

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

«Московский государственный технический университет имени Н.Э. Баумана (национальный исследовательский университет)» (МГТУ им. Н.Э. Баумана)

ФАКУЛЬТЕТ ИНФОРМАТИКА 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

Вариант № 5

Название лабораторной работы: Spark

Дисциплина: Языки программирования для работы с большими данными

Студент гр.	ИУ6-22М		Р.Г. Гаделия
		(Подпись, дата)	(И.О. Фамилия)
Преподавател	Ь		П.В. Степанов
_		(Подпись, дата)	(И.О. Фамилия)

Введение

Целью лабораторной работы является приобретение навыков работы со Spark на языке программирования Java.

Практическая часть

Задание 1

- 1) Выбрать любой датасет на kaggle.com
- 2) Сделать 10 выборок данных по выбранной предметной области Код написанной программы представлен в листинге 1.

Листинг 1 – Программа для первого задания

```
package org.example;
import org.apache.spark.sql.Dataset;
import org.apache.spark.sql.Row;
import org.apache.spark.sql.SparkSession;
public class Main {
    public static void main(String[] args) {
        SparkSession spark = SparkSession.builder()
                .appName("Lab10")
                .master("local")
                .getOrCreate();
        Dataset<Row> flights = spark.read()
                .option("header", true)
                .option("inferSchema", true)
.csv("/Users/ivangorshkov/Documents/BMSTU/BigDataJava/Java-IU6-
12M/Lab10/russian air service CARGO AND PARCELS.csv");
        flights.createOrReplaceTempView("flights");
        spark.sql("SELECT * FROM flights").show();
        spark.sql("select * from flights where Year = '2020'").show();
        spark.sql("select * from flights where Year = '2010'").show();
        spark.sql("select * from flights where May <</pre>
September").show();
        spark.sql("select * from flights where AirportName =
'Trip'").show();
        spark.sql("select AirportName, Year from flights order by
February desc").show();
        spark.sql("select AirportName, Year, March from flights where
January > 150.0 order by November desc").show();
        spark.sql("select AirportName, Year, April from flights order
by May desc").show();
        spark.sql("select AirportName, Year, May, June from flights
where Year = '2014' and AirportName = 'Trip' order by May
desc").show();
        spark.sql("select avg(January) from flights").show();
        spark.stop();
```

Результат выполнения программы показан на рисунке 1.

	.+												
/ AirportName Year	∖January f ·+	-ebruary +	March <i>A</i>	April Ma +	y June ++	July A	ugust Se +	ptember 00	tober No	vember D +-	ecember Wi	iole yea	ar Airport coordinates +
Abakan 2020		9.0	8.8		.01 0.01		8.8	0.0	0.01	0.0	0.01		0 (Decimal('91.3997
Aikhal 2026 Loss 2026		0.0 0.0			0.8 8. 0.8 8.		0.0 0.0	0.0 0.0	0.0 0.0	0.8 0.8	0.0 0.0		0 (Decimal('111.543 0 (Decimal('125.398
Amderma 2020	8.81	0.0			.0 0.0		0.0	8.01	0.01	0.01	0.01		0 (Decimal('61.5774
Anadyr (Carbon) 2026 Anapa (Vitiazevo) 2026		0.0 0.0			.0 0.0 0.0 0.		0.0 0.0	0.0 0.0	0.0 0.0	0.0 0.0	0.0 0.0		0 (Decimal('177.738 .0 (Decimal('37.3415
Apatite (Khibiny) 2026			8.8		.8 6.8		0.0	8.0	0.01	8.8	0.0		0 (Decimal('33.5819
Arkhangelsk (Vask 2020			8.81		.0 0.0		0.0	8.01	0.0	0.01	0.01		0 (Decimal('40.7867
Arkhangelsk (Talagy) 2020 Astrakhan (Narima 2020			0.0 0.0		0.0 0.0 0.		0.0 0.0	0.0 0.0	0.0 0.0	0.0 0.0	0.0 0.0		0 (Decimal('40.7148 .0 (Decimal('47.9998
Trip 2026			8.81	0.0 0.	.01 0.01	0.01	0.0	0.0	8.8	0.0	0.01		0 (Decimal('138.042
Baykit 2020 Barnaul (Titov Name) 2020			0.0 0.0		.0.0 0.0 .0.0 0.0		0.0 0.0	0.0 0.0	0.0 0.0	0.0 0.0	0.0 0.0		0 (Decimal('96.3667 .0 (Decimal('83.5477
In Salah 2026					0.0		0.0	8.01	0.01	0.01	0.01		0 (Decimal('130.399
white Mountain 2020 Belgorod 2020					.0.0 0.0 .0.0 0.		0.0 0.0	0.0 0.0	0.0 0.0	0.0 0.0	0.0 0.0		0 (Decimal('146.228 .0 (Decimal('36.5705
Novy Urengoy 2026					.0 0.0		0.0	8.01	0.01	0.0	0.01		0 (Decimal('66.6945
Belushi 2020 Usinsk 2020		0.0 0.0	0.0 0.0		0.8 8. 0.8 8.		0.0 0.0	0.0 0.0	0.0 0.0	0.0 0.0	0.0 0.0		0 (Decimal('47.6234 0 (Decimal('65.0461
Beringovskiy 2020			8.8		.0 6.0		0.0	8.01	0.0	0.0	0.0		0 (Decimal('179.293
+													
only showing top 20 rows													
AirportName Year													++ ar Airport coordinates
+													
Abakan 2018 Aikhal 2018		192.2 5		8.7 118 0.8 8				181.4	88.5 0.8	83.3 0.8	110.6		.8 (Decimal('91.3997 .0 (Decimal('111.543
Loss 2010	0.01	0.0	0.0	8.8 8	.0 0.6	0.0	8.0	0.01	8.8	0.01	0.0		.0 (Decimal('125.398
Amderma 2018 Anadyr 2018		0.0 119.0 15		0.8 0 7.8 213				0.0 130.0	0.8 191.8	0.0 147.8	0.0 352.0		.0 (Decimal('61.5774 .0 (Decimal('177.738
Anapa (Vitjazevo) 2010					66 59.19	189.76		62.42	47.53	31.6	29.72		48 (Decimal('37.3415
Antypayuta 2010	0.01				.0 0.6 .0 0.6			0.0 0.0	0.8 0.8	0.8 0.8	0.8 0.8		.0 Not found .0 (Decimal('33.5819
Apatite (Khibiny) 2018 Arkhangelsk (Vask 2018					.0 0.6 .0 0.6			0.0 	8.81	0.0 0.8	0.8		.0 (Decimal('40.7067
Arkhangelsk (Talagy) 2010	122.8	156.8 18						200.22 76.19	191.3	353.91 63.541	479.2 71.06		52 (Decimal('40.7148
Astrakhan (Narima 2018 Achinsk 2018				97 51. 0.8 8				76.19 0.8	71.23 8.8	63.54 0.8	71.06		89 (Decimal('47.9998 .0 (Decimal('90.5664
Trip 2018			0.0	8.8 8	.0 0.0			9.8	8.8	0.0	0.8		.0 (Decimal('138.042 .0 (Decimal('113.479
Auctions 2010 Baykit 2010					.0 0.6 .0 0.6			0.0 0.0	8.8 8.8	0.8 0.8	0.8 0.8		.0 (Decimal('113.479 .0 (Decimal('96.3667
Barnaul (Titov Name) 2018	0.01	0.8	0.0	8.8 8	.0 0.0	0.8	8.0	0.0	0.81	0.01	0.8		.0 (Decimal('83.5477
In Salah 2010 white Mountain 2010				0.8 0 0.8 0	.0 0.6 .0 0.6			0.8 0.8	8.8 8.8	0.8 0.8	0.8 0.8		.0 (Decimal('130.399 .0 (Decimal('146.228
Belgorod 2010	0.01	0.8	0.0	0.0 0	.0 0.0	0.0	8.9	0.0	0.81	0.01	0.8		.0 (Decimal('36.5705
Novy Urengoy 2018	15.5 ++	15.1 1	16.7 4 +-	7.7 66	.9 52.7 +	7 49.1 +	43.2 +	57.0 	54.8 +-	85.8	55.1 +-		.6 (Decimal('66.6945
only showing top 20 rows													
	+	+-	+		-++-	+			+	+	+		++
AirportName Year	January F +	ebruary M	iarch A		/ June J -++-		gust Sep	tember Oct 	ober Nove	ember Dec	ember Who	le year	Airport coordinates
Abakan 2028	43.58				0.01		0.0	0.0	0.0	8.0	8.8		(Decimal('91.3997
Aikhal 2028 Loss 2028		0.0 0.0	0.0 0.0		9 0.8 9 0.8		0.8 0.8	0.0 0.0	0.0 0.0	8.0 8.0	8.8 8.8		(Decimal('111.543 (Decimal('125.398
Amderma 2028		0.01	8.8	0.8 8.6			0.0	0.0	0.0	8.0	0.0		(Decimal('61.5774
Anadyr (Carbon) 2028		0.0	0.01	0.0 8.0	9.0	9.8	0.0	0.0	0.0	8.0	0.0		(Decimal('177.738
Anadyr (Larbon) 2026 Anapa (Vitjazevo) 2028 Apatite (Khibiny) 2028	49.08	0.0 0.0 0.0	0.0 0.0 0.0	0.8 8.6	9.0	8.8	0.0 0.0 0.0	0.0 0.0 0.0	0.0 0.0 0.0	8.0 8.0 8.0	8.8 8.8 8.8	0.0	(Decimal('177.738 (Decimal('37.3415 (Decimal('33.5819
Anapa (Vitjazevo) 2028 Apatite (Khibiny) 2028 Arkhangelsk (Vask 2028	49.88 0.8 0.8	0.8 0.8 0.8	0.0 0.0 0.0	0.8 8.0 0.8 8.0 0.8 8.0	0.0 0.0 0.0 0.0 0.0	0.0 0.0 0.0	0.0 0.0 0.0	0.0 0.0 0.0	0.0 0.0 0.0	6.0 6.0 8.0	8.8 8.8 8.8	6.8 6.8	(Decimal('37.3415 (Decimal('33.5819 (Decimal('40.7067
Anapa (Vitjazevo) 2028 Apatite (Khibiny) 2028	49.88 0.8 0.8	0.0 0.0	0.8 0.8	0.8 8.0 0.8 8.0 0.8 8.0	9 0.8 3 0.8 3 0.8	0.0 0.0 0.0 0.0	0.0 0.0	0.0 0.0	0.0 0.0	8.0 8.0	8.8 8.8	8.8 8.8 8.8	(Decimal('37.3415 (Decimal('33.5819
Anapa (Vitjazevo) 2028 Apatite (Khibiny) 2028 Arkhangelsk (Vask 2028 Arkhangelsk (Talagy) 2028 Astrakhan (Narima 2028 Trip 2028	49.08 0.0 0.0 96.0 53.55	0.0 0.0 0.0 0.0 0.0	8.81 8.81 8.81 8.81 8.81	0.8 8.6 0.8 8.6 0.8 8.6 0.8 8.6	3 0.8 3 0.8 3 0.8 3 0.8 3 0.8 3 0.8	8.81 8.81 8.81 8.81 8.81	0.0 0.0 0.0 0.0 0.0 0.0	0.0 0.0 0.0 0.0 0.0	0.0 0.0 0.0 0.0 0.0	6.0 6.0 6.0 6.0 6.0 6.0	8.8 8.8 8.8 8.8 8.8	8.8 8.8 8.8 8.8	(Decimal('37.3415 (Decimal('33.5819 (Decimal('40.7867 (Decimal('40.7148 (Decimal('47.9998 (Decimal('138.842
Anapa (Vitjazevo) 2028 Apatite (khibiny) 2028 Arkhangelsk (Vask.,12028 Arkhangelsk (Talagy) 2028 Astrakhan (Narima 2028 Trip 2028 Baykit 2028	49.88 0.8 0.8 96.8 53.55 1.5	0.8 0.8 0.8 0.8	8.8 8.8 8.8 8.8	0.8 8.6 0.8 8.6 0.8 8.6 0.8 8.6 0.8 8.6	0.0 0.0 0.0 0.0 0.0 0.0	0.01 0.01 0.01 0.01 0.01 0.01	0.0 0.0 0.0 0.0	0.0 0.0 0.0 0.0 0.0	0.0 0.0 0.0 0.0 0.0 0.0	6.0 6.0 6.0 6.0 6.0 6.0	8.6 8.6 8.6 8.6 8.6 8.6	6.6 6.6 6.6 6.6 6.8	(Checimal('37.3415 (Decimal('33.5819 (Decimal('48.7967 (Decimal('48.7148 (Decimal('47.9998
Anapa (Vitjazevo) 2028 Apatite (khibiny) 2028 Arkhangelsk (Vask 2028 Arkhangelsk (Talagy) 2028 Astrakhan (Narisa 2028 Trja 2028 Baykit 2028 Barnaul (Titov Name) 2020 In Salah 2028	49.88 0.8 0.8 96.8 53.55 1.5 0.8 163.4 0.8	0.0 0.0 0.0 0.0 0.0 0.0 0.0	8.61 8.61 8.61 8.61 8.61 8.61 8.61	0.8 8.6 0.8 8.6 0.8 8.6 0.8 8.6 0.8 8.6 0.8 8.6	3 0.0 3 0.0	8.81 8.81 8.61 8.61 8.61 8.61 8.61	0.8 0.8 0.8 0.8 0.8 0.8 0.8	0.0 0.0 0.0 0.0 0.0 0.0 0.0	6.0 6.0 6.0 6.0 6.0 6.0 6.0 6.0	6.0 6.0 6.0 6.0 6.0 6.0 6.0	6.6 6.6 6.6 6.6 6.6 6.6 6.6 6.6	6.6 6.6 6.6 6.6 6.6 6.6	(Decimal('37.3415 (Decimal('33.5819 (Decimal('48.7867 (Decimal('48.7148 (Decimal('47.9998 (Decimal('47.9998 (Decimal('98.3667 (Decimal('98.3667 (Decimal('38.399 (Decimal('138.399
Anapa (Vitjazevo) 2028 Apaite (Khibiny) 2028 Arkhangelsk (Vask 2020 Arkhangelsk (Talagy) 2028 Astrakhan (Narima 2020 Baykit 2020 Baykit 2020 Barnaul (Titov Name) 2020 In Salah 2020 white Hountain 2020	49.88 0.8 0.8 0.8 0.8 96.8 53.55 1.5 0.8 163.4 0.8 0.8 0.8	0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0	0.61 0.61 0.61 0.61 0.61 0.61 0.61	0.8 8.6 0.8 8.6 0.8 8.6 0.8 8.6 0.8 8.6 0.8 8.6 0.8 8.6	3 0.0 3 0.0	6.6 6.6 6.6 6.6 6.6 6.6 6.6 6.6	0.8 0.8 0.8 0.8 0.8 0.8 0.8 0.8	0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0	8.9 8.9 8.9 8.9 9.9 8.9 8.9 8.9	6.9 6.0 6.0 6.0 6.0 6.0 6.0 6.0 6.0 6.0	6.6 6.6 6.6 6.6 6.6 6.6 6.6 6.6	6.6 6.6 6.6 6.6 6.6 6.6	[Oecinal('37.3415 [Oecinal('33.5819 [Oecinal('48.7867 [Oecinal('48 (Oecinal('48 (Oecinal('47.9998 [Oecinal('58.3642 (Oecinal('38.3642 (Oecinal('38.3647 (Oecinal('38.3647 (Oecinal('438.399 (Oecinal('446.228
Anapa (Vitjazevo) 1202 Apatie (Khibiny) 1202 Arkhangelsk (Vask 12020 Arkhangelsk (Talagy) 1202 Astrakhan (Nariam 12020 Barkil 12020 Baykil 12020 Barnaul (Titov Name) 12020 In Salah 12020 Bathe 10untain 12020 Betgordi 12020 Novy Urengoy) 12020 Novy Urengoy) 12020	49.88 0.8 0.8 0.8 0.8 96.8 1.5 1.5	0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0	8.6 8.6 8.6 8.6 8.6 8.6 8.6 8.6	0.010.1 0.010.1 0.010.1 0.010.1 0.010.1 0.010.1 0.010.1 0.010.1 0.010.1 0.010.1	9 0.0 0.	8.8 8.8	0.8 0.8 0.8 0.8 0.8 0.8 0.8 0.8	0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0	6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8	6.0 6.0 6.0 6.0 6.0 6.0 6.0 6.0	6.0 6.0 6.0 6.0 6.0 6.0 6.0 6.0	6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8	[Operimal(137, 3415] [Operimal(148, 7867] [Operimal(148, 7867] [Operimal(148, 7348] [Operimal(148, 7348] [Operimal(148, 847] [Operimal(138, 842] [Operimal(148, 847]
Anapa (Vitjazevo) 12026 Apatite (Khibiny) 12026 Arkhangelsk (Vask12026 Arkhangelsk (Talagy) 12026 Astrakhan (Narima12026 Baylitl 2026 Baynaul (Titov Name) 12026 In Salahi 12026 white Mountain 12026 Belgorod 12028 Novy Unengoy 12026 Belshil 12026	49.08 0.0 0.0 96.0 53.55 1.5 0.0 163.4 0.8 0.8 5.24 5.24	0.01 0.01 0.01 0.01 0.01 0.01 0.01 0.01 0.01 0.01 0.01	8.8 8.8 8.8 8.8 8.8 8.8 8.8 6.8 8.8 8	0.010.1 0.010.1 0.010.1 0.010.1 0.010.1 0.010.1 0.010.1 0.010.1 0.010.1 0.010.1 0.010.1	3 0.0 3 0.	8.8 8.8	0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0	6.0 6.0 6.0 6.0 6.0 6.0 6.0 6.0	0.0 0.0	6.6 6.6 6.6 6.6 6.6 6.6 6.6 6.6	6.91 6.91 6.91 6.91 6.91 6.91 6.91 6.91 6.91 6.91	6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8	[Oecinal('37, 3415] [Oecinal('37, 3415] [Oecinal('48, 7867] [Oecinal('48, 7867] [Oecinal('48, 7348] [Oecinal('38, 842] [Oecinal('338, 842] [Oecinal('338, 842] [Oecinal('338, 849] [Oecinal('338, 399] [Oecinal('348, 258] [Oecinal('358, 5795] [Oecinal('36, 5795]
Anapa (Vitjazevo) 1202 Apaite (Khibiny) 1202 Arkhangelsk (Vask1202 Arkhangelsk (Talagy) 1202 Arkhangelsk (Talagy) 1202 Astrakhan (Nariam1202 Baykit 1202 Baykit 1202 Barnaul (Titov Name) 1202 In Salah 1202 White Hountain 1202 Belgorod 1202 Roy Urengoy 1202 Belushi 1202 Belushi 1202 Belushi 1202 Belushi 1202 Beringovskiy 1202	49.88 0.8 0.8 96.8 53.55 1.5 0.8 163.4 0.8 0.8 0.8 0.8 0.8	0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0	8.6 8.6 8.6 8.6 8.6 8.6 8.6 8.6	0.818.1 0.818.1 0.818.1 0.818.1 0.818.1 0.818.1 0.818.1 0.818.1 0.818.1 0.818.1 0.818.1 0.818.1	9 0.0 0.	8.6 8.6 8.6 8.6 8.6 8.6 8.6 8.6	0.8 0.8 0.8 0.8 0.8 0.8 0.8 0.8	0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0	6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8	6.0 6.0 6.0 6.0 6.0 6.0 6.0 6.0	6.0 6.0 6.0 6.0 6.0 6.0 6.0 6.0	6.6 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8	[Operimal(137, 3415] [Operimal(148, 7867] [Operimal(148, 7867] [Operimal(148, 7348] [Operimal(148, 7348] [Operimal(148, 847] [Operimal(138, 842] [Operimal(148, 847]
Anapa (Vitjazevo) 1282 Apatie (Khibiny) 1282 Arkhangelsk (Vask., 1282 Arkhangelsk (Talayy) 1282 Astrakhan (Nariam., 1282 Baykit 12828 Baykit 12828 Barnaul (Titov Namo) 1282 In Salah 12828 Bathar 12828 Novy Urengoy 12828 Betapord 2	49.88 0.8 0.8 96.8 53.55 1.5 0.8 163.4 0.8 0.8 0.8 0.8 0.8	6.0 6.0 6.0 6.0 6.0 6.0 6.0 6.0	8.8 8.8 8.8 8.8 8.8 8.8 8.8 8.8	0.818.1 0.818.1 0.818.1 0.818.1 0.818.1 0.818.1 0.818.1 0.818.1 0.818.1 0.818.1 0.818.1 0.818.1	31 0.91 31 0.91	8.6 8.6 8.6 8.6 8.6 8.6 8.6 8.6	0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0	6.9 9.9 9.9 9.9 9.9 9.9 9.9 9.9	9.9 9.9	6.0 6.0 6.0 6.0 6.0 6.0 6.0 6.0	6.91 6.91 6.91 6.91 6.91 6.91 6.91 6.91	6.6 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8	[Oecinal('37, 3415] [Oecinal('38, 589)] [Oecinal('48, 7867] [Oecinal('48, 7148] [Oecinal('78, 7148] [Oecinal('78, 7998] [Oecinal('58, 5867] [Oecinal('58, 5867] [Oecinal('68, 5867] [Oecinal('68, 5867] [Oecinal('66, 6985] [Oecinal('66, 6985] [Oecinal('66, 6985] [Oecinal('66, 6985]
Anapa (Vitjazevo) 1202 Apaite (Khibiny) 1202 Arkhangelsk (Vask1202 Arkhangelsk (Talagy) 1202 Arkhangelsk (Talagy) 1202 Astrakhan (Nariam1202 Baykit 1202 Baykit 1202 Barnaul (Titov Name) 1202 In Salah 1202 White Hountain 1202 Belgorod 1202 Roy Urengoy 1202 Belushi 1202 Belushi 1202 Belushi 1202 Belushi 1202 Beringovskiy 1202	49.08 0.0 96.0 76.0 53.55 1.5 1.63.4 0.0 163.4 0.0 5.24 5.24 5.5 0.0	6.9 6.9 6.0	6.6 8.6	0.0 0.1 0.0 0.1	31 0.01 31 0.01	6.6 6.6 6.6 6.6 6.6 6.6 6.6 6.6	0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0	6.9 6.0 6.0 6.0 6.9 6.9 6.9 6.9 6.9 6.9 6.9 6.9	9.9 9.9 9.9 9.9 9.9 9.9 9.9 9.9	6.0 6.0 6.0 6.0 6.0 6.0 6.0 6.0	8.81 8.01 8.81 8.81 8.81 8.81 8.81 8.81	6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8	[Oecinal('37, 3415] [Oecinal('38, 589)] [Oecinal('48, 7867] [Oecinal('48, 7148] [Oecinal('78, 7148] [Oecinal('78, 7998] [Oecinal('58, 5867] [Oecinal('58, 5867] [Oecinal('68, 5867] [Oecinal('68, 5867] [Oecinal('66, 6985] [Oecinal('66, 6985] [Oecinal('66, 6985] [Oecinal('66, 6985]
Anapa (Vitjazevo) 12026 Apatite (Khibiny) 12026 Arkhangelsk (Vask12026 Arkhangelsk (Talagy) 12026 Arkhangelsk (Talagy) 12026 Baylatil 12026 Baylatil 12026 Barnaul (Titov Name) 12026 In Salahi 12026 White Mountain 12026 Belgardel 12028 Rovy Urengey 12028 Belgardel 12028 Beringovskiyl 2028 Inly showing top 29 rows	49.08 0.08 0.08 96.08 95.08 53.55 0.08 1.51 0.08 1.524 5.24 5.51 0.08 0.08	6.9 6.0	8.8 8.8	6.818.6 6.818.6 6.818.6 6.818.6 6.818.6 6.818.6 6.818.6 6.818.6 6.818.6 6.818.6 6.818.6 6.818.6 6.818.6	31 0.91 31 0.91	6.6 6.6 6.6 6.6 6.6 6.6 6.6 6.6	0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0	6.9 6.0 6.9 6.9 6.9 6.9 6.9 6.9 6.0 6.9 6.0 6.0 6.0	9.9 9.9 9.9 9.9 9.9 9.9 9.9 9.9	8.8 8.8 8.8 8.8 8.8 8.8 8.8 8.8	8.8 8.8	6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8	(Decimal('37, 3415] (Decimal('37, 3415] (Decimal('48, 7867] (Decimal('48, 7867] (Decimal('48, 7868] (Decimal('48, 7868] (Decimal('318, 842] (Decimal('318, 842] (Decimal('318, 5477] (Decimal('318, 5477] (Decimal('48, 5478] (Decimal('48, 5785] (Decimal('48, 5785] (Decimal('47, 6234] (Decimal('47, 6234] (Decimal('47, 6234] (Decimal('47, 6234] (Decimal('47, 6234]
Anapa (Vitjazevo) 1202 Apatite (Khibiny) 1202 Arkhangelsk (Vask 1202 Arkhangelsk (Talayy) 1202 Astrakhan (Narian 1202 Astrakhan (Narian 1202 Barnaul (Titov Name) 1202 In Salah 1202 White Hountain 1202 Belspored 1202 Belspo	49.88 6.8 96.8 96.8 53.55 1.5 0.8 163.4 0.8 163.4 0.8 163.4 0.8 0.8 0.8 0.8 0.8	6.9 6.9	8.8 9.8 9.8 9.8 9.8 9.8 8.8 8.8 9.8 9.8	6.8 6.1 6.8 6.1	31 0.91 31 0.91	6.6 6.6	0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0	6.9 6.9 6.9 6.9 6.9 6.9 6.9 6.9	9.9 9.9 9.9 9.9 9.9 9.9 9.9 9.9	8.8 8.8 8.8 8.8 8.8 8.8 8.8 8.8	e.el e.el e.el e.el e.el e.el e.el e.el	6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8	(Decimal('37, 3415 (Oecimal('36, 7867 (Oecimal('46, 7867 (Oecimal('46, 7867 (Oecimal('46, 7346 (Oecimal('48, 7346 (Oecimal('48, 7346 (Oecimal('33, 842 (Oecimal('33, 842 (Oecimal('33, 847 (Oecimal('34, 547 (Oecimal('44, 228 (Oecimal('44, 228 (Oecimal('47, 6234 (Oeci
Anapa (Vitjazevo) 12026 Apatite (Khibiny) 12026 Arkhangelsk (Vask12026 Arkhangelsk (Talagy) 12026 Arkhangelsk (Talagy) 12026 Baylatil 12026 Baylatil 12026 Barnaul (Titov Name) 12026 In Salahi 12026 White Mountain 12026 Belgardel 12028 Rovy Urengey 12028 Belgardel 12028 Beringovskiyl 2028 Inly showing top 29 rows	49.88 49.88 0.81 0.81 0.81 0.81	6.9 6.0	8.8 8.8	0.010.1 0.010.1 0.010.1 0.010.1 0.010.1 0.010.1 0.010.1 0.010.1 0.010.1 0.010.1	3 0.9 3 0.9	6.6 6.9 6.9 6.9 6.9 6.9 6.9 6.9	0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0	6.9 6.9 6.9 6.9 6.9 6.9 6.9 6.9	9.9 9.9 9.9 9.9 9.9 9.9 9.9 9.9	8.8 8.8 8.8 8.8 8.8 8.8 8.8 8.8	8.8 8.8	6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8	(Decimal('37, 3415] (Decimal('37, 3415] (Decimal('48, 7867] (Decimal('48, 7867] (Decimal('48, 7868] (Decimal('48, 7868] (Decimal('318, 842] (Decimal('318, 842] (Decimal('318, 5477] (Decimal('318, 5477] (Decimal('48, 5478] (Decimal('48, 5785] (Decimal('48, 5785] (Decimal('47, 6234] (Decimal('47, 6234] (Decimal('47, 6234] (Decimal('47, 6234] (Decimal('47, 6234]
Anapa (Vitjazevo) 1222 Apatite (Khibiny) 1202 Aprikangelsk (Vask 1202 Arkhangelsk (Vask 1202 Arkhangelsk (Talagy) 1202 Astrakhan (Narian 1262 Earnaul (Titov Name) 1202 Barnaul (Titov Name) 1202 Betspored 120	49.88 49.88 6.81 6.81 6.81 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85 19.85	6.9 6.9	8.9 8.9 8.1 8.2 8.3 8.4 8.4 8.4	0.9 6.1 0.9 6.1	3) 0.9 4) 0.9 4) 0.9 6 0.9 6 0.9	8.8 8.8 9.8	0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0	6.9 6.9 9.9 9.9 9.9 9.9 9.9 9.9	6.9 9.9 9.9 9.9 9.9 9.9 9.9 9.9	6.6 6.6 6.6 6.6 6.6 6.6 6.6 6.6	0.9 0.0 0.0 0.0 0.0 0.0 0.0 0.0	8.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.9 6.9	(Decimal('37, 3415 (Decimal('35, 5819 (Decimal('48, 7867 (Decimal('48, 7867 (Decimal('48, 7868 (Decimal('35, 842 (Decimal('35, 842 (Decimal('35, 842 (Decimal('35, 847 (Decimal('35, 547 (Decimal('36, 548 (Decimal('36, 548 (Decimal('47, 6234 (Decimal('138, 842 (38, 640 (38
Anapa (Vitjazevo) 1202 Apatite (Khibiny) 1202 Arkhangelsk (Vask 1202 Arkhangelsk (Vask 1202 Arkhangelsk (Talagy) 1202 Astrakhan (Narima 1202 Battil 1202 Baykit 1202 Baykit 1202 Barnaul (Titov Name) 1202 Battil 120	49.08 0.01 0.01 9.08 19.08 153.55 1.5 0.01 163.41 0.01 	6.9 9.9 9.9 9.9 9.9 9.9 9.9 9.9	8.8 8.8 8.8 8.8 8.8 8.8 8.8 8.8	0.8 0.1 0.9 0.1	30 0.01 31 0.01 31 0.01 31 0.01 33 0.01 33 0.01 33 0.01 33 0.01 33 0.01 33 0.01 33 0.01 33 0.01 33 0.01 34 0.01 35 0.01 36 0.01 36 0.01 37 0.01 38 0.01 38 0.01 39 0.01 30 0.0	8.81 8.81 8.81 8.81 8.81 8.81 8.81 8.81 8.81 8.81 8.81 8.81 8.81 2.21 2.21 2.21 2.21 2.31	0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0	6.9 6.9 6.9 6.9 6.9 6.9 6.9 6.9	6.9 9.9	6.6 6.6 6.6 6.6 6.6 6.6 6.6 6.6	e.el e.el e.el e.el e.el e.el e.el e.el	8.9 6.9 6.9 6.9 6.9 6.9 6.9 6.9 6.9 7 7 8	(Decimal('37, 3415] (Decimal('38, 5819] (Decimal('48, 7867] (Decimal('48, 7867] (Decimal('48, 7867] (Decimal('48, 7868] (Decimal('78, 7868] (Decimal('38, 842] (Decimal('38, 5477] (Decimal('48, 5477] (Decimal('48, 5477] (Decimal('48, 5478] (Decimal('48, 5785] (Decimal('47, 6234] (Decimal('47, 6234] (Decimal('47, 6234] (Decimal('47, 6234] (Decimal('47, 6234] (Decimal('47, 6234] (Decimal('179, 273] 1.31 (Decimal('138, 842] 0.81 (Decimal('138, 842]
Anapa (Vitjazevo) 1202 Apatite (Khibiny) 1202 Aprikangelsk (Vask 1202 Arkhangelsk (Vask 1202 Arkhangelsk (Talagy) 1202 Astrakhan (Narian 1262 Earnaul (Titov Name) 1202 Barnaul (Titov Name) 1202 Betspored 120	49.88 49.88 0.81 0.81	6.9 6.9	8.8 8.8 9.8 9.8 9.8 8.8 8.8 8.8	0.8 0.1 0.8 0.1	3] e.el 3] e.e	8.81 8.81	0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0	6.9 6.9 6.9 6.9 6.9 6.9 6.9 6.9	6.9 9.9 9.9 9.9 9.9 9.9 9.9 9.9	6.6 6.6 6.6 6.6 6.6 6.6 6.6 6.6	0.01 0.01 0.01 0.01 0.01 0.01 0.01 0.01	6.8 6.8 6.8 6.9 6.9 6.9 6.9 6.9 6.9 6.9 6.9 6.9	(Decimal('37, 3415 (Decimal('35, 5819 (Decimal('48, 7867 (Decimal('48, 7867 (Decimal('48, 7868 (Decimal('35, 842 (Decimal('35, 842 (Decimal('35, 842 (Decimal('35, 847 (Decimal('35, 547 (Decimal('36, 548 (Decimal('36, 548 (Decimal('47, 6234 (Decimal('138, 842 (38, 640 (38
Anapa (Vitjazevo) 2028 Apatite (Khibiny) 2028 Arkhangelsk (Vask 2024 Arkhangelsk (Vask 2024 Arkhangelsk (Talagy) 2028 Astrakhan (Narima 2028 Barhaul (Titov Name) 2028 Barhaul 2028 Barhau	49.88 49.88 6.91	6.9 9.9 9.9 9.9 9.9 9.9 9.9 9.9	9.9 9.0 9.0 9.0 9.0 9.0 9.0 9.0	e.e e.te.e e.te.	3] 0.0 13]	9.91 9.91 9.91 9.91 9.91 9.91 9.91 9.91	9.8 9.8 9.8 9.8 9.8 9.8 9.8 9.8	e.9i e.9i e.9i e.9i e.9i e.9i e.9i e.9i	6.0 6.0 6.0 6.0 6.0 6.0 6.0 6.0	8.0 8.0 8.0 8.0 8.0 8.0 8.0 8.0 8.0 8.0 4.89 8.0 8.0 8.0	e.el e.el e.el e.el e.el e.el e.el e.el	6.8 6.8 6.9 6.8 6.8 6.8 6.8 6.8 6.8 6.8	(Decimal('37, 3415] (Decimal('38, 5819] (Decimal('48, 7867] (Decimal('48, 7867] (Decimal('48, 7867] (Decimal('48, 7868] (Decimal('78, 7868] (Decimal('38, 842] (Decimal('38, 5477] (Decimal('88, 5477] (Decimal('88, 5477] (Decimal('88, 5477] (Decimal('48, 228] (Decimal('48, 228] (Decimal('47, 2334] (Decimal('47, 2334] (Decimal('47, 2334] (Decimal('47, 2334] (Decimal('47, 2334] (Decimal('48, 842] (Decimal('138, 842] (Bel) (Decimal('138, 842]
Anapa (Vitjazevo) 222 Apatite (Khibiny) 222 Apatite (Khibiny) 222 Arkhangelsk (Vask 2028 Arkhangelsk (Vask 2028 Arkhangelsk (Talayy) 222 Astrakhan (Narian 2624 Tripl2228 Barnaul (Titov Name) 222 Barnaul Titov Name Venezal 122 Barnaul Titov Name Venezal	49.88 49.88 19.8	0.9 0.9	6.9 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8	9.810.1 9.810.1		9.91 9.91 9.91 9.91 9.91 9.91 9.91 9.91	9.81 9.81	e.el e.el e.el e.el e.el e.el e.el e.el	9.81 9.81	8.6 0.0	9.61 9.91 9.91 9.91 9.91 9.91 9.91 9.91	6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 7	(Decimal('37, 3415] (Oecimal('38, 5819] (Oecimal('48, 7867] (Oecimal('48, 7148] (Oecimal('48, 7148] (Oecimal('48, 7148] (Oecimal('38, 842] (Oecimal('38, 842] (Oecimal('38, 842] (Oecimal('38, 5477] (Oecimal('38, 5477] (Oecimal('38, 5477] (Oecimal('36, 5485] (Oecimal('36, 5485] (Oecimal('36, 5485] (Oecimal('47, 6234] (Oecimal('47, 6234] (Oecimal('47, 6234] (Oecimal('179, 293] 0ecimal('179, 293] 0ecimal('179, 293] 0ecimal('179, 293] 0ecimal('18, 842] (Oecimal('18, 842] 0.8] (Oecimal('18, 842]
Anana (Vitjazevo) 1222 Apatite (Khibiny) 1202 Apatite (Khibiny) 1202 Arkhangelsk (Vask 1202 Arkhangelsk (Vask 1202 Arkhangelsk (Talayy) 1202 Astrakhan (Nariam 1262 Baykit 1202 Baykit 1202 Baykit 1202 Baykit 1202 Baykit 1202 Belspored 120	49.88 49.88 69.81	0.01 0.01	6.61 6.61 6.61 6.61 6.61 6.61 6.61 6.61	9.810.1 9.8		9.91 9.91	9.81 9.81	9.9 9.9 9.9 9.9 9.9 9.9 9.9 9.9	9.81 9.81	e.ei e.ei e.ei e.ei e.ei e.ei e.ei e.ei	9.8 9.9 9.9 9.9 9.9 9.9 9.9 9.9	6.6 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 7 8.8	(Decimal('37, 3415] (Decimal('38, 5819] (Decimal('48, 7867] (Decimal('48, 7867] (Decimal('48, 7867] (Decimal('48, 7868] (Decimal('38, 842] (Decimal('38, 842] (Decimal('38, 5477] (Decimal('38, 5477] (Decimal('38, 5477] (Decimal('38, 5478] (Decimal('48, 228] (Decimal('48, 228] (Decimal('48, 258] (Decimal('47, 2334] (Decimal('47, 2334] (Decimal('47, 2334] (Decimal('47, 2334] (Decimal('138, 842] (Becimal('138, 842] 8.8] (Decimal('138, 842]
Anapa (Vitjazevo) 222 Apatite (Khibiny) 222 Apatite (Khibiny) 222 Arkhangetsk (Vask 222 Arkhangetsk (Vask 222 Astrakhan (Narian 222 Astrakhan (Narian 222 Barnaul (Titov Namo) 222 Barnaul (Titov Namo) 222 Betigorodi	49.88 49.88 49.88	8.51 9.01 9.01 9.01 9.01 9.01 9.01 9.01 9.0	6.61 6.61 6.61 6.61 6.61 6.61 6.61 6.61	9.91		9.91 9.91	9.81 9.81 9.81 9.81 9.81 9.81 9.81 9.81	0.01	9.81 9.81	8.6 0.0	e.si e.si e.si e.si e.si e.si e.si e.si	6.6 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8	(Decimal('37, 3415] (Oecimal('38, 5819] (Oecimal('48, 7867] (Oecimal('48, 7148] (Oecimal('48, 7148] (Oecimal('48, 7148] (Oecimal('38, 842] (Oecimal('38, 842] (Oecimal('38, 842] (Oecimal('38, 5477] (Oecimal('38, 5477] (Oecimal('38, 5477] (Oecimal('36, 5485] (Oecimal('36, 5485] (Oecimal('36, 5485] (Oecimal('47, 6234] (Oecimal('47, 6234] (Oecimal('47, 6234] (Oecimal('179, 293] 0ecimal('179, 293] 0ecimal('179, 293] 0ecimal('179, 293] 0ecimal('18, 842] (Oecimal('18, 842] 0.8] (Oecimal('18, 842]
Anapa (Vitjazevo) 222 Apatite (Khibiny) 2828 Apatite (Khibiny) 2828 Arkhangelsk (Task 2828 Arkhangelsk (Task 2828 Astrakhan (Nariam 2828 Barhall (Titov Name) 2828 Barhall (Titov Name) 2828 Barhall (Titov Name) 2828 Barhall 2828 Belingovskiy 2828 Belingovskiy 2828 Belingovskiy 2828 In Salel 2828 Belingovskiy 2828 Belingovskiy 2828 In Salel 2828 Belingovskiy 2828 Belin	49.88 49.88 69.81	8.81	6.61 6.61 6.61 6.61 6.61 6.61 6.61 6.61	9.810.1 9.810.1	31 0.01 31 0.0	6.01 6.01 6.01 6.01 6.01 6.01 6.01 6.01	9.81 9.81	6.01 9.01	9.81 9.81	8.ei 0.0	e.si	6.6 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 7 7 8 8	(Decimal('37, 3415] (Decimal('38, 5819] (Decimal('48, 7867] (Decimal('48, 7867] (Decimal('48, 7867] (Decimal('48, 7868] (Decimal('78, 7868] (Decimal('78, 7868] (Decimal('38, 842] (Decimal('38, 5477] (Decimal('38, 5477] (Decimal('38, 5477] (Decimal('38, 5478] (Decimal('38, 5478] (Decimal('48, 228] (Decimal('47, 6234] (Decimal('47, 6234] (Decimal('47, 6234] (Decimal('179, 293] Bearl Airport coordinate 8.8] (Decimal('138, 842]
Anana (Vitjazevo) 222 Apatite (Khibiny) 222 Apatite (Khibiny) 222 Arkhangelsk (Vask., 228 Arkhangelsk (Vask., 228 Arkhangelsk (Talayy) 222 Astrakhan (Nariam., 222 Baykiti2228 Baykiti2228 Barnaul (Titov Namo) 222 Bathi 222	49.88 49.88 49.88 9.01 9.01 9.02	8.51 9.01 9.01 9.01 9.01 9.01 9.01 9.01 9.0	6.61 6.61 6.61 6.61 6.61 6.61 6.61 6.61	9.816.1 9.819.1 9.819.1 9.819.1 9.819.1 9.819.1 9.819.1 9.819.1 9.819.1 9.819.1 9.819.1 9.819.1 9.819.1 9.819.1 9.819.1 9.819.1 9.819.1 9.819.1 9.819.1 9.819.1	31 0.81 31 31 0.81 31 0.81 31 31 0.81 31 31 0.81 31 31 0.81 31 31 0.81 31 31 0.81 31 31 0.81 31 31 0.81 31 31 0.81 31 31 0.81 31 31 0.81 31 31 31 31 31 31 31 31 31 31 31 31 31	e.ei e.ei e.ei e.ei e.ei e.ei e.ei e.ei	9.81 9.81	0.01 0.01	9.81 9.81	e.ei e.ei e.ei e.ei e.ei e.ei e.ei e.ei	e.si e.si	6.8 6.8 6.8 6.9 6.9 6.9 6.9 6.9 6.9 6.9 7 7 8	(Decimal('37, 3415] (Decimal('38, 5812] (Decimal('48, 7867] (Decimal('48, 7146] (Decimal('48, 7146] (Decimal('48, 7146] (Decimal('38, 842] (Decimal('38, 842] (Decimal('38, 842] (Decimal('38, 842] (Decimal('38, 842] (Decimal('38, 842] (Decimal('48, 842] (Decimal('148, 842] (Decimal('148, 842] (Decimal('148, 842] (Decimal('138, 842]
Anana (Vitjazevo) 222 Apatite (Khibiny) 222 Apatite (Khibiny) 222 Arkhangetsk (Vask 222 Arkhangetsk (Vask 222 Astrakhan (Narian 222 Earnal (Titov Namo) 222 Barnal (Titov Namo) 222 Barnal (Titov Namo) 222 Betspored 222 Betspor	49.88 49.88 49.88	8.81	6.61 6.01 6.01 6.01 6.01 6.01 6.01 6.01	9.9 a, 9.9	31 0.01 13 0.0	e.ei e.ei e.ei e.ei e.ei e.ei e.ei e.ei	9.81 9.81	0.01	9.81 9.81	e.ei e.ei e.ei e.ei e.ei e.ei e.ei e.ei	e.si e.si e.si e.si e.si e.si e.si e.si	6.6 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8	(Decimal('37, 3415] (Decimal('38, 5819] (Decimal('48, 7867] (Decimal('48, 7867] (Decimal('48, 7867] (Decimal('48, 7868] (Decimal('78, 7868] (Decimal('78, 7868] (Decimal('38, 842] (Decimal('38, 5477] (Decimal('38, 5477] (Decimal('38, 5477] (Decimal('38, 5478] (Decimal('38, 5478] (Decimal('48, 228] (Decimal('47, 6234] (Decimal('47, 6234] (Decimal('47, 6234] (Decimal('179, 293] Bearl Airport coordinate 8.8] (Decimal('138, 842]
Anana (Vitjazevo) 222 Apatite (Khibiny) 222 Apatite (Khibiny) 222 Apatite (Khibiny) 222 Arkhangelsk (Vask., 228 Arkhangelsk (Vask., 228 Arkhangelsk (Talayy) 222 Astrakhan (Nariam., 222 Baykiti 222 Belanal (Titov Namo) 222 Belapord 222 Belapord	49.88 49.88 49.88 9.01 9.02 9.02	0.01 0.01	e.ei e.ei e.ei e.ei e.ei e.ei e.ei e.ei	9.910.1 9.101.	3] 0.8] 3] 3] 0.8] 3] 3] 3] 3] 3] 3] 3] 3] 3] 3] 3] 3] 3]	e.el e.el e.el e.el e.el e.el e.el e.el	9.81 9.81 9.81 9.81 9.81 9.81 9.81 9.81	0.01 0.01	9.81 9.81 9.81 9.81 9.81 9.81 9.81 9.81	e.ei e.ei	e.si e.si	6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8	(Decimal('37, 3415] (Decimal('38, 5812] (Decimal('48, 7867] (Decimal('48, 7146] (Decimal('48, 7146] (Decimal('48, 7146] (Decimal('38, 842] (Decimal('38, 842] (Decimal('38, 842] (Decimal('38, 842] (Decimal('38, 842] (Decimal('38, 842] (Decimal('148, 842] (Deci
Anana (Vitjazevo) [222 Apatite (Khibiny) [222 Apatite (Khibiny) [222 Arkhangelsk (Vask] [228 Arkhangelsk (Vask] [228 Arkhangelsk (Talayy)] [222 Astrakhan (Nariam] [227 Baykit] [222 Baykit] [222 Baykit] [222 In Salah] [222 Belapord	49.88 49.88 49.88	0.01 0.01 0.01 0.01 0.01 0.01 0.01 0.01	e.ei e.ei e.ei e.ei e.ei e.ei e.ei e.ei	9.9 0. 0. 0. 0. 0. 0. 0. 0	a) e.e a) a a a a a a a a a a a a a	8.81 8.81 8.81 8.81 8.81 8.81 8.81 8.81	9.81 9.81	0.01 0.01 0.01 0.01 0.01 0.01 0.01 0.01	9.81 9.81 9.81 9.81 9.81 9.81 9.81 9.81	e.ei e.ei e.ei e.ei e.ei e.ei e.ei e.ei	e. si e.	6.0 6.0 6.0 6.0 6.0 6.0 6.0 6.0	(Decimal('37, 3415] (Decimal('38, 5812] (Decimal('48, 7867] (Decimal('48, 7146] (Decimal('48, 7146] (Decimal('48, 7146] (Decimal('38, 842] (Decimal('38, 842] (Decimal('38, 5477] (Decimal('38, 5477] (Decimal('38, 5477] (Decimal('48, 228] (Decimal('19, 248] (Decimal('19, 248] (Decimal('1138, 842] (Becimal('1138, 842] (Becimal('138, 842] (Becimal('1
Anana (Vitjazevo) 2828 Apatite (Khibiny) 2828 Apatite (Khibiny) 2828 Arkhangetsk (Vask 2828 Arkhangetsk (Vask 2828 Arkhangetsk (Talayy) 2828 Astrakhan (Narian 2828 Barnaul (Titrov Namo) 2828 Barnaul (Titrov Namo) 2828 Barnaul (Titrov Namo) 2828 Betiaprod 2829 Betiaprod 2829 Betiaprod 2829 Betiaprod 2829 Betiaprod 2828 Betiaprod 2828 Beringovskiy 2828 Inly showing top 28 rows AirportName Year January Trip 2828 1.5 Trip 2818 6.8 Trip 2819 1.35 Trip 2819 6.8 Trip 2811 6.8 Trip 2814 6.8 Trip 2815 6.8 Trip 2816 6.8 Trip 2889 6.8 Trip 2889 6.8 Trip 2888 6.8	49.88 49.88 49.88 6.81	B. 61	e.ei e.ei e.ei e.ei e.ei e.ei e.ei e.ei	9.810.1 9.810.1	a) e.el al c.el al c.e	6.81 6.81 8.81 8.81 8.81 8.81 8.81 8.81	9.81 9.81 9.81 9.81 9.81 9.81 9.81 9.81	6.61 9.02 9.03 9.04 9.05	9.81 9.81 9.81 9.81 9.81 9.81 9.81 9.81	9.01 9.01 9.01 9.01 9.01 9.01 9.01 9.01	e.si	6.8 9.9 9.8 9.8 9.8 9.8 9.8 9.8 9	(Decimal('37, 3415] (Oecimal('38, 593] (Oecimal('38, 7146] (Oecimal('48, 7146] (Oecimal('48, 7146] (Oecimal('48, 7146] (Oecimal('38, 842] (Oecimal('38, 842] (Oecimal('38, 842] (Oecimal('38, 842] (Oecimal('38, 843] (Oecimal('38, 843] (Oecimal('38, 843] (Oecimal('47, 6234] (Oecimal('47, 6234] (Oecimal('47, 6234] (Oecimal('47, 6234] (Oecimal('47, 6234] (Oecimal('138, 842] (Oecima
Anana (Vitjazevo) [222 Apatite (Khibiny) [222 Apatite (Khibiny) [222 Arkhangelsk (Vask] [228 Arkhangelsk (Vask] [228 Arkhangelsk (Talayy)] [222 Astrakhan (Nariam] [227 Baykit] [222 Baykit] [222 Baykit] [222 In Salah] [222 Belapord	49.88 49.88 49.88 69.81	B. 61	e.ei e.ei e.ei e.ei e.ei e.ei e.ei e.ei	9.8 0. 0. 0. 0. 0. 0. 0.	a) e.el al e.e	8.81 8.81 8.81 8.81 8.81 8.81 8.81 8.81	9.81 9.81 9.81 9.81 9.81 9.81 9.81 9.81	6.81 9.91	9.81 9.81 9.81 9.81 9.81 9.81 9.81 9.81	9.01 9.01 9.01 9.01 9.01 9.01 9.01 9.01	e.si	6.0 6.0 6.0 6.0 6.0 6.0 6.0 6.0	(Decimal('37, 3415] (Decimal('38, 5812] (Decimal('48, 7867] (Decimal('48, 7146] (Decimal('48, 7146] (Decimal('48, 7146] (Decimal('38, 842] (Decimal('38, 842] (Decimal('38, 5477] (Decimal('38, 5477] (Decimal('38, 5477] (Decimal('48, 228] (Decimal('19, 248] (Decimal('19, 248] (Decimal('1138, 842] (Becimal('1138, 842] (Becimal('138, 842] (Becimal('1
Anana (Vitjazevo) 222 Anana (Vitjazevo) 222 Apatite (Khibiun) 222 Ankhangelsk (Vask 222 Arkhangelsk (Vask 222 Arkhangelsk (Vask 222 Astrakhan (Narian 222 Astrakhan (Narian 222 Baykiti 222 Bayki	49.88 49.88 19.68 9.69 9.69	e.si	e.ei e.	0.8 0.1 0.102.1 0.103.1 0.1	3) 0.8 31 0.8	8.81 8.81 8.81 8.81 8.81 8.81 8.81 8.81	9.81 9.81 9.81 9.81 9.81 9.81 9.81 9.81	e.si e.si	9.81 9.81 9.81 9.81 9.81 9.81 9.81 9.81	e.ei e.ei	e. si e.	6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8	(Decimal('37, 3415] (Decimal('36, 73415] (Decimal('40, 73467] (Decimal('40, 73467] (Decimal('40, 73468] (Decimal('40, 73468] (Decimal('33, 3492] (Decimal('33, 3492] (Decimal('34, 3492] (Decimal('34, 3492] (Decimal('34, 3492] (Decimal('34, 3492] (Decimal('34, 3492] (Decimal('34, 3492] (Decimal('40, 3492] (Decimal('40, 3492] (Decimal('40, 3492] (Decimal('40, 3492] (Decimal('41, 3492] (Decimal('113, 3492.
Anapa (Vitjazevo) 2828 Apatite (Khibiny) 2828 Aphanelsis (Vasik 2828 Arkhangelsis (Vasik 2828 Arkhangelsis (Vasik 2828 Arkhangelsis (Vasik 2828 Arkhangelsis (Talayy) 2828 Barnaul (Titov Namo) 2828 Barnaul (Titov Namo) 2828 Barnaul (Titov Namo) 2828 Barnaul (Titov Namo) 2828 Betiagrodi 2828 Betiagrodi 2828 Betiagrodi 2828 Betiagrodi 2828 Beringovskiy 2828 Beringovskiy 2828 Inly showing top 28 rows AirportName Year January AirportName Year January Trip 2828 1.5 Trip 2818 6.8 Trip 2817 6.8 Trip 2815 6.8 Trip 2815 6.8 Trip 2815 6.8 Trip 2815 6.8 Trip 2816 6.8 Trip 2817 6.8 Trip 2818 6.8 Trip 2818 6.8 Trip 2818 6.8 Trip 2887 6.8 Trip 2887 6.8	49.88 49.88 9.88 9	B. 61	e.ei e.	9.8 9. 9. 9. 9. 9. 9. 9.	a) e.el 10 10 10 10 10 10 10 1	8.81 8.81 8.81 8.81 8.81 8.81 8.81 8.81	9.81 9.81 9.81 9.81 9.81 9.81 9.81 9.81	0.01 0.01 0.01 0.01 0.01 0.01 0.01 0.01	9.81 9.81 9.81 9.81 9.81 9.81 9.81 9.81	e.ei e.ei e.ei e.ei e.ei e.ei e.ei e.ei	e.si	6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8	(Decimal('37, 3415] (Oecimal('38, 5919] (Oecimal('38, 7146] (Oecimal('48, 7146] (Oecimal('48, 7146] (Oecimal('48, 7146] (Oecimal('48, 7146] (Oecimal('38, 342] (Oecimal('38, 342] (Oecimal('38, 5477] (Oecimal('38, 5477] (Oecimal('38, 5477] (Oecimal('36, 5765] (Oecimal('48, 5765] (Oecimal('48, 5765] (Oecimal('48, 542] (Oecimal('47, 6234] (Oecimal('47, 6234] (Oecimal('47, 6234] (Oecimal('47, 6234] (Oecimal('47, 6234] (Oecimal('179, 293] ear Airport coordinate 0.8 (Oecimal('138, 842] 0.8 (Oecimal('38, 842] 0.8 (
Anana (Vitjazevo) 2228 Apatite (Khibiny) 2228 Apatite (Khibiny) 2228 Arkhangelsk (Vask 2228 Arkhangelsk (Vask 2228 Astrakhan (Nariam 2229 Astrakhan (Nariam 2229 Baykiti 2228	49.88 49.88 9.81	e.si	e.ei e.	9.8 9. 9. 9. 9. 9. 9. 9.	3] 0.8] 3] 3] 0.8] 3] 3] 3] 3] 3] 3] 3] 3] 3] 3] 3] 3] 3]	8.9 8.9 8.9 8.9 8.9 8.9 8.9 8.9	9.81 9.81 9.81 9.81 9.81 9.81 9.81 9.81	6.01 6.01 9.02 9.03 9.04 9.05	9.81 9.81 9.81 9.81 9.81 9.81 9.81 9.81	8.01 0.01 0.01 0.01 0.01 0.01 0.01 0.01	e. al e. al	6.8 6.9 6.8 6.8 6.8 6.8 6.8 6.8 7 8.8 8.8 8.8 8.8 8.8 8.8 8.8	(Decimal('37, 3415] (Oecimal('35, 5819] (Oecimal('36, 7867] (Oecimal('48, 7867] (Oecimal('48, 7146] (Oecimal('48, 7146] (Oecimal('48, 7146] (Oecimal('38, 842] (Oecimal('38, 842] (Oecimal('38, 5477] (Oecimal('38, 5477] (Oecimal('38, 5477] (Oecimal('36, 5765] (Oecimal('46, 228] (Oecimal('46, 228] (Oecimal('47, 223] (Oecimal('4
Anapa (Vitjazevo) 222 Apatite (Khibiny) 202 Apatite (Khibiny) 202 Arkhangelsk (Vask 2020 Arkhangelsk (Vask 2020 Arkhangelsk (Talagy) 202 Astrakhan (Narian 2020 Barnaul (Titov Name) 202 Barnaul (Titov Name) 202 Belspored 202 Usinsk 202 Belspored	49.88 49.88 9.81 9.81 9.81 9.81 9.81 9.81 9.81 9.81	9.01 0.01 0.01 0.01 0.01 0.01 0.01 0.01	e.ei e.ei e.ei e.ei e.ei e.ei e.ei e.ei	9.819.1 9.819.1 9.819.1 9.819.1 9.819.1 9.819.1 9.819.1 9.819.1 9.819.1 9.819.1 9.819.1 9.819.1 9.819.1 9.819.1	a) e.el a a a a a a a a a a a a a	8.9 8.9 8.9 9.9 9.9 9.9 9.9 9.9	9.81 9.81 9.81 9.81 9.81 9.81 9.81 9.81	6.81	9.81 9.81 9.81 9.81 9.81 9.81 9.81 9.81	9.01 9.01 9.01 9.01 9.01 9.01 9.01 9.01	e. si	8.8 8.8 8.8 8.8 8.8 8.8 8.8 8.8 8.8 8.8	(Decimal('37, 3415] (Decimal('38, 5819] (Decimal('48, 7867] (Decimal('48, 7867] (Decimal('48, 7867] (Decimal('48, 7868] (Decimal('48, 7868] (Decimal('38, 842] (Decimal('38, 842] (Decimal('38, 5477] (Decimal('38, 5477] (Decimal('38, 5477] (Decimal('38, 5478] (Decimal('48, 5785] (Decimal('47, 6234] (Decimal('47, 6234] (Decimal('47, 6234] (Decimal('47, 6234] (Decimal('18, 842] (Decim
Anana (Vitjazevo) 222 Apatite (Khibiny) 202 Apatite (Khibiny) 202 Arkhangelsk (Vask 202 Arkhangelsk (Vask 202 Astrakhan (Narian 202 Astrakhan (Narian 202 Barnaul (Titov Namo) 202 Barnaul (Titov Namo) 202 Belspored 203 Belspored	49.88 49.88 49.88	B. si B. si	e.ei e.	9.8 10.	a) e.el a a a a a a a a a a a a a	8.9[8.0] 8.0]	9.81 9.81 9.81 9.81 9.81 9.81 9.81 9.81	6.81	9.8i 9.8i 9.8i 9.8i 9.8i 9.8i 9.8i 9.8i	9.01 0.01 0.01 0.01 0.01 0.01 0.01 0.01	e. si	6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8	[Oecimal('37, 3415] [Oecimal('38, 5919] [Oecimal('48, 7967] [Oecimal('48, 7967] [Oecimal('48, 7968] [Oecimal('48, 7968] [Oecimal('48, 7968] [Oecimal('48, 5967] [Oecimal('38, 647] [Oecimal('38, 547] [Oecimal('38, 547] [Oecimal('38, 547] [Oecimal('44, 228] [Oecimal('44, 228] [Oecimal('47, 6234] [Oecimal('47, 6234] [Oecimal('47, 6234] [Oecimal('47, 6234] [Oecimal('47, 6234] [Oecimal('179, 293] **Bell [Oecimal('138, 842] **Bell [Oe
Anapa (Vitjazevo) 2828 Apatite (Khibiny) 2828 Apatite (Khibiny) 2828 Arkhangelsk (Vask 2828 Arkhangelsk (Vask 2828 Arkhangelsk (Vask 2828 Astrakhan (Narima 2828 Astrakhan (Narima 2828 Baykiti 2828 B	49.88 49.88 9.08 9.08 9.08 9.08	B. si B. si	e.ei e.	e.e a. e.e	3] 0.8] 3] 3] 3] 3] 3] 3] 3] 3] 3] 3] 3] 3] 3]	8.9 9.0 9.0 9.0 9.0 9.0 9.0 9.0 9	9.81 9.81 9.81 9.81 9.81 9.81 9.81 9.81	6.81 6.91 6.	9.8i 9.8i 9.8i 9.8i 9.8i 9.8i 9.8i 9.8i	e.ei e.	e. al e. al	6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8	[Oecimal('37, 3415] [Oecimal('38, 5919] [Oecimal('48, 7967] [Oecimal('48, 7967] [Oecimal('48, 7968] [Oecimal('48, 7968] [Oecimal('48, 7968] [Oecimal('48, 5967] [Oecimal('38, 647] [Oecimal('38, 547] [Oecimal('38, 547] [Oecimal('38, 547] [Oecimal('44, 228] [Oecimal('44, 228] [Oecimal('47, 6234] [Oecimal('47, 6234] [Oecimal('47, 6234] [Oecimal('47, 6234] [Oecimal('47, 6234] [Oecimal('179, 293] **Bell [Oecimal('138, 842] **Bell [Oe
Anana (Vitjazevo) 2228 Apatite (Khibiny) 2228 Apatite (Khibiny) 2228 Arkhangelsk (Vask. 1228 Arkhangelsk (Vask. 1228 Arkhangelsk (Talayy) 2228 Astrakhan (Narima. 1229 Baykit1 2228 Baykit1 2228 Baykit1 2228 Baykit1 2228 Baykit1 2228 Baykit1 2228 Belapord 2238 Belapord 2338 Alaport 2338 Belapord 2338	49.88 49.88 9.08 9.08 9.08 9.08 9.08 9.08	B. 61 0.81 0.81 0.81 0.81 0.81 0.81 0.81 0.8	e.ei e.	e.e e.i e.e e.	3] 0.8] 3] 3] 0.8] 3] 3] 3] 3] 3] 3] 3] 3] 3] 3] 3] 3] 3]	8.9[9.0] 8.9[9.0] 9.9[9.0[9.0] 9.9[9.0[9.0] 9.9[9.0] 9.9[9.0[9.0] 9.9[9.0[9.0] 9.9[9.0[9.0] 9.9[9.0[9.0] 9.9[9.0[9.0] 9.9[9.0[9.0[9.0] 9.9[9.0[9.0[9.0[9.0[9.0[9.0[9.0[9.81 9.81 9.81 9.81 9.81 9.81 9.81 9.81	6.81	9.8i 9.8i 9.8i 9.8i 9.8i 9.8i 9.8i 9.8i	e.ei e.	e. si e.	6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8	[Oecimal('37, 3415] [Oecimal('35, 5819] [Oecimal('40, 7867] [Oecimal('40, 7867] [Oecimal('40, 7867] [Oecimal('40, 7868] [Oecimal('40, 7868] [Oecimal('38, 842] [Oecimal('38, 5647] [Oecimal('38, 5847] [Oecimal('38, 5847] [Oecimal('40, 228] [Oecimal('40, 228] [Oecimal('40, 228] [Oecimal('40, 228] [Oecimal('47, 224] [Oecimal('47, 224] [Oecimal('47, 224] [Oecimal('47, 224] [Oecimal('47, 224] [Oecimal('47, 224] [Oecimal('148, 842] [Oecimal('
Anapa (Vitjazevo) 222 Apatite (Khibiny) 222 Aptite (Khibiny) 223 Arkhangelsk (Vask 2028 Arkhangelsk (Vask 2028 Arkhangelsk (Vask 2028 Arkhangelsk (Talayy) 222 Astrakhan (Narian 2620 Barnaul (Titov Name) 2222 Barnaul (Titov Name) 2222 Barnaul (Titov Name) 2222 Belshi 2228 Belshi 2228 Belshi 2228 Beringovskiy 2228 Beringovskiy 2228 Beringovskiy 2228 Beringovskiy 2228 Bringovskiy 2228 Bringovskiy 2228 Inty showing top 28 rows AirportName Year January Airport	49.88 49.88 9.81	B. 61 0.81 0.81 0.81 0.81 0.81 0.81 0.81 0.8	e.ei e.	9.8 0. 0. 0. 0. 0. 0. 0.	3 0.8 3 0	8.9[8.9] 8.9[8.9[8.9] 8.9[8.9] 8.9[8.9[8.9] 8.9[8.9] 8.9[8.9[8.9] 8.9[8.9[8.9] 8.9[8.9] 8.9[8.9[8.9] 8.9[8.9[8.9] 8.9[8.9[8.9] 8.9[8.9[8.9] 8.9[8.9[8.9] 8.9[8.9[8.9] 8.9[8.9[8.9] 8.9[8.9[8.9] 8.9[8.9[8.9] 8.9[8.9[8.9] 8.9[8.9[8.9] 8.9[8.9[8.9] 8.9[8.9[8.9] 8.9[8.9[8.9[8.9] 8.9[8.9[8.9] 8.9[8.9[8.9] 8.9[8.9[8.9[8.9] 8.9[8.9[8.9[8.9] 8.9[8.9[8.9[8.9] 8.9[8.9[8.9[8.9] 8.9[8.9[8.9[8.9] 8.9[8.9[8.9[8.9] 8.9[8.9[8.9[8.9[8.9] 8.9[8.9[8.9[8.9[8.9] 8.9[8.9[8.9[8.9[8.9[8.9[8.9[8.9[9.81 9.81 9.81 9.81 9.81 9.81 9.81 9.81	6.81	9.8i 9.8i 9.8i 9.8i 9.8i 9.8i 9.8i 9.8i	e.ei e.ei e.ei e.ei e.ei e.ei e.ei e.ei	6. 81	6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8	[Oecimal('37, 3415] [Oecimal('38, 5819] [Oecimal('38, 7146] [Oecimal('48, 7146] [Oecimal('48, 7146] [Oecimal('48, 7146] [Oecimal('48, 7146] [Oecimal('38, 3842] [Oecimal('38, 3842] [Oecimal('38, 5477] [Oecimal('38, 5477] [Oecimal('38, 5477] [Oecimal('38, 5875] [Oecimal('36, 5865] [Oecimal('36, 5865] [Oecimal('47, 6234] [Oecimal('47, 6234] [Oecimal('47, 6234] [Oecimal('47, 6234] [Oecimal('17, 6234] [Oecimal('18, 842] [Oecimal('18, 8
Anapa (Vitjazevo) 2228 Apatite (Khibiny) 2228 Apatite (Khibiny) 2228 Arkhangelsk (Vask 2228 Arkhangelsk (Vask 2228 Arkhangelsk (Vask 2228 Astrakhan (Nariam 2229 Barnaul (Titov Name) 2229 Barnaul (Titov Name) 2229 Usinski 2228 Belgarod 2229 Belgarod 2	49.88 49.88 9.88 9.88 9.88 9.8	e.si	e.ei e.	e.e e.i e.e e.	3 0.8 3 0	8.9 8.9 8.9 8.9 8.9 8.9 8.9 8.9	9.81 9.81 9.81 9.81 9.81 9.81 9.81 9.81	6.81	9.8 9.8 9.8 9.8 9.8 9.8 9.8 9.8	8.01 0.01 0.01 0.01 0.01 0.01 0.01 0.01	6. 81	6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8	(Decimal('37, 3415] (Oecimal('35, 5819] (Oecimal('36, 7146] (Oecimal('40, 7146] (Oecimal('40, 7146] (Oecimal('40, 7146] (Oecimal('40, 7146] (Oecimal('40, 7146] (Oecimal('35, 842] (Oecimal('36, 5467] (Oecimal('37, 6234] (Oecimal('47, 6234] (Oecimal('47, 6234] (Oecimal('47, 6234] (Oecimal('179, 293] earl Airport coordinate 0.8 (Oecimal('138, 842) 0.8 (
Anana (Vitjazevo) 2228 Apatite (Khibiny) 2228 Apatite (Khibiny) 2228 Arkhangelsk (Yask 2228 Arkhangelsk (Yask 2228 Arkhangelsk (Yask 2228 Arkhangelsk (Yask 2228 Astrakhan (Nariana 2229 Baykati 2228 Baykati 2228 Baykati 2228 Baykati 2228 Barnaul (Titov Name) 2228 Belgaroad 2228 Belgaroad 2228 Belgaroad 2228 Belgaroad 2228 Belgaroad 2228 Belgaroad 2228 Beringovskiy 2228 Inly showing top 28 rows AirportName Year January	4,9.88 9.08 9.08 9.08 9.08 9.08	e.si	e.ei e.	e.e e.i e.e e.	3) e.el 3] e.e	8.9 8.9 8.9 8.9 8.9 8.9 8.9 8.9	9.81 9.81 9.81 9.81 9.81 9.81 9.81 9.81	6.81	9.8 9.8 9.8 9.8 9.8 9.8 9.8 9.8	B.el B.	e. si e.	6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8	(Decimal('37, 3415] (Decimal('35, 5819] (Decimal('40, 7867] (Decimal('40, 7867] (Decimal('40, 7867] (Decimal('40, 7868] (Decimal('40, 7868] (Decimal('38, 842] (Decimal('38, 5877] (Decimal('38, 5877] (Decimal('38, 5877] (Decimal('46, 228] (Decimal('46, 228] (Decimal('46, 228] (Decimal('46, 248] (Decimal('46, 248] (Decimal('46, 248] (Decimal('46, 248] (Decimal('46, 248] (Decimal('47, 248] (Decimal('47, 248] (Decimal('17, 248] (Decimal('17, 248] (Decimal('17, 248] (Decimal('17, 248] (Decimal('17, 248] (Decimal('17, 248] (Decimal('18, 842] (Decimal('18
Anapa (Vitjazevo) 222 Apatite (Khibiny) 202 Apatite (Khibiny) 202 Apatite (Khibiny) 202 Arkhangelsk (Vask 202 Arkhangelsk (Vask 202 Astrakhan (Narian 202 Astrakhan (Narian 202 Barnaul (Titov Namo) 202 Barnaul (Titov Namo) 202 Belsporod 202 Be	49.88 49.88 49.88	e.si e.	e.ei e.	e.e a. e.e	3] 0.8] 3] 3] 0.8] 3] 3] 3] 3] 3] 3] 3] 3] 3] 3] 3] 3] 3]	8.9 8.0	9.8i 9.8i 9.8i 9.8i 9.8i 9.8i 9.8i 9.8i	6.81	9.8 9.8 9.8 9.8 9.8 9.8 9.8 9.8	8.01 0.01 0.01 0.01 0.01 0.01 0.01 0.01	6.81 6.91 9.91 9.91 9.91 9.91 9.91 9.91 9.9	6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8	(Decimal('37, 3415] (Oecimal('38, 5919] (Oecimal('38, 7146] (Oecimal('48, 7146] (Oecimal('48, 7146] (Oecimal('48, 7146] (Oecimal('48, 7146] (Oecimal('48, 7146] (Oecimal('38, 842] (Oecimal('48, 85, 847] (Oecimal('38, 842] (Oecimal('38, 847] (Oecimal('48, 847] (Oecimal('138, 847] (Oecima
Anapa (Vitjazevo) 2228 Apatite (Khibiny) 2228 Apatite (Khibiny) 2228 Arkhangelsk (Vask 2228 Arkhangelsk (Vask 2228 Arkhangelsk (Vask 2228 Astrakhan (Narima 2229 Barnaul (Titov Namo) 2229 Barnaul (Titov Namo) 2229 Barnaul (Titov Namo) 2229 Belgorod 2228 Arkhane Year 3228 Artip 2228 Belgorod 2228 Artip 2228 Belgorod 2228 Belgorod 2228 Artip 2228 Belgorod 2228 Artip 2228 Artip 2228 Artip 2228 Belgorod 2228 Artip 2228 Artip 2228 Artip 2228 Artip 2228 Artip 2228 Artip 2228 Belgorod 2228 Artip 2	49.88 49.88 9.88 9.88 9.88 9.88 9.88 9.88	B. 61 0.81 0.81 0.81 0.81 0.81 0.81 0.81 0.8	e.ei e.	e.e a. e.e	3) 0.8] 3] 3] 3] 3] 3] 3] 3] 3] 3] 3] 3] 3] 3]	8.9 8.0	9.8i 9.8i 9.8i 9.8i 9.8i 9.8i 9.8i 9.8i	6.81	9.8 9.8 9.8 9.8 9.8 9.8 9.8 9.8	8.01 0.01 0.01 0.01 0.01 0.01 0.01 0.01	6. 81	6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8 6.8	[Opecimal('37, 3415] [Opecimal('38, 5819] [Opecimal('48, 7867] [Opecimal('48, 7867] [Opecimal('48, 7146] [Opecimal('49, 7146] [Opecimal('38, 5847] [Opecimal('38, 5847] [Opecimal('38, 5847] [Opecimal('38, 5847] [Opecimal('38, 5847] [Opecimal('38, 5847] [Opecimal('46, 228] [Opecimal('46, 228] [Opecimal('46, 248] [Opecimal('47, 223] [Opecimal('47, 223] [Opecimal('47, 223] [Opecimal('47, 223] [Opecimal('47, 223] [Opecimal('179, 233] Pearl Airport coordinate 8.8][Opecimal('138, 842] [Opecimal('138, 842] [Opecim

```
AirportName|Year|
|Moscow (Sheremety...|2012|
|Moscow (Sheremety...|2019|
|Moscow (Sheremety...|2018|
|Moscow (Sheremety...|2017|
|Moscow (Sheremety...|2011|
|Moscow (Sheremety...|2016|
|Moscow (Sheremety...|2013|
|Moscow (Sheremety...|2014|
|Moscow (Sheremety...|2008|
| Moscow (Domodedovo)|2013|
|Moscow (Sheremety...|2010|
| Moscow (Domodedovo)|2011|
|Moscow (Sheremety...|2007|
| Moscow (Domodedovo)|2014|
| Moscow (Domodedovo)|2012|
|Moscow (Sheremety...|2015|
| Moscow (Domodedovo)|2018|
| Moscow (Domodedovo)|2008|
| Moscow (Domodedovo)|2015|
|Moscow (Sheremety...|2009|
only showing top 20 rows
|Moscow (Sheremety...|2019| 26703.6|
|Moscow (Sheremety...|2017| 25153.5|
|Moscow (Sheremety...|2016| 17486.6|
 | Moscow (Domodedovo)|2012| 13593.3|
 | Moscow (Domodedovo)|2013|14225.51|
| Moscow (Domodedovo)|2014|14015.58|
 | Moscow (Domodedovo)|2010|11290.43|
only showing top 20 rows
|Moscow (Sheremety...|2012| 25820.2|
|Moscow (Sheremety...|2019| 28634.2|
|Moscow (Sheremety...|2018| 24709.4|
| Moscow (Domodedovo)|2013|13922.46|
| Moscow (Domodedovo)|2014|14121.12|
| Moscow (Domodedovo)|2008|11837.18|
| Moscow (Domodedovo)|2018| 9397.39|
only showing top 20 rows
 |AirportName|Year|May|June|
         ----+---+
```

Trip|2014|0.0| 0.0|

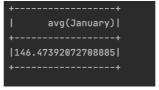


Рисунок 1 – Результат выполнения программы

Вывод: В результате выполнения лабораторной работы были приобретены навыки работы со Spark на языке программирования Java.