

# DataOps.onETL



Unit#4

Объекты подключения и манипуляции данными с использованием структуры DataFrame



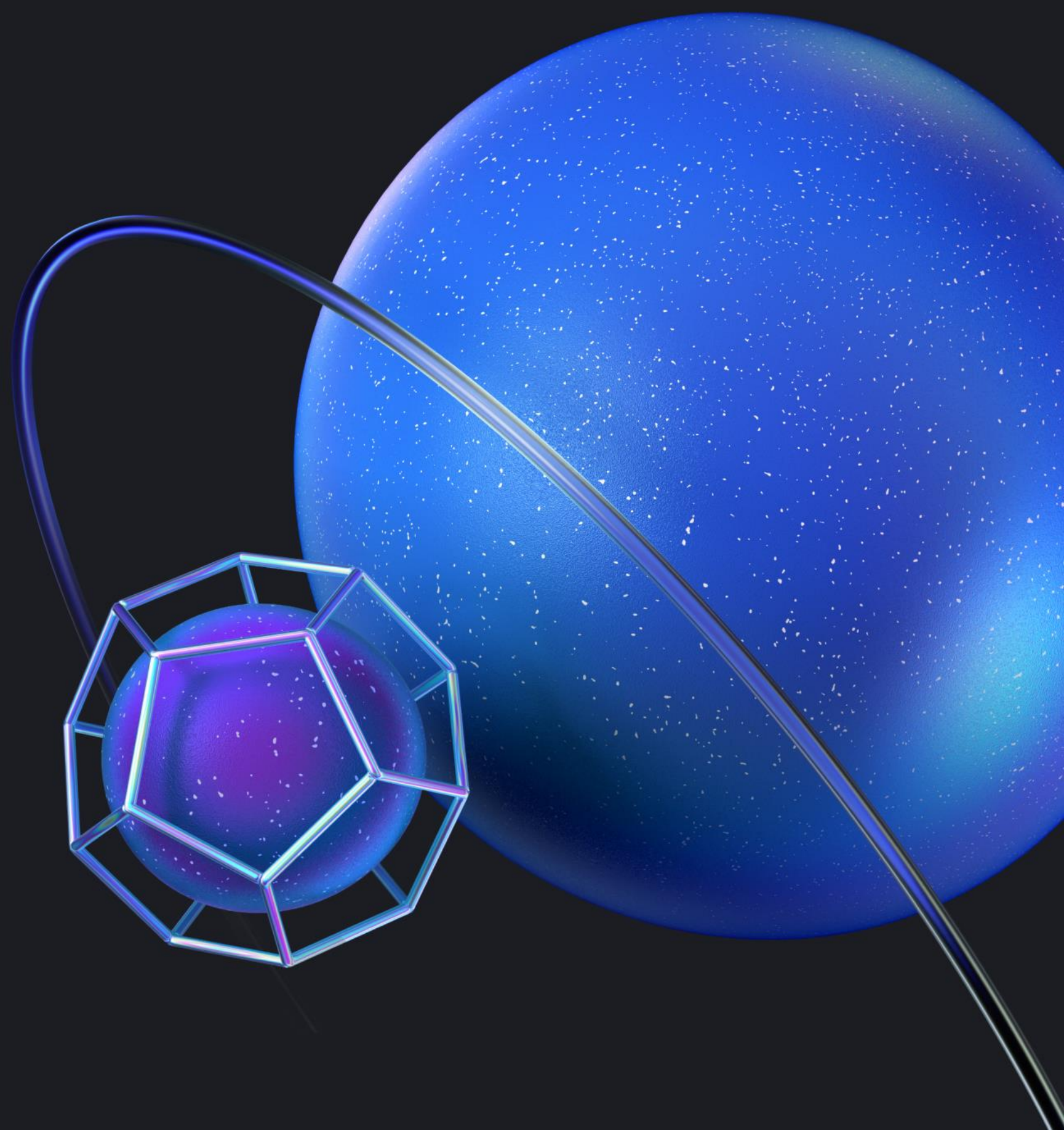
**Саттар Гюльмамедов**

РО команды ETL

МТС Тета

х

DataOps Platform



# onETL урок # 4

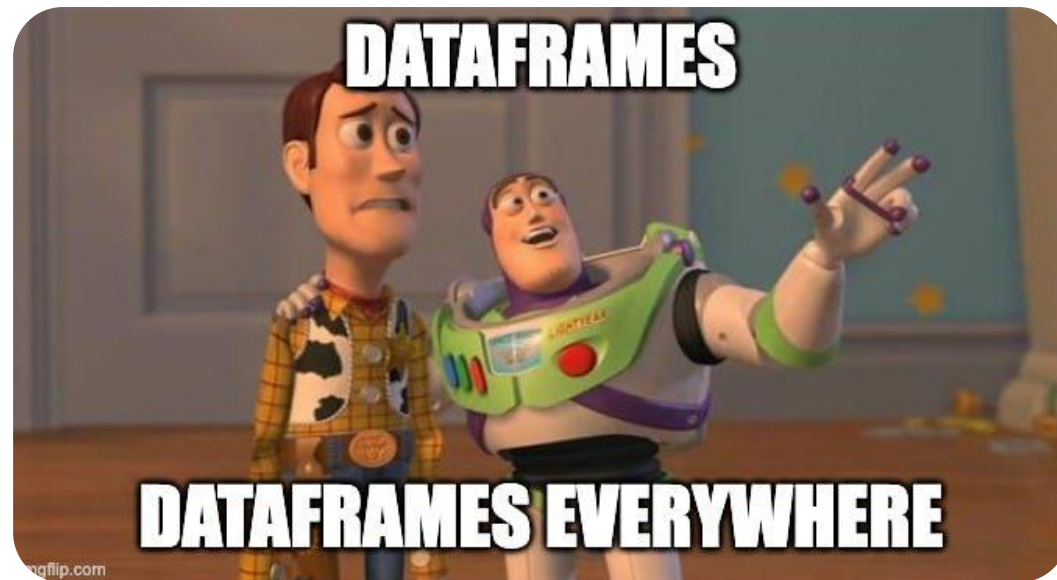


- File DataFrame Connections
- FileDF Reader
- FileDF Writer
- Где найти документацию



# DataFrame

`spark.df != pandas.df`





# FileDataFrameConnections



→ LocalFS

→ HDFS

→ S3



	A	B	C	D	E	F	G	H	I	J	K	L
2	ABW	0,810230588	0,749301039	0,69161526	0,637959162	0,590062487	0,537295706	0,494795263	0,451969659	0,134255302	-0,045044623	-0,086392283
3	AFE	2,740404834	2,780206834	2,774990291	2,802585821	2,72815853	2,655671599	2,688371428	2,691134381	2,678183716	2,607471694	2,543756525
4	AFG	4,077627729	3,466788304	3,657576065	3,121341229	2,581549399	2,866492148	2,885207973	2,908529093	3,134746908	2,851357654	2,534498317
5	AFW	2,812850935	2,761838977	2,750730568	2,723317337	2,71305851	2,706266076	2,669238668	2,633982069	2,61564614	2,573377402	2,539799442
6	AGO	3,758702522	3,735525429	3,684428967	3,617677518	3,586211007	3,550986636	3,464456984	3,395277792	3,268348405	3,166029862	3,096752671
7	ALB	-0,16515104	-0,183211385	-0,207047	-0,291205787	-0,159880412	-0,091972294	-0,246732042	-0,426007367	-0,574206959	-0,926918062	-1,21579032
8	AND	0,630034575	0,497261876	0,355274942	0,17437769	1,10060299	1,772182713	1,580147087	1,757491257	1,761891295	1,702288225	0,994607149
9	ARB	2,157196193	2,299823498	2,259225675	2,155966313	2,109697128	2,068738846	2,096194015	2,062686901	1,757899007	1,623335187	1,78833854
10	ARE	1,041344608	0,997641826	0,956397624	0,911950036	0,863868923	0,819744473	0,789450108	0,779087177	0,817694386	0,834812788	0,808075151
11	ARG	1,136905694	1,1191092	1,099461091	1,078001344	1,057181578	1,037133897	1,015808349	0,993397493	0,970053991	0,947490959	0,925835474
12	ARM	-0,49817016	-0,448296309	-0,395592881	-0,392995249	-0,444257167	-0,486625263	-0,540251082	-0,564065543	-0,533006607	-0,522963243	-0,377102034
13	ASM	-1,146298234	-1,304782022	-1,478945711	-1,639270188	-1,807230769	-1,971818747	-2,122935674	-2,304138699	-2,42124971	-2,530170926	-1,706495481
14	ATG	1,083707561	0,934326293	0,831589248	0,786935419	0,690288329	0,610956145	0,554872966	0,534443401	0,592053943	0,597151595	0,581875833
15	AUS	1,745820001	1,72115144	1,491566489	1,439216653	1,561940498	1,646827797	1,49599671	1,477490711	1,233428251	0,140895083	1,238638797
16	AUT	0,455937472	0,589387255	0,781541633	1,120992502	1,081396299	0,694621104	0,487071925	0,444673639	0,415176732	0,435671684	0,956287797
17	AZE	1,328763827	1,293447027	1,248208997	1,191209858	1,11785721	0,981261805	0,866316916	0,846646726	0,684365325	0,441197736	0,039507866
18	BDI	3,529972848	3,551107704	3,345863383	2,18870609	1,629024754	2,287301447	2,983823553	3,26424857	2,867082195	2,672478462	2,66016084
19	BEL	0,620163579	0,471340144	0,443929288	0,57944624	0,506300002	0,385228018	0,455184695	0,540461327	0,430996815	0,411601997	0,856132433
20	BEN	2,915012122	2,899920609	2,92622898	2,951249123	2,949828144	2,946321561	2,922392187	2,887074085	2,829138113	2,759705179	2,702041089
21	BFA	3,031880077	3,008421955	2,979778499	2,972345608	2,934811289	2,865655421	2,768681163	2,703875567	2,688787667	2,650375915	2,559888391
22	BGD	1,243571246	1,267157297	1,245960208	1,191061127	1,230795354	1,249724066	1,161378463	1,11317249	1,14420974	1,149318476	1,074836743
23	BGR	-0,579220596	-0,559647217	-0,568389264	-0,638069468	-0,701382098	-0,730443175	-0,722080405	-0,703905641	-0,600241518	-0,814846472	-6,187252982
24	BHR	1,055562986	2,954757119	3,845379359	3,816604913	3,429078922	3,291633803	2,072370164	0,459362585	-1,125242678	-0,96602484	0,611005553
25	BHS	1,081836998	0,934994066	0,898582547	0,912227428	0,831528122	0,765793759	0,720688948	0,657441085	0,471994445	0,352416997	0,508137884
26	BIH	-1,854258767	-1,558328902	-1,293477828	-1,317606589	-1,237305784	-1,18362668	-1,166594811	-1,166081335	-1,266771023	-1,440652596	-1,150513985
27	BLR	-0,1566176	-0,03838	0,056151568	0,132853237	0,087721088	-0,109782323	-0,213824207	-0,201786601	-0,42347525	-0,828232721	-0,804228427
28	BLZ	2,256631567	2,238167082	2,194276752	2,116320634	2,046871171	1,98926184	1,948634454	1,823016119	1,486221516	1,285629921	1,301640251
29	BMU	0,361775877	0,312791571	0,210544188	0,151869636	-1,05247098	-1,060534596	0,070427498	-0,010952131	-0,028168133	-0,202104145	-0,36450519

[20]: df.show(31, truncate=False)

	country_code	ly2012	ly2013	ly2014	ly2015	ly2016	ly2017	ly2018	ly2019	ly2020	ly2021	ly2022
1	ABW	0.810230587788097	0.749301039347625	0.691615260101675	0.63795916173479	0.590062487333077	0.537295706145068	0.494795263046989	0.451969658863314	0.134255302359952	-0.0450446230906329	-0.0863922826549927
2	AFE	2.740404834239655	2.780206833556065	2.7749902913343476	2.80258582087383	2.7281585299875815	2.6556715992301463	2.6883714280742765	2.6911343805639802	2.6781837164190136	2.6074716940955796	2.5437565245545812
3	AFG	4.07762772858795	3.46678830410982	3.65757606530364	3.12134122890873	2.58154939895659	2.86649214750073	2.88520797304849	2.90852909268743	3.13474690779163	2.85135765449931	2.53449831666733
4	AFW	2.8128509346212525	2.7618389767510365	2.7507305679848315	2.723317337417113	2.7130585103698763	2.7062660758635104	2.669238668027347	2.6339820691134435	2.6156461398579154	2.5733774016879494	2.539799442479179
5	AGO	3.75870252194126	3.73552542857565	3.68442896683492	3.61767751762095	3.58621100683242	3.5509866361995	3.46445698360985	3.39527779197097	3.26834840458456	3.16602986153025	3.09675267067326
6	ALB	-0.165151040121679	-0.183211384606402	-0.207046999760594	-0.291205786840436	-0.159880412127734	-0.0919722937442495	-0.246732042281782	-0.426007367832238	-0.574206959244013	-0.926918061639398	-1.21579032012532
7	AND	0.630034574524577	0.497261875887559	0.355274942185653	0.174377690367456	1.10060298979862	1.77218271287863	1.58014708656145	1.75749125736997	1.76189129507448	1.70228822534561	0.994607149162366
8	ARB	2.1571961933289856	2.299823498004862	2.2592256749777135	2.155966313224255	2.109697128451259	2.0687388460712413	2.0961940152885745	2.0626869010707622	1.7578990069759186	1.6233351872385242	1.7883385403946193
9	ARE	1.04134460783085	0.997641825826328	0.956397623844095	0.911950036242349	0.863868922846647	0.819744473334198	0.789450107726376	0.779087177439076	0.817694385518168	0.834812788038646	0.808075150558113
10	ARG	1.13690569378967	1.11910919999772	1.09946109101836	1.07800134449604	1.05718157762752	1.03713389686514	1.01580834908002	0.993397493293046	0.970053991438048	0.947490959333989	0.925835473662058
11	ARM	-0.498170160314836	-0.44829630947071	-0.395592881362063	-0.392995248824912	-0.444257166904663	-0.486625263038029	-0.54025108197947	-0.564065542556681	-0.533006606500067	-0.522963242641951	-0.377102034390101
12	ASM	-1.1462982343373	-1.30478202241628	-1.47894571116667	-1.63927018758043	-1.80723076890422	-1.97181874740715	-2.12293567392473	-2.30413869890523	-2.4212497099907	-2.53017092647207	-1.70649548061895
13	ATG	1.0837075607137	0.934326293295174	0.831589247616481	0.786935419380774	0.690288328652049	0.6109561448905	0.55487296546274	0.534443400511753	0.592053942582792	0.597151594527623	0.581875832853236
14	AUS	1.74582000067359	1.721151439989	1.49156648912847	1.43921665259925	1.56194049810399	1.64682779735017	1.49599671041328	1.47749071058438	1.23342825131378	0.14089508329287	1.23863879657531
15	AUT	0.45593747231891	0.589387254692989	0.781541632525841	1.12099250236958	1.08139629872802	0.694621103604983	0.487071924768999	0.444673639004876	0.415176731565556	0.435671684056384	0.956287797076634
16	AZE	1.32876382733755	1.293447027000591	1.24820899735588	1.19120985807018	1.11785721013911	0.98126180450475	0.866316916370841	0.846646726414021	0.684365325388407	0.441197735816311	0.039507866238038
17	BDI	3.52997284776032	3.55110770404353	3.3458633833966	2.18870609045616	1.62902475388378	2.28730144695888	2.98382355303175	3.26424857025887	2.86708219532744	2.67247846179476	2.66016084004174
18	BEL	0.620163579337378	0.471340143966627	0.44392928847521	0.579446241677791	0.506300002451376	0.385228018737269	0.455184695275022	0.540461327098122	0.430996815099722	0.411601996532257	0.856132433382837
19	BEN	2.91501212217311	2.89992060914312	2.92622897894691	2.95124912266827	2.94982814392549	2.94632156118807	2.92239218681089	2.88707408493347	2.82913811282206	2.75970517933377	2.70204108933956
20	BFA	3.03188007722477	3.00842195462805	2.97977849923872	2.97234560821631	2.93481128912484	2.86565542084027	2.76868116258008	2.70387556725288	2.68878766730686	2.65037591522147	2.55988839127434
21	BGD	1.24357124646827	1.26715729721878	1.24596020808779	1.19106112723824	1.23079535447054	1.24972406577957	1.16137846270492	1.1131724900167	1.14420974001656	1.14931847631384	1.0748367431585
22	BGR	-0.579220596364389	-0.559647217410859	-0.568389264052903	-0.638069468160719	-0.701382097841331	-0.730443175298289	-0.722080405222853	-0.703905641117455	-0.60024151846542	-0.814846471528947	-6.18725298204958
23	BHR	1.05556298562381	2.95475711852975	3.84537935863732	3.81660491261165	3.42907892150581	3.29163380324261	2.07237016433842	0.459362584678761	-1.12524267833833	-0.966024840061971	0.611005553402689
24	BHS	1.08183699838334	0.934994065798907	0.898582547625484	0.912227428271637	0.8315281222187361	0.765793758966215	0.720668948364001	0.657441084996389	0.471994444906908	0.352416997336031	0.508137884368941
25	BIH	-1.8542587666769	-1.55832890226336	-1.29347782804997	-1.31760658922995	-1.23730578353764	-1.1836267966933	-1.16659481128762	-1.16608133465016	-1.26677102311527	-1.44065259604193	-1.15051398459433
26	BLR	-0.156617600248906	-0.0383800002667846	0.056151567510169	0.132853236725342	0.0877210878706449	0.19897232950581	-0.213824260996771	-0.201786600970723	-0.42347525027949	-0.82823272086807	-0.80422846880244
27	BLZ	2.55663156677437	2.23816708234299	2.19427675238731	2.1163286041839	2.04687117105776	1.9878278402512	1.9486345386443	1.82301611899652	1.48622151557362	1.285629290655272	1.3016420216226742
28	BMU	0.651735877335588	0.312791570929084	0.210544438038464	0.158189636092964	1.05427897968271	1.0605345906762	0.070427470850419	0.01952133739413	0.0281681326396052	0.202104414558527	0.364505179849171
29	BOL	1.36645398784646	1.6295731991515085	1.60315110672219	1.5731494041843	1.5240782963473	1.52010808105833	1.48747419672442	1.45750439147179	1.3397922043111	1.19344861786783	1.19027478312628
30	BRA	0.9803436601223	0.868364510617205	0.85783482872841	0.84599260134206	0.8112564782781	0.79226340042272	0.793768153848091	0.766107802804815	0.656176682246802	0.52845952199416	0.455958348488331
31	BRB	0.25145634388319	0.24156434388319	0.225668487744075	0.21239263810218	0.2032394889215748	0.19228892188534	0.172887977476607	0.17575598089405	0.182992163306018	0.180841453458187	0.15945836707803
Only showing top 31 rows												

# LocalFS

```
from onetl.connection import SparkLocalFS
from pyspark.sql import SparkSession

spark = SparkSession.builder \
    .master("local") \
    .appName("spark-app-name") \
    .getOrCreate()

local_fs = SparkLocalFS(spark=spark)
local_fs.check()
```

Python



[https://onetl.readthedocs.io/en/stable/connection/file\\_df\\_connection/spark\\_local\\_fs.html](https://onetl.readthedocs.io/en/stable/connection/file_df_connection/spark_local_fs.html)

MTC Teta

x

DataOps Platform

# HDFS

```
from onetl.connection import Hive
from pyspark.sql import SparkSession

spark = (
    SparkSession.builder.appName("spark-app-name")
        .option(
            "spark.kerberos.access.hadoopFileSystems",
            "hdfs://namenode1.domain.com:8020",
        )
        .option("spark.kerberos.principal", "user")
        .option("spark.kerberos.keytab", "/path/to/keytab")
        .enableHiveSupport()
        .getOrCreate()
)

hdfs = SparkHDFS(
    host="namenode1.domain.com",
    cluster="cluster_name",
    spark=spark,
)

hdfs.check()
```

Python



[https://onetl.readthedocs.io/en/stable/connection/file\\_df\\_connection/spark\\_hdfs/index.html](https://onetl.readthedocs.io/en/stable/connection/file_df_connection/spark_hdfs/index.html)

```
from onetl.connection import SparkS3
from pyspark.sql import SparkSession

maven_packages = SparkS3.get_packages(spark_version="3.5.0")
excluded_packages = [
    "com.google.cloud.bigdataoss:gcs-connector",
    "org.apache.hadoop:hadoop-aliyun",
    "org.apache.hadoop:hadoop-azure-datalake",
    "org.apache.hadoop:hadoop-azure",
]

spark = (
    SparkSession.builder.appName("spark-app-name")
    .config("spark.jars.packages", ", ".join(maven_packages))
    .config("spark.jars.excludes", ", ".join(excluded_packages))
    .config("spark.hadoop.fs.s3a.committer.magic.enabled", "true")
    .config("spark.hadoop.fs.s3a.committer.name", "magic")
    .config(
        "spark.hadoop.mapreduce.outputcommitter.factory.scheme.s3a",
        "org.apache.hadoop.fs.s3a.commit.S3ACommitterFactory",
    )
    .config(
        "spark.sql.parquet.output.committer.class",
        "org.apache.spark.internal.io.cloud.BindingParquetOutputCommitter",
    )
    .config(
        "spark.sql.sources.commitProtocolClass",
        "org.apache.spark.internal.io.cloud.PathOutputCommitProtocol",
    )
    .getOrCreate()
)

s3 = SparkS3(
    host="domain.com",
    protocol="http",
    bucket="my-bucket",
    access_key="ACCESS_KEY",
    secret_key="SECRET_KEY",
    extra={
        "path.style.access": True,
    },
    spark=spark,
)

s3.check()
s3.close()
```

## Обычные:

- host
- [port]
- protocol (https | http)
- spark

## Специфические

- [access\_key]
- [secret\_key]
- [session\_token]
- [region]
- bucket



[https://onetl.readthedocs.io/en/stable/connection/file\\_df\\_connection/spark\\_s3/index.html](https://onetl.readthedocs.io/en/stable/connection/file_df_connection/spark_s3/index.html)

MTC Тета

x

DataOps Platform





# S3() - extra

```
extra = {  
    "path.style.access": True,  
    "committer.magic.enabled": True,  
    "committer.name": "magic",  
    "connection.timeout": 300000,  
}
```

Python

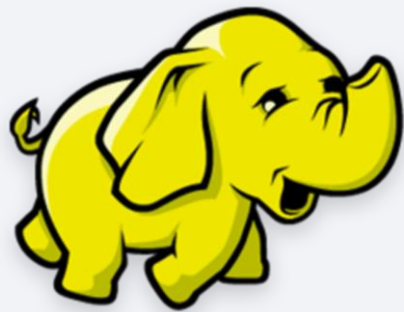


[https://hadoop.apache.org/docs/current/hadoop-aws/tools/hadoop-aws/index.html#General\\_S3A\\_Client\\_configuration](https://hadoop.apache.org/docs/current/hadoop-aws/tools/hadoop-aws/index.html#General_S3A_Client_configuration)

# S3 — Troubleshooting

```
spark.sparkContext.setLogLevel("debug")
```

Python



[https://onetl.readthedocs.io/en/stable/connection/file\\_df\\_connection/spark\\_s3/troubleshooting.html](https://onetl.readthedocs.io/en/stable/connection/file_df_connection/spark_s3/troubleshooting.html)

MTC Teta

x

DataOps Platform

# Files Formats



- CSV
- Excel
- JSON
- JSONLine
- XML
- Avro
- ORC
- Parquet

MTC Tera

x

DataOps Platform



# CSV



## Параметры:

- sep
- encoding
- quote
- escape
- header
- lineSep

## Методы:

- -



[https://onetl.readthedocs.io/en/stable/file\\_df/file\\_formats/csv.html](https://onetl.readthedocs.io/en/stable/file_df/file_formats/csv.html)



<https://spark.apache.org/docs/latest/sql-data-sources-csv.html>

MTC Teta

x

DataOps Platform

# Excel

## Параметры:

- header
- inferSchema

## Методы:

- get\_packages



[https://onel.readthedocs.io/en/stable/file\\_df/file\\_formats/excel.html](https://onel.readthedocs.io/en/stable/file_df/file_formats/excel.html)



<https://github.com/crealytics/spark-excel>

MTC Тета

×

DataOps Platform

# JSON

## Параметры:

- multiLine
- encoding
- lineSep

## Методы:

- -



[https://oneth.readthedocs.io/en/stable/file\\_df/file\\_formats/json.html](https://oneth.readthedocs.io/en/stable/file_df/file_formats/json.html)



<https://spark.apache.org/docs/latest/sql-data-sources-json.html>



## Параметры:

- multiLine
- encoding
- lineSep

## Методы:

- -



[https://onetl.readthedocs.io/en/stable/file\\_df/file\\_formats/jsonline.html](https://onetl.readthedocs.io/en/stable/file_df/file_formats/jsonline.html)



<https://spark.apache.org/docs/latest/sql-data-sources-json.html>

# XML

## Параметры:

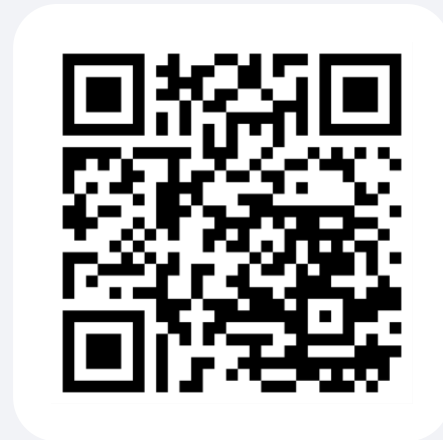
→ rowTag

## Методы:

→ get\_packages



[https://onetl.readthedocs.io/en/stable/file\\_df/file\\_formats/xml.html](https://onetl.readthedocs.io/en/stable/file_df/file_formats/xml.html)



<https://github.com/databricks/spark-xml>

# Avro

## Параметры:

- avroSchema
- avroSchemaUrl

## Методы:

- get\_packages



[https://onetl.readthedocs.io/en/stable/file\\_df/file\\_formats/avro.html](https://onetl.readthedocs.io/en/stable/file_df/file_formats/avro.html)



<https://spark.apache.org/docs/latest/sql-data-sources-avro.html>

MTC Тета

x

DataOps Platform



# ORC

Параметры:

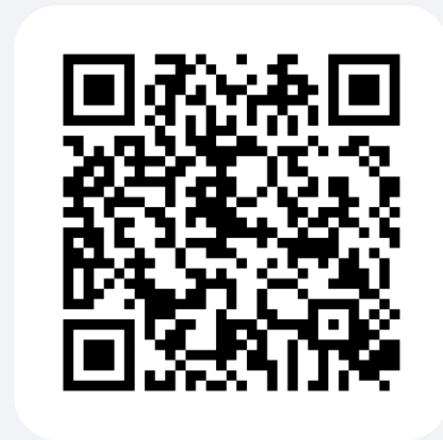
→ -

Методы:

→ -



[https://onetl.readthedocs.io/en/stable/file\\_df/file\\_formats/orc.html](https://onetl.readthedocs.io/en/stable/file_df/file_formats/orc.html)



<https://spark.apache.org/docs/latest/sql-data-sources-orc.html>

MTC Тета

×

DataOps Platform

# Parquet



Параметры:

→ -

Методы:

→ -



[https://onetl.readthedocs.io/en/stable/file\\_df/file\\_formats/parquet.html](https://onetl.readthedocs.io/en/stable/file_df/file_formats/parquet.html)



<https://spark.apache.org/docs/latest/sql-data-sources-parquet.html>

MTC Teta

×

DataOps Platform

# FileDF Reader



→ FileDFReader()

→ run()

MTC Тета

×

DataOps Platform



# FileDFReader()

- connection
- source\_path
- format



MTC Teta

x

DataOps Platform

# df\_schema

Python

```
from pyspark.sql.types import (  
    DoubleType,  
    IntegerType,  
    StringType,  
    StructField,  
    StructType,  
    TimestampType,  
)  
  
df_schema = StructType(  
    [  
        StructField("_id", IntegerType()),  
        StructField("text_string", StringType()),  
        StructField("hwm_int", IntegerType()),  
        StructField("hwm_datetime",  
TimestampType()),  
        StructField("float_value", DoubleType()),  
    ],  
)
```

# options



[https://onetl.readthedocs.io/en/stable/file\\_df/file\\_df\\_reader/options.html](https://onetl.readthedocs.io/en/stable/file_df/file_df_reader/options.html)



<https://spark.apache.org/docs/latest/sql-data-sources-generic-options.html>

```
from onetl.connection import SparkLocalFS
from onetl.file import FileDFReader
from onetl.file.format import CSV
```

```
csv = CSV(delimiter=",")
local_fs = SparkLocalFS(spark=spark)
```

```
reader = FileDFReader(
    connection=local_fs,
    format=csv,
    source_path="/path",
)
df = reader.run()
```

Python

# FileDF Writer



→ FileDFWriter()

→ run(df)

MTC Тета

×

DataOps Platform



# FileDFWriter()

- connection
- target\_path
- format



MTC Teta

x

DataOps Platform

# options



[https://onetl.readthedocs.io/en/stable/file\\_df/file\\_df\\_writer/options.html](https://onetl.readthedocs.io/en/stable/file_df/file_df_writer/options.html)



<https://spark.apache.org/docs/latest/sql-data-sources-load-save-functions.html>

MTC Teta

x

DataOps Platform

# run(df)

```
from onetl.connection import SparkLocalFS
from onetl.file import FileDFWriter
from onetl.file.format import CSV

csv = CSV(delimiter=",")
local_fs = SparkLocalFS(spark=spark)

writer = FileDFWriter(
    connection=local_fs,
    format=csv,
    target_path="/path/to/directory",
    options=FileDFWriter.Options(if_exists="replace_entire_directory"),
)

writer.run(df)
```

Python

# onETL урок # 4



- File DataFrame Connections
- FileDF Reader
- FileDF Writer
- Где найти документацию

# Спасибо!

М Т  
С



[onetools@mts.ru](mailto:onetools@mts.ru)  
<https://t.me/c/1511728757/5>

МТС Тета

x

DataOps Platform

