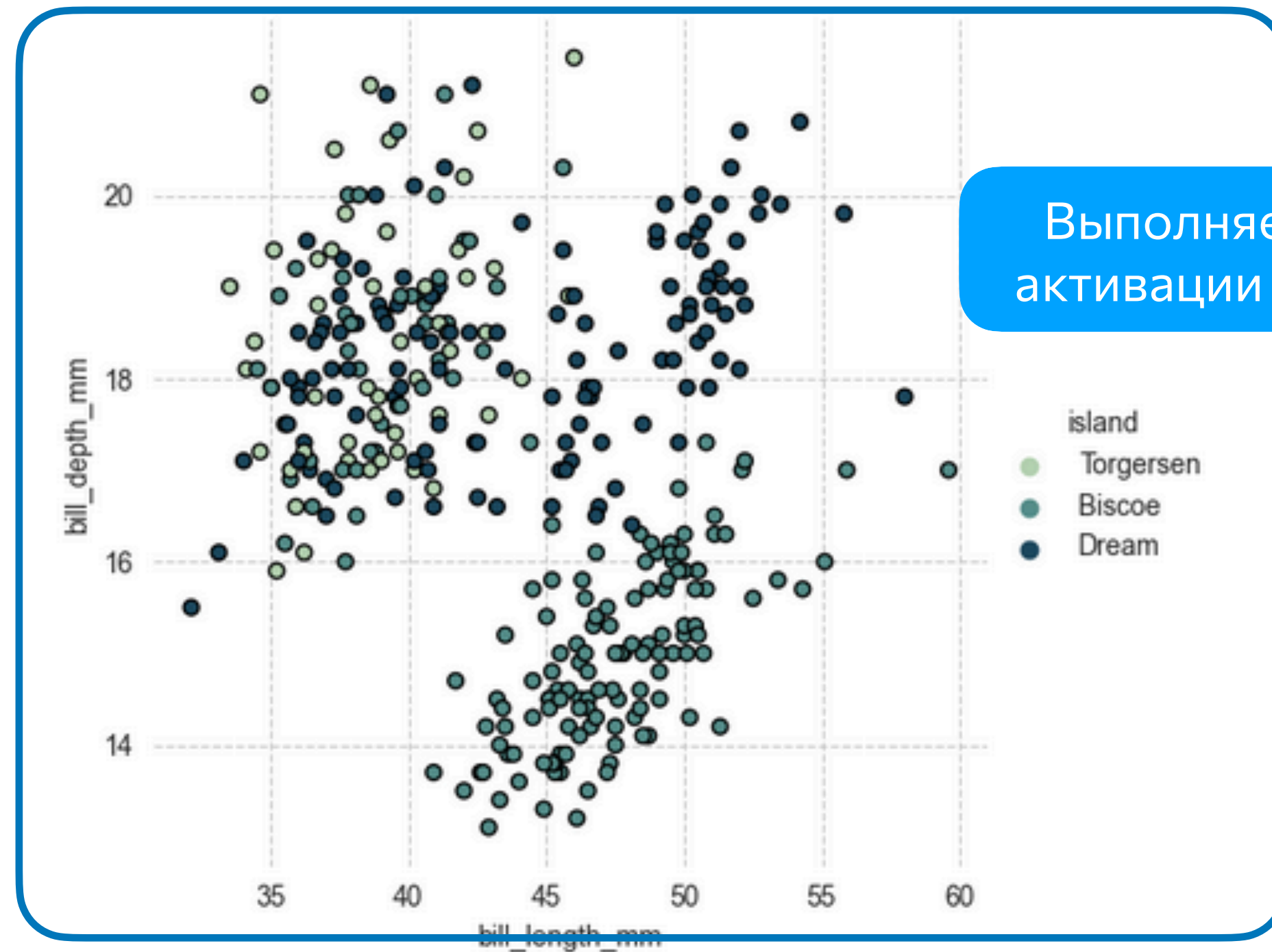
Пример из 0.1.7

mlibs проектная работа

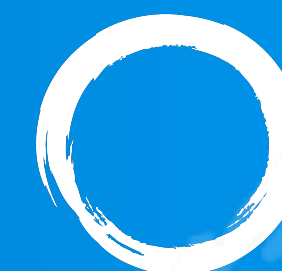
```
1 '''  
2  
3 Seaborn Plots  
4  
5 '''  
6  
7 req = ''  
8 create seaborn relplot  
9 x: bill_length_mm  
10 y: bill_depth_mm  
11 hue: island  
12 using penguins  
13 '''  
14  
15 # interpreter['create seaborn boxplot using housing y AGE x RAD']  
16 # interpreter['create seaborn relplot x: bill_length_mm y: bill_depth_mm  
    hue island col=island alpha=1.0 s:50 mew: 1 using penguins']  
17 # interpreter['create seaborn relplot x: bill_length_mm y: bill_depth_mm  
    hue island using penguins']  
18  
19 interpreter.setpp({'fill':True,'alpha':1,'mew':1})  
20 interpreter[req]
```

```
[7]: interpreter.setpp({'fill':True,'alpha':1,'mew':1})
      interpreter[req]

plot parameter updated!
[note] NER used to clean input text!
create seaborn relplot
[note] found relevant global task [srelplot] w/ [0.91] certainty!
[note] one data source token has been set!
```



Выполняется функция активации метки **srelplot**



NER в mlibs

mlibs проектная работа



Упрощенный

NER в mlibs

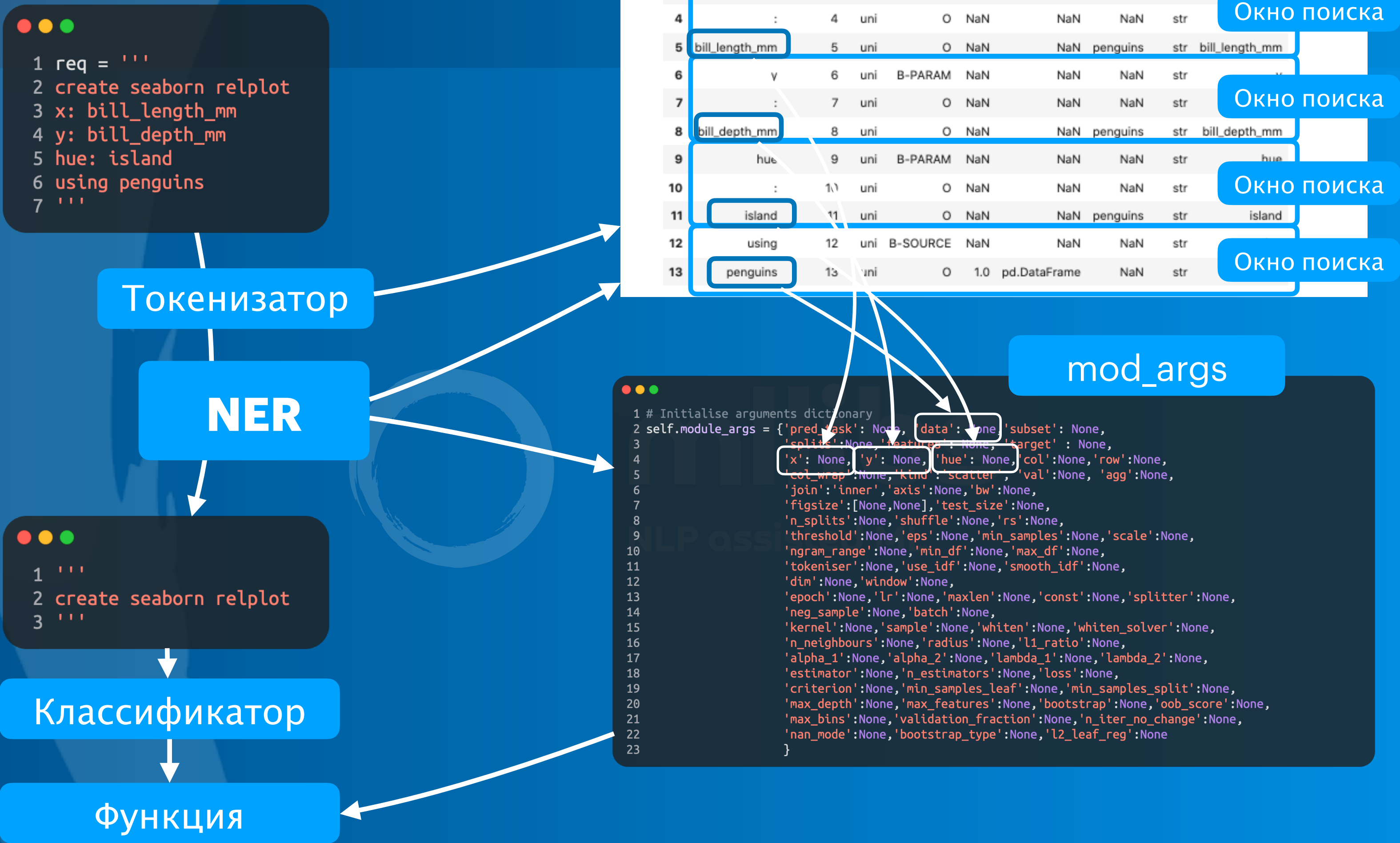
mlibs проектная работа



Упрощенный

NER в mlibs

mlibs проектная работа



Подведем итоги

mllibs проектная работа

Какие есть подзадачи для интерпретатора

- > Создания функции активации на разные темы МО
- > Классификатор для подборки функции активации
- > Генератор для подборки функции активации
- > NER tagger для извлечения информации
- > Рекомендательные системы по итогу выполнения проекта

Какие есть подзадачи для генератора

- > Ответы на вопросы по темам ML
- > Гиды по проектам ML
- > ...