

# Министерство науки и высшего образования Российской Федерации Федеральное государственное бюджетное образовательное учреждение

## высшего образования Лосковский госуларственный технический

# «Московский государственный технический университет имени Н.Э. Баумана

(национальный исследовательский университет)» (МГТУ им. Н.Э. Баумана)

# ФАКУЛЬТЕТ **ИНФОРМАТИКА, ИСКУССТВЕННЫЙ ИНТЕЛЛЕКТ И СИСТЕМЫ** УПРАВЛЕНИЯ

КАФЕДРА КОМПЬЮТЕРНЫЕ СИСТЕМЫ И СЕТИ (ИУ6)

НАПРАВЛЕНИЕ ПОДГОТОВКИ **09.04.01 Информатика и вычислительная техника** МАГИСТЕРСКАЯ ПРОГРАММА **09.04.01/07 Интеллектуальные системы анализа, обработки и интерпретации больших данных** 

# ОТЧЕТ

по лабораторной работе №10

Название:	Работа со Spark		
Дисциплина	: Языки программи данными	рования для работн	ы с большими
Студент	<u>ИУ6-22М</u> (Группа)	(Подпись, дата)	М.А. Зотов (И.О. Фамилия)
Преподавате	ель	(Подпись, дата)	П.В. Степанов (И.О. Фамилия)

**Цель:** получить опыт работы со Spark.

#### Залание:

- 1. Выбрать любой датасет на kaggle.com
- 2. Сделать 10 выборок данных по выбранной предметной области

### Выполнение:

Был выбран датасет «Статистика по пользователям интернета в мире» (https://www.kaggle.com/datasets/ashishraut64/internet-users).

Подключение Spark и предобработка датасета:

# Запросы:

1. Среднее количество пользователей интернета по годам

Рисунок 1 – Результат работы 1 запроса

2. Прирост количества пользователей интернета относительно прошлого года

```
, 2) as internet users num increase "
"from tmp "
"order by 1").show()
          +---+
          |year|internet_users_num_increase|
          +---+
          |1990|
                              227402.89
          |1991|
                              116022.89|
          |1992|
                              175923.67
          1993
                              174597.56
          1994
                              805490.33|
          |1995|
                             1377929.44|
          1996
                             2286187.89|
          11997
                             1974325.11
```

Рисунок 2 – Результат работы 2 запроса

3. Страны, где доля пользователей интернета с 2000 по 2020 год была выше 40%

```
spark.sql("select entity, round(avg(internet users perc), 2)
avg internet users perc "
         "from src "
         "where internet users perc > 40 "
         " and code <> 'Region' "
            and year between 2000 and 2020 "
         "group by 1 "
         "having count(*) = 21 "
         "order by 2 desc").show(truncate=False)
                  |entity
                              |avg_internet_users_perc|
                  +----+
                  |Iceland
                              88.66
                             87.32
                  |Norway
                                                   Sweden
                           84.83
                  |Netherlands |82.26
                  |South Korea |80.11
                                                   Т
                              178.99
                  l Canada
                  . . . .
```

Рисунок 3 – Результат работы 3 запроса

4. Список регионов по максимальной доле пользователей мобильной связи

+++			
num entity		avg_mobile_users	
++			
1	European Union	74.98	
12	High income	71.85	
3	Europe and Central Asia	70.25	
4	North America	61.98	
5	Latin America and Caribbean	52.76	

Рисунок 4 – Результат работы 4 запроса

5. Страны с долей пользователей интернета за 2020 год выше среднего

```
spark.sql("select entity, internet users perc "
         "from src "
         "where year = 2020"
         " and code <> 'Region' "
" and internet_users_perc > (select avg(internet_users_perc)
from src where year = 2020) "
         "order by 2 desc").show(truncate=False)
                      +----+
                                  |internet_users_perc|
                      +-----
                      |Bahrain
                                  99.66999817
                      |Qatar
                                 99.65284729
                      |Kuwait
                                 |99.10588074
                                 199
                      lIceland
                      |Luxembourg | 98.82242584
                      |Saudi Arabia | 97.86000061
```

Рисунок 5 – Результат работы 5 запроса

6. Года, когда в России доля пользователей интернета было ниже среднего по миру

```
spark.sql("with tmp as ("
              select year, avg(internet users perc) as
avg internet users perc "
          " from src "
" group by 1) "
          "select src.year "
          "from src join tmp "
          "on src.year = tmp.year "
          "where entity = 'Russia' "
          " and internet users_perc < avg_internet_users_perc "</pre>
          "order by 1").show()
                                       +---+
                                       |year|
                                       |1990|
                                       11991
                                       |1992|
                                       |1993|
                                       |1994|
                                       |1995|
```

Рисунок 6 – Результат работы 6 запроса

#### 7. Топ 3 страны по доле пользователей мобильной связи

```
spark.sql("with tmp as ("
             select entity, max(mobile users) as mobile users "
             from src "
            where code <> 'Region' "
        group by 1)"
        "select src.entity, src.year, src.mobile users "
        "from src join tmp "
        "on src.entity = tmp.entity and src.mobile users = tmp.mobile users
        "order by 3 desc "
        "limit 3").show(truncate=False)
                      +----+
                      entity
                               |year|mobile_users|
                      +-----
                               |2018|99.97170258 | |
                      |Mongolia |2013|99.85401154 |
                      |Uzbekistan|2020|99.75439453 |
                      +-----
```

Рисунок 7 – Результат работы 7 запроса

8. Количество стран, у которых доля пользователей интернета была выше 30, по годам

```
spark.sql("select year, count(*) as count of countries "
         "from src "
         "where internet users perc > 30 "
         " and code <> 'Region' "
         "group by 1 "
         "order by 1").show()
                         +---+
                         |year|count_of_countries|
                         +----+
                         1998
                                             4
                         |1999|
                                            12|
                         120001
                                            20|
                         120011
                                            271
                                            291
                         2002
                         2003
                                            36 l
```

Рисунок 8 – Результат работы 8 запроса

9. Последние страны подключившиеся к интернету

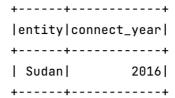


Рисунок 9 – Результат работы 9 запроса

# 10. Страны-рекордсменки по доле пользователей интернета по годам

```
spark.sql("select distinct year, "
        " first value(entity) over (partition by year order by
internet users perc desc) as country, "
       max(internet_users_perc) over (partition by year) as
internet users perc "
        "from src "
        "order by 1").show(truncate=False)
            +---+
            |year|country
                                    |internet_users_perc|
            +----+
                                    0.784728527
            |1990|United States
                                   |5.22E-05
|8.55E-05
|2.78399086
            |1991|Thailand
            |1992|South Asia
            |1993|Norway
            |1994|Iceland
                                   6.794811726
                              19.237088203
            |1995|United States
            |1996|Netherlands
                                     9.649068832
```

Рисунок 10 – Результат работы 10 запроса

**Вывод:** в ходе выполнения лабораторной работы был получен опыт работы со Spark.