

Департамент образования и науки города Москвы
Государственное автономное образовательное учреждение высшего
образования города Москвы
«Московский городской педагогический университет»
Институт цифрового образования
Департамент информатики, управления и технологий

ДИСЦИПЛИНА:

Проектный практикум по разработке ETL-решений

Лабораторная работа №4

Тема:

«Работа с XML и JSON файлами»

Выполнил(а): Морозова Валерия АДЭУ-211

Преподаватель:

Москва

2025

Вариант 10. `Parking_Violations_Issued_-_Fiscal_Year_2017.csv`.

Практическая работа 4.1. Анализ данных с помощью DASK

Задание.

4.1.1. Настроить среду и рабочий каталог.

```
2.1.1. Настроить среду и рабочий каталог.

[1] 1 from google.colab import drive
    2 drive.mount('/content/drive')

Mounted at /content/drive

[2] 1 ls

drive/ sample_data/

1 !pip install "dask[complete]"

Requirement already satisfied: dask[complete] in /usr/local/lib/python3.11/dist-packages (2024.12.1)
Requirement already satisfied: click>=8.1 in /usr/local/lib/python3.11/dist-packages (from dask[complete]) (8.1.8)
Requirement already satisfied: cloudpickle>=3.0.0 in /usr/local/lib/python3.11/dist-packages (from dask[complete]) (3.1.1)
Requirement already satisfied: fsspec>=2021.09.0 in /usr/local/lib/python3.11/dist-packages (from dask[complete]) (2024.10)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.11/dist-packages (from dask[complete]) (24.2)
Requirement already satisfied: partd>=1.4.0 in /usr/local/lib/python3.11/dist-packages (from dask[complete]) (1.4.2)
Requirement already satisfied: pyyaml>=5.3.1 in /usr/local/lib/python3.11/dist-packages (from dask[complete]) (6.0.2)
Requirement already satisfied: toolz>=0.10.0 in /usr/local/lib/python3.11/dist-packages (from dask[complete]) (0.12.1)
Requirement already satisfied: importlib_metadata>=4.13.0 in /usr/local/lib/python3.11/dist-packages (from dask[complete]) (8.1.0)
Requirement already satisfied: pyarrow>=14.0.1 in /usr/local/lib/python3.11/dist-packages (from dask[complete]) (18.1.0)
Collecting lz4>=4.3.2 (from dask[complete])
  Downloading lz4-4.4.4.3-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (3.8 kB)
```

Рисунок 1. Среда и рабочий каталог настроены

```
[4] 1 # import libraries
    2 import sys
    3 import os
    4
    5 ## import dask libraries
    6 import dask.dataframe as dd
    7 from dask.diagnostics import ProgressBar
    8
    9 # import libraries
   10 import pandas as pd

1 cwd = os.getcwd()
2
3 # print
4 print('', sys.executable)
5 print('', cwd)

/usr/bin/python3
/content
```

Рисунок 2. Импортированы библиотеки

4.1.2. Загрузить данные.

2.1.2. Загрузить данные.

```
[9] 1 df = dd.read_csv('/content/drive/MyDrive/Проектный практикум по ETL/Parking_Violations_Issued_-_Fiscal_Year_2017.csv')
    2 df
```

Dask DataFrame Structure:

	Summons Number	Plate ID	Registration State	Plate Type	Issue Date	Violation Code	Vehicle Body Type	Vehicle Make	Issuing Agency	Street Code1	Street Code2	Street Code3	Vehicle Expiration Date	Violation Location	Violation Precinct	Issuer Precinct
npartitions=32	int64	string	string	string	string	int64	string	string	string	int64	int64	int64	int64	float64	int64	int64
	***	***	***	***	***	***	***	***	***	***	***	***	***	***	***	***
	***	***	***	***	***	***	***	***	***	***	***	***	***	***	***	***
	***	***	***	***	***	***	***	***	***	***	***	***	***	***	***	***
	***	***	***	***	***	***	***	***	***	***	***	***	***	***	***	***
Dask Name: to_string_dtype, 2 expressions																

Рисунок 3. Данные загружены

Dask DataFrame отображает только структуру данных, а не сами данные, чтобы избежать загрузки всех данных в память при выводе на экран.

При попытке отобразить первые 10 строк данных вышла ошибка (рисунок 4).

```
ValueError: Mismatched dtypes found in `pd.read_csv`/`pd.read_table`.
```

Column	Found	Expected
House Number	object	float64
Time First Observed	object	float64

The following columns also raised exceptions on conversion:

- House Number
ValueError("could not convert string to float: '150-34'")
- Time First Observed
ValueError("could not convert string to float: '1138A'")

Usually this is due to dask's dtype inference failing, and *may* be fixed by specifying dtypes manually by adding:

```
dtype={'House Number': 'object',  
      'Time First Observed': 'object'}
```

to the call to `read_csv`/`read_table`.

Рисунок 4. Ошибка определения типов данных

```
[9] 1 dtype = {  
    2 |     'House Number': 'object',  
    3 |     'Time First Observed': 'object'  
    4 | }
```

```
[10] 1 df = dd.read_csv(  
    2 |     '/content/drive/MyDrive/Проектный практикум по ETL/Parking_Violations_Issued_-_Fiscal_Year_2017.csv',  
    3 |     dtype=dtype  
    4 | )
```

Рисунок 5. Изменение типа данных

Согласно рекомендации колаб по решению проблемы, ошибка была исправлена и данные перезагружены с учетом изменений.

1 df.head(10)

	Summons Number	Plate ID	Registration State	Plate Type	Issue Date	Violation Code	Vehicle Body Type	Vehicle Make	Issuing Agency	Street Code1	...	Vehicle Color	Unregistered Vehicle?	Vehicle Year	Meter Number	Feet From Curb	Violation Post Code	Violation Description	No Standing or Stopping Violation	Hydrant Violation
0	5092469481	GZH7067	NY	PAS	07/10/2016	7	SUBN	TOYOT	V	0	...	GY	NaN	2001	<NA>	0	<NA>	FAILURE TO STOP AT RED LIGHT	NaN	NaN
1	5092451658	GZH7067	NY	PAS	07/08/2016	7	SUBN	TOYOT	V	0	...	GY	NaN	2001	<NA>	0	<NA>	FAILURE TO STOP AT RED LIGHT	NaN	NaN
2	4006265037	FZX9232	NY	PAS	08/23/2016	5	SUBN	FORD	V	0	...	BK	NaN	2004	<NA>	0	<NA>	BUS LANE VIOLATION	NaN	NaN
3	8478629828	66623ME	NY	COM	06/14/2017	47	REFG	MITSU	T	10610	...	WH	NaN	2007	<NA>	0	04	47-Double PKG-Midtown	NaN	NaN
4	7868300310	37033JV	NY	COM	11/21/2016	69	DELV	INTER	T	10510	...	WHITE	NaN	2007	<NA>	0	316	69-Failure to Disp Muni Recpt	NaN	NaN
5	5096917368	FZD8593	NY	PAS	06/13/2017	7	SUBN	ME/BE	V	0	...	WH	NaN	2012	<NA>	0	<NA>	FAILURE TO STOP AT RED LIGHT	NaN	NaN
6	1413609545	X20DCM	NJ	PAS	08/03/2016	40	SDN	TOYOT	P	54070	...	WHITE	0.0	0	-	1	<NA>	<NA>	NaN	NaN
7	4628525523	326SF9	MA	PAS	12/21/2016	36	UT	BMW	V	0	...	<NA>	NaN	2001	<NA>	0	<NA>	PHOTO SCHOOL ZN SPEED	NaN	NaN

Рисунок 6. Отображение первых 10 строк

4.1.3. Проверить качество данных (например, отсутствующие значения и выбросы).

2.1.3. Проверить качество данных (например, отсутствующие значения и выбросы).

```
[10] 1 # count missing values
      2 missing_values = df.isnull().sum()
      3 missing_values
```

↗ Dask Series Structure:

```
npartitions=1
Date First Observed    int64
Violation Time          ...
Dask Name: sum, 5 expressions
Expr=(~ NotNull(frame=ArrowStringConversion(frame=FromMapProjectable(4254b71))))).sum()
```

Рисунок 7. Проверка качества данных

4.1.4. Удалить столбцы (множество пропусков в значениях, бесполезные столбцы для анализа).

```
[15] 1 result = df.compute()
      2 print(result)
```

↗ /usr/local/lib/python3.11/dist-packages/dask/dataframe/io/csv.py:199: DtypeWarning: Columns (18,38) have mixed types. Specify dtype option on import or set low_memory=False.

```
df = reader(bio, **kwargs)

Summons Number Plate ID Registration State Plate Type Issue Date \
0      5092469481  GZH7067          NY      PAS    07/10/2016
1      5092451658  GZH7067          NY      PAS    07/08/2016
2      4006265037  FZX9232          NY      PAS    08/23/2016
3      8478629828  66623ME          NY      COM    06/14/2017
4      7868300310  37033JV          NY      COM    11/21/2016
...      ...      ...
352313  1415891400  HGK6453          NJ      PAS    11/02/2068
352314  1384716543  GRA6240          NY      PAS    07/12/2069
352315  1413536554  RC8528          PA      PAS    08/14/2069
352316  1415514203  HGU9544          NY      PAS    11/15/2069
352317  141595370   GPP1608          NY      PAS    11/19/2069

Violation Code Vehicle Body Type Vehicle Make Issuing Agency \
0              7          SUBN      TOYOT      V
1              7          SUBN      TOYOT      V
2              5          SUBN      FORD      V
3              47         REFG      MITSU      T
4              69         DELV      INTER      T
...      ...      ...
352313          21          SDN      HONDA      S
352314          20          SUBN      TOYOT      X
352315          46          SUBN      <NA>      P
352316          40          SUBN      JEEP      P
```

Рисунок 8. Выведение результата

```

[16] 1 # calculate percent missing values
      2 mysize = df.index.size
      3 missing_count = ((missing_values / mysize) * 100)
      4 missing_count

Dask Series Structure:
npartitions=1
Date First Observed    float64
...
Violation Time         ...
Dask Name: mul, 9 expressions
Expr=(~ NotNull(frame=ArrowStringConversion(frame=FromMapProjectable(4d4d1e8))).sum() / Index(frame=ArrowStringConversion(frame=FromMapProjectable(4d4d1e8))).size()) * 100

```

Рисунок 9. Расчет процента пропущенных значений

1 # запуск вычисления, используя метод подсчета	
2 with ProgressBar():	
3 missing_count_percent = missing_count.compute()	
4 missing_count_percent	
	Violation Location 19.183510
	Violation Precinct 0.000000
	Issuer Precinct 0.000000
	Issuer Code 0.000000
	Issuer Command 19.093212
	Issuer Squad 19.101506
	Violation Time 0.000583
	Time First Observed 92.217488
	Violation County 0.366073
	Violation In Front Of Or Opposite 20.005826
	House Number 21.184968
	Street Name 0.037110
	Intersecting Street 68.827675

Рисунок 10. Вычисление с помощью метода подсчёта

2.1.4. Удалить столбцы (пропуски в значениях, бесполезные столбцы для анализа).

```

1 # Получаем список столбцов, которые нужно удалить
2 columns_to_drop = missing_count_percent[missing_count_percent > 60].index
3
4 # Проверяем, что columns_to_drop не пуст
5 if not columns_to_drop.empty:
6     print("Столбцы для удаления:", columns_to_drop)
7
8     # Удаляем столбцы
9     with ProgressBar():
10         df_dropped = df.drop(columns=columns_to_drop).compute()
11 else:
12     print("Нет столбцов для удаления.")
13     df_dropped = df.compute()

Столбцы для удаления: Index(['Time First Observed', 'Intersecting Street', 'Violation Legal Code',
                              'Unregistered Vehicle?', 'Meter Number',
                              'No Standing or Stopping Violation', 'Hydrant Violation',
                              'Double Parking Violation'],
                              dtype='object')

[#####] | 18% Completed | 29.94 s/usr/local/lib/python3.11/dist-packages/dask/data
df = reader(bio, **kwargs)
[#####] | 100% Completed | 132.93 s

```

Рисунок 11. Удаление бесполезных столбцов и пропущенных значений

Удаление столбцов осуществляется при условии процента пустых значений больше 60. Чтобы избежать ошибки была добавлена проверка `if not columns_to_drop.empty`, чтобы убедиться, что `columns_to_drop` не пуст перед выполнением операции `drop`. Это предотвращает попытку удаления несуществующих столбцов.

4.2.1. Визуализировать DAG с одним узлом и зависимостями.

Пример

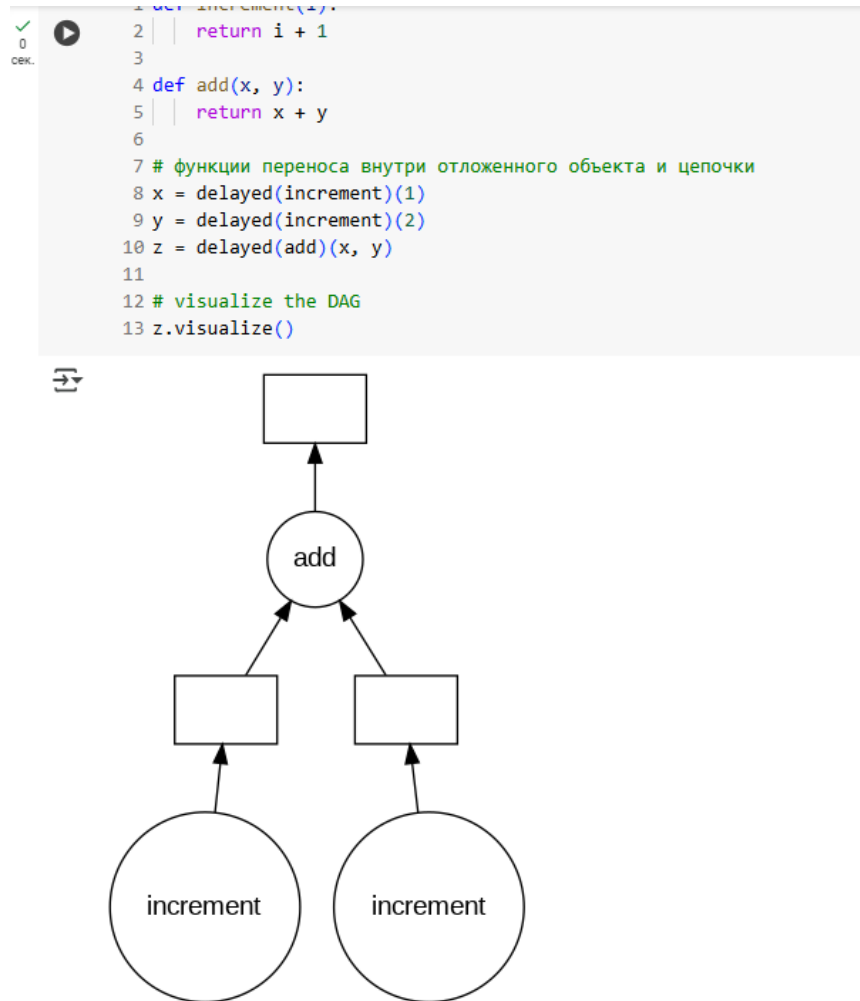


Рисунок 12. Пример визуализации DAG с одним узлом и зависимостями

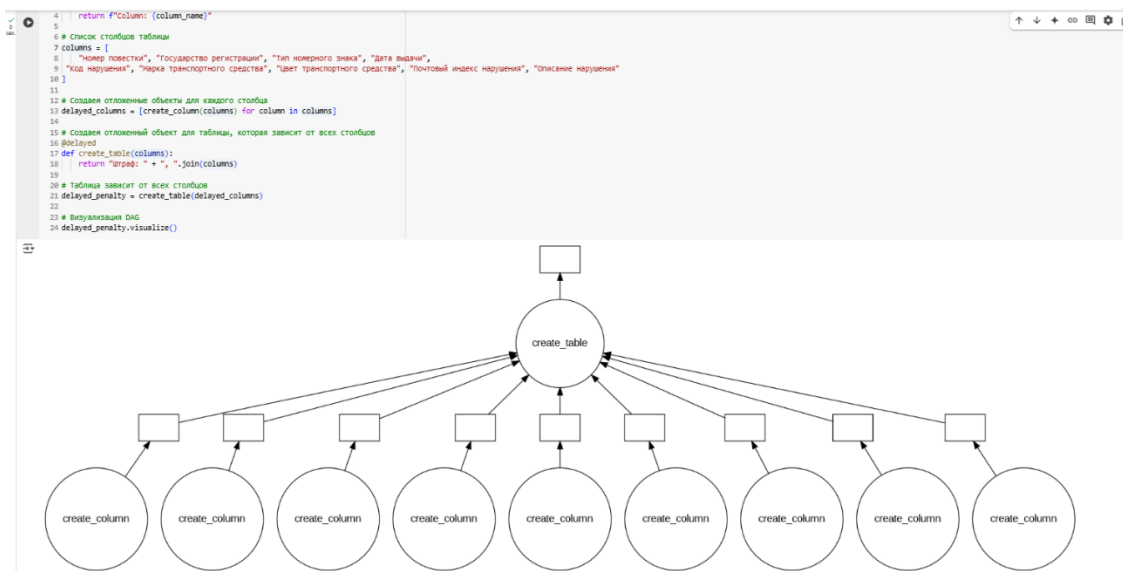


Рисунок 13. Свой вариант визуализации DAG с одним узлом и зависимостями

4.2.1. Визуализировать DAG с более чем одним узлом и зависимостями.

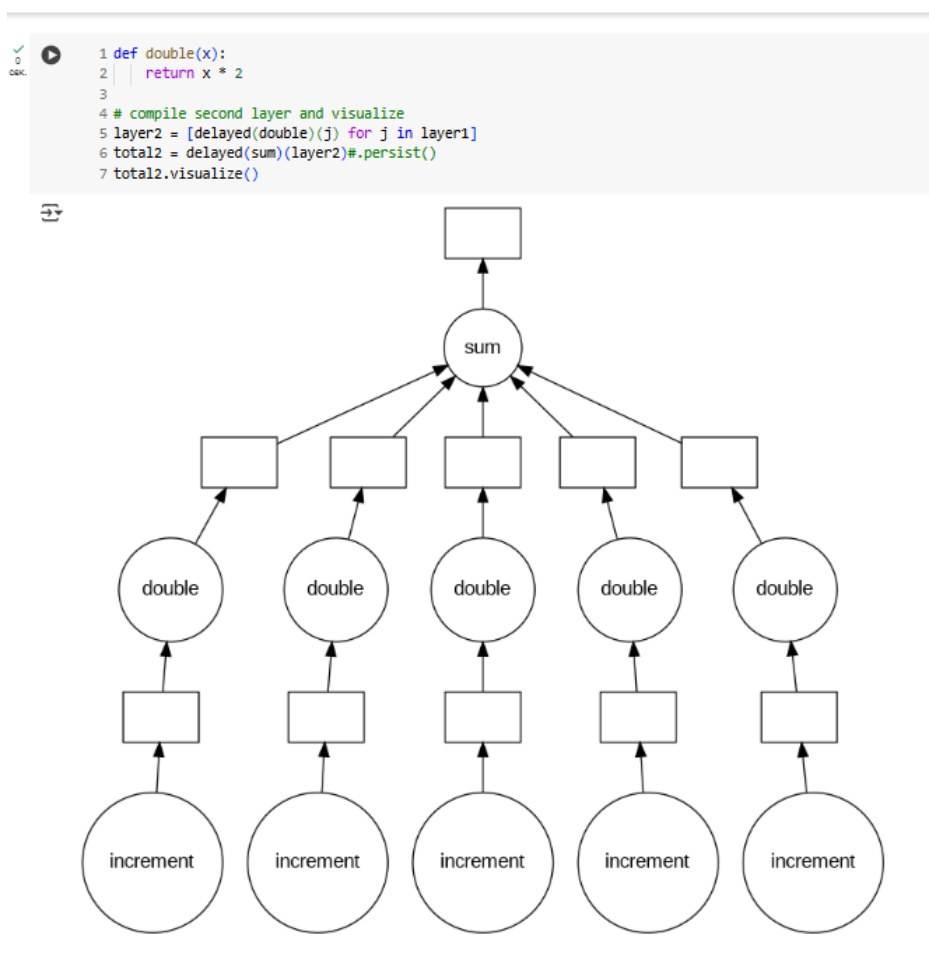


Рисунок 14. Пример визуализации DAG с более чем одним узлом и зависимостями

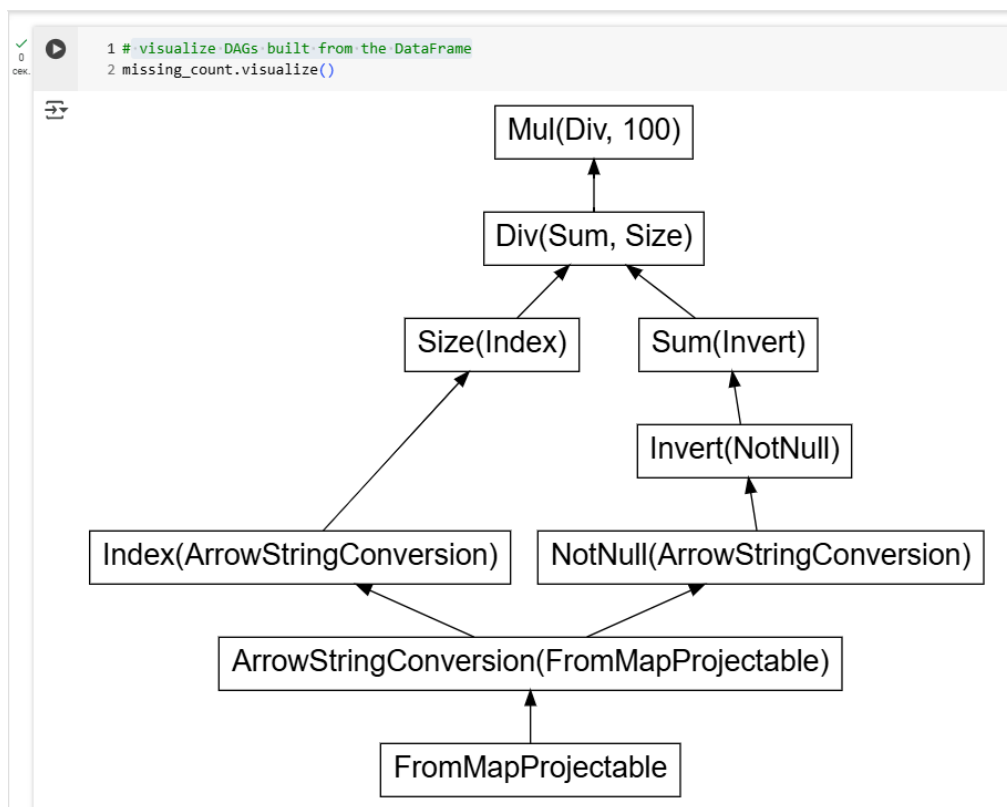


Рисунок 15. Визуализация, построенная на основании data frame