

PYSPARK ASSIGNMENT

Group 13 - 20BDA47

ETL using Pyspark

Extract
<ul style="list-style-type: none"><li>Read input from different datasource - json, csv, parquet</li><li>from data structures - dataframe, array, others</li></ul>
Transform
<ul style="list-style-type: none"><li>Data preprocessing/cleaning</li><li>RDD-&gt;transformation - map, filter</li><li>Dataframe function - groupby, filter and so on</li></ul>
Load
<ul style="list-style-type: none"><li>Store as file</li><li>store to DB</li></ul>

```
In [3]: from random import random
import os
from pyspark.sql import SparkSession
import pandas as pd
```

Load the data : Reading data as pandas dataframe

```
In [2]: datapd = pd.read_csv("D:/ADATA/netflix_titles.csv")

In [3]: datapd

Out[3]:
```

	show_id	type	title	director	cast	country	date_added	release_year	rating	duration	listed_in	description
0	s1	TV Show	Blood & Water	NaN	Ama Qamata, Khosi Ngema, Gail Mabasa, Thabani...	South Africa	September 24, 2021	2020	TV-MA	2 Seasons	International TV Shows, TV Dramas, TV Mysteries	After crossing paths at a party, a Cape Town L...
2	s3	TV Show	Ganglands	Julien Leclercq	Samir Bouajila, Tracy Gotsas, Samuel Jouy, Nabil...	NaN	September 24, 2021	2021	TV-MA	1 Season	Crime TV Shows, International TV Shows, TV Acti...	To protect his family from a powerful drug lor...
3	s4	TV Show	Jailbirds New Orleans	NaN	NaN	NaN	September 24, 2021	2021	TV-MA	1 Season	Docuseries, Reality TV	Feuds, filtrations and toilet talk go down amo...
4	s5	TV Show	Kota Factory	NaN	Mayur More, Jitendra Kumar, Ranjan Raj, Alam K...	India	September 24, 2021	2021	TV-MA	2 Seasons	International TV Shows, Romantic TV Shows, TV ...	In a city of coaching centers known to train L...
...	...	...	...	...	...	...	...	...	...	...	...	...
8802	s803	Movie	Zodiac	David Fincher	Mark Ruffalo, Jake Gyllenhaal, Robert Downey J...	United States	November 20, 2009	2007	R	158 min	Cut Movies, Dramas, Thrillers	A political cartoonist, a crime reporter and a...
8803	s804	TV Show	Zombie Dumb	NaN	NaN	NaN	July 1, 2019	2019	TV-Y7	2 Seasons	Kids' TV, Korean TV Shows, TV Comedies	While living alone in a spooky town, a young g...
8804	s805	Movie	Zombieland	Ruben Fleischer	Jesse Eisenberg, Woody Harrisson, Emma Stone...	United States	November 1, 2019	2009	R	88 min	Comedies, Horror Movies	Looking to survive in a world taken over by z...
8805	s806	Movie	Zooan	Peter Hewitt	Tim Allen, Courtney Cox, Chevy Chase, Kate Ma...	United States	January 11, 2020	2006	PG	88 min	Children & Family Movies, Comedies	Dragged from civilian life, a former superhero...
8806	s807	Movie	Zumaan	Mozes Singh	Vicky Kaushal, Sarah-Jane Dias, Raaghav Chahal...	India	March 2, 2019	2015	TV-14	111 min	Dramas, International Movies, Music & Musicals	A scrappy but poor boy worrms his way into a...

8807 rows x 12 columns

Creating spark dataframe

```
In [4]: spark_session = SparkSession.builder.master("local").\
        .appName("SparkApplication")\
        .config("spark.driver.bindAddress", "localhost")\
        .config("spark.ui.port", "4041")\
        .getOrCreate()
sc = spark_session.sparkContext

In [5]: data = spark_session.read.csv("D:/ADATA/netflix_titles.csv", sep=',', inferSchema=True, header=True, nullValue='')

In [6]: print(data)

DataFrame[show_id: string, type: string, title: string, director: string, cast: string, country: string, date_added: string, release_year: string, rating: string, duration: string, listed_in: string, description: string]

In [7]: print(data.show())

+-----+
|show_id|type|title|director|cast|country|date_added|release_year|rating|duration|listed_in|
+-----+
|s1|Movie|Dick Johnson Is Dead|Kirsten Johnson|null|United States|September 25, 2021|2020|PG-13|90 min|Documentaries|As her fa
ther nears the end of his life, filmm...
|s2|TV Show|Blood & Water|null|Ama Qamata, Khosi...|South Africa|September 24, 2021|2021|TV-MA|2 Seasons|International TV Shows, TV Dramas, TV Mysteries|After crossing paths at a party, a Cape Town L...
|s3|TV Show|Ganglands|Julien Leclercq|Samir Bouajila, Tracy Gotsas, Samuel Jouy, Nabil...|NaN|September 24, 2021|2021|TV-MA|1 Season|Crime TV Shows, I...|To protec
t his fa...
|s4|TV Show|Jailbirds New Orleans|null|null|NaN|September 24, 2021|2021|TV-MA|1 Season|Docuseries, Reality TV|Feuds, f
iltrations...
|s5|TV Show|Kota Factory|null|Mayur More, Jiten...|India|September 24, 2021|2021|TV-MA|2 Seasons|International TV Shows, Romantic TV Shows, TV Hor...|In a city
of coac...
|s6|TV Show|Midnight Mass|Mike Flanagan|Kate Siegel, Zach...|null|September 24, 2021|2021|TV-MA|1 Season|TV Dramas, TV Hor...|The arriv
al of a ...
|s7|Movie|My Little Pony: A...|Robert Cullen, Jo...|Vanessa Hudgens, ...|null|September 24, 2021|2021|PG|91 min|Children & Family...|Equestri
a's divid...
|s8|Movie|Sankofa|Haile Gerima|Kofi Ghanaba, Oya...|United States, Gh...|September 24, 2021|1993|TV-MA|125 min|Dramas, Independe...|On a phot
o shoot ...
|s9|TV Show|The Great British...|Andy Devonshire|Mel Giedroyc, Sue...|United Kingdom|September 24, 2021|2021|TV-MA|9 Seasons|British TV Shows, ...|A talente
d batch ...
|s10|Movie|The Starling|Theodore Melfi|Melissa McCarthy,...|United States|September 24, 2021|2021|PG-13|184 min|Comedies, Dramas|A woman a
djusting...
|s11|TV Show|Vendetta: Truth, ...|null|null|null|September 24, 2021|2021|TV-MA|1 Season|Crime TV Shows, D...|Sicily b
oasts a ...
|s12|TV Show|Bangkok Breaking|Kongkiat Komesiri|Sukollawat Kanaro...|null|September 23, 2021|2021|TV-MA|1 Season|Crime TV Shows, I...|Strugglin
g to ear...
|s13|Movie|Je Suis Karl|Christian Schwochow|Luna Wedler, Jann...|Germany, Czech Re...|September 23, 2021|2021|TV-MA|127 min|Dramas, Internati...|After mos
t of her...
|s14|Movie|Confessions of an...|Bruno Garotti|Klara Castanho, L...|null|September 22, 2021|2021|TV-PG|91 min|Children & Family...|When the
clever b...
|s15|TV Show|Crime Stories: In...|null|null|null|September 22, 2021|2021|TV-MA|1 Season|British TV Shows, ...|Cameras f
ollowing...
|s16|TV Show|Dear White People|null|Logan Browning, B...|United States|September 22, 2021|2021|TV-MA|4 Seasons|TV Comedies, TV D...|Students
of colo...
|s17|Movie|Europe's Most Dan...|Pedro de Echave G...|null|null|September 22, 2021|2020|TV-MA|67 min|Documentaries, In...|Declassif
ied docu...
|s18|TV Show|Falsa identidad|null|Luis Ernesto Fran...|Mexico|September 22, 2021|2020|TV-MA|2 Seasons|Crime TV Shows, S...|Strangers
Diego a...
|s19|Movie|Intrusion|Adam Salky|Freida Pinto, Log...|null|September 22, 2021|2021|TV-14|94 min|Thrillers|After a d
eadly ho...
|s20|TV Show|Jaguar|null|Blanca Suárez, Iv...|null|September 22, 2021|2021|TV-MA|1 Season|International TV ...|In the 19
60s, a H...
+-----+
only showing top 28 rows
None
```

```
In [8]: data.printSchema()

root
 |-- show_id: string (nullable = true)
 |-- type: string (nullable = true)
 |-- title: string (nullable = true)
 |-- director: string (nullable = true)
 |-- cast: string (nullable = true)
 |-- country: string (nullable = true)
 |-- date_added: string (nullable = true)
 |-- release_year: string (nullable = true)
 |-- rating: string (nullable = true)
 |-- duration: string (nullable = true)
 |-- listed_in: string (nullable = true)
 |-- description: string (nullable = true)
```

```
In [9]: print(data.columns)

['show_id', 'type', 'title', 'director', 'cast', 'country', 'date_added', 'release_year', 'rating', 'duration', 'listed_in', 'description']
```

```
In [10]: data.count()

8809
```

```
Out[10]: 8809

In [11]: header = data.first()

In [12]: print(header)

Row(show_id='s1', type='Movie', title='Dick Johnson Is Dead', director='Kirsten Johnson', cast=None, country='United States', date_added='September 25, 2021', release_year='2020', rating='PG-13', duration='90 min', listed_in='Documentaries', description='As her father nears the end of his life, filmmaker Kirsten Johnson stages his death in inventive and comical ways to help them both face the inevitable.')
```

```
In [13]: print(data.dtypes)

[('show_id', 'string'), ('type', 'string'), ('title', 'string'), ('director', 'string'), ('cast', 'string'), ('country', 'string'), ('date_added', 'string'), ('release_year', 'string'), ('rating', 'string'), ('duration', 'string'), ('listed_in', 'string'), ('description', 'string')]
```

```
In [14]: print(data.head())

Row(show_id='s1', type='Movie', title='Dick Johnson Is Dead', director='Kirsten Johnson', cast=None, country='United States', date_added='September 25, 2021', release_year='2020', rating='PG-13', duration='90 min', listed_in='Documentaries', description='As her father nears the end of his life, filmmaker Kirsten Johnson stages his death in inventive and comical ways to help them both face the inevitable.')
```

```
In [15]: data.describe().show()

+-----+
|summary|duration|show_id|listed_in|description|title|director|cast|country|date_added|release_year|
+-----+
|count|8809|8809|8808|8806|8807|6173|7983|7977|8796|8807| |
|mean|8804|1994.0|null|1124.7692307692307|8806|6173|7983|7977|8796|8807|
|stddev|2016.8|1994.0|null|null|4.66666666666667|8806|6173|7983|7977|8796|8807|
|s1|TV Show|Ganglands|Julien Leclercq|Samir Bouajila, Tracy Gotsas, Samuel Jouy, Nabil...|NaN|September 24, 2021|2021|TV-MA|1 Season|Crime TV Shows, I...|To protec
t his fa...
|s4|TV Show|Jailbirds New Orleans|null|null|NaN|September 24, 2021|2021|TV-MA|1 Season|Docuseries, Reality TV|Feuds, f
iltrations...
|s5|TV Show|Kota Factory|null|Mayur More, Jiten...|India|September 24, 2021|2021|TV-MA|2 Seasons|International TV Shows, Romantic TV Shows, TV Hor...|In a city
of coac...
|s6|TV Show|Midnight Mass|Mike Flanagan|Kate Siegel, Zach...|null|September 24, 2021|2021|TV-MA|1 Season|TV Dramas, TV Hor...|The arriv
al of a ...
|s7|Movie|My Little Pony: A...|Robert Cullen, Jo...|Vanessa Hudgens, ...|null|September 24, 2021|2021|PG|91 min|Children & Family...|Equestri
a's divid...
|s8|Movie|Sankofa|Haile Gerima|Kofi Ghanaba, Oya...|United States, Gh...|September 24, 2021|1993|TV-MA|125 min|Dramas, Independe...|On a phot
o shoot ...
|s9|TV Show|The Great British...|Andy Devonshire|Mel Giedroyc, Sue...|United Kingdom|September 24, 2021|2021|TV-MA|9 Seasons|British TV Shows, ...|A talente
d batch ...
|s10|Movie|The Starling|Theodore Melfi|Melissa McCarthy,...|United States|September 24, 2021|2021|PG-13|184 min|Comedies, Dramas|A woman a
djusting...
|s11|TV Show|Vendetta: Truth, ...|null|null|null|September 24, 2021|2021|TV-MA|1 Season|Crime TV Shows, D...|Sicily b
oasts a ...
|s12|TV Show|Bangkok Breaking|Kongkiat Komesiri|Sukollawat Kanaro...|null|September 23, 2021|2021|TV-MA|1 Season|Crime TV Shows, I...|Strugglin
g to ear...
|s13|Movie|Je Suis Karl|Christian Schwochow|Luna Wedler, Jann...|Germany, Czech Re...|September 23, 2021|2021|TV-MA|127 min|Dramas, Internati...|After mos
t of her...
|s14|Movie|Confessions of an...|Bruno Garotti|Klara Castanho, L...|null|September 22, 2021|2021|TV-PG|91 min|Children & Family...|When the
clever b...
|s15|TV Show|Crime Stories: In...|null|null|null|September 22, 2021|2021|TV-MA|1 Season|British TV Shows, ...|Cameras f
ollowing...
|s16|TV Show|Dear White People|null|Logan Browning, B...|United States|September 22, 2021|2021|TV-MA|4 Seasons|TV Comedies, TV D...|Students
of colo...
|s17|Movie|Europe's Most Dan...|Pedro de Echave G...|null|null|September 22, 2021|2020|TV-MA|67 min|Documentaries, In...|Declassif
ied docu...
|s18|TV Show|Falsa identidad|null|Luis Ernesto Fran...|Mexico|September 22, 2021|2020|TV-MA|2 Seasons|Crime TV Shows, S...|Strangers
Diego a...
|s19|Movie|Intrusion|Adam Salky|Freida Pinto, Log...|null|September 22, 2021|2021|TV-14|94 min|Thrillers|After a d
eadly ho...
|s20|TV Show|Jaguar|null|Blanca Suárez, Iv...|null|September 22, 2021|2021|TV-MA|1 Season|International TV ...|In the 19
60s, a H...
+-----+
only showing top 28 rows
```

```
In [16]: print(data.distinct().count())

8809
```

```
In [17]: from pyspark.sql import functions as F
data.filter(F.col("show_id").cast("int").isNull()).count()

8809
```

```
Out[17]: 8809

In [18]: data.describe().show()

+-----+
|summary|duration|show_id|type|listed_in|description|title|director|cast|country|date_added|release_year|
+-----+
|count|8809|8809|8808|8806|8807|6173|7983|7977|8796|8807| |
|mean|8804|1994.0|null|1124.7692307692307|8806|6173|7983|7977|8796|8807|
|stddev|2016.8|1994.0|null|null|4.66666666666667|8806|6173|7983|7977|8796|8807|
|s1|TV Show|Ganglands|Julien Leclercq|Samir Bouajila, Tracy Gotsas, Samuel Jouy, Nabil...|NaN|September 24, 2021|2021|TV-MA|1 Season|Crime TV Shows, I...|To protec
t his fa...
|s4|TV Show|Jailbirds New Orleans|null|null|NaN|September 24, 2021|2021|TV-MA|1 Season|Docuseries, Reality TV|Feuds, f
iltrations...
|s5|TV Show|Kota Factory|null|Mayur More, Jiten...|India|September 24, 2021|2021|TV-MA|2 Seasons|International TV Shows, Romantic TV Shows, TV Hor...|In a city
of coac...
|s6|TV Show|Midnight Mass|Mike Flanagan|Kate Siegel, Zach...|null|September 24, 2021|2021|TV-MA|1 Season|TV Dramas, TV Hor...|The arriv
al of a ...
|s7|Movie|My Little Pony: A...|Robert Cullen, Jo...|Vanessa Hudgens, ...|null|September 24, 2021|2021|PG|91 min|Children & Family...|Equestri
a's divid...
|s8|Movie|Sankofa|Haile Gerima|Kofi Ghanaba, Oya...|United States, Gh...|September 24, 2021|1993|TV-MA|125 min|Dramas, Independe...|On a phot
o shoot ...
|s9|TV Show|The Great British...|Andy Devonshire|Mel Giedroyc, Sue...|United Kingdom|September 24, 2021|2021|TV-MA|9 Seasons|British TV Shows, ...|A talente
d batch ...
|s10|Movie|The Starling|Theodore Melfi|Melissa McCarthy,...|United States|September 24, 2021|2021|PG-13|184 min|Comedies, Dramas|A woman a
djusting...
|s11|TV Show|Vendetta: Truth, ...|null|null|null|September 24, 2021|2021|TV-MA|1 Season|Crime TV Shows, D...|Sicily b
oasts a ...
|s12|TV Show|Bangkok Breaking|Kongkiat Komesiri|Sukollawat Kanaro...|null|September 23, 2021|2021|TV-MA|1 Season|Crime TV Shows, I...|Strugglin
g to ear...
|s13|Movie|Je Suis Karl|Christian Schwochow|Luna Wedler, Jann...|Germany, Czech Re...|September 23, 2021|2021|TV-MA|127 min|Dramas, Internati...|After mos
t of her...
|s14|Movie|Confessions of an...|Bruno Garotti|Klara Castanho, L...|null|September 22, 2021|2021|TV-PG|91 min|Children & Family...|When the
clever b...
|s15|TV Show|Crime Stories: In...|null|null|null|September 22, 2021|2021|TV-MA|1 Season|British TV Shows, ...|Cameras f
ollowing...
|s16|TV Show|Dear White People|null|Logan Browning, B...|United States|September 22, 2021|2021|TV-MA|4 Seasons|TV Comedies, TV D...|Students
of colo...
|s17|Movie|Europe's Most Dan...|Pedro de Echave G...|null|null|September 22, 2021|2020|TV-MA|67 min|Documentaries, In...|Declassif
ied docu...
|s18|TV Show|Falsa identidad|null|Luis Ernesto Fran...|Mexico|September 22, 2021|2020|TV-MA|2 Seasons|Crime TV Shows, S...|Strangers
Diego a...
|s19|Movie|Intrusion|Adam Salky|Freida Pinto, Log...|null|September 22, 2021|2021|TV-14|94 min|Thrillers|After a d
eadly ho...
|s20|TV Show|Jaguar|null|Blanca Suárez, Iv...|null|September 22, 2021|2021|TV-MA|1 Season|International TV ...|In the 19
60s, a H...
+-----+
only showing top 28 rows
```

Exploratory data analysis using spark df : Unique showld count

```
In [24]: data.select('show_id').distinct().count()

8809
```

```
In [25]: from pyspark.sql.functions import countDistinct
gr = data.agg(countDistinct("show_id"))
gr.show()

+-----+
|count(show_id)|
+-----+
|8809|
+-----+
```

GroupBy type and count of showld

```
In [26]: gr = data.groupBy("type").agg(countDistinct("show_id"))
gr.show()

+-----+
|type|count(show_id)|
+-----+
|TV Show|2876|
|Movie|6131|
|William Wyler|1|
|null|1|
+-----+
```

GroupBy type, release\_year and count of showld

```
In [27]: gr = data.groupBy("type", "release_year").agg(countDistinct("show_id"))
gr.show()

+-----+
|type|release_year|count(show_id)|
+-----+
|Movie|June 12, 2021|1| |
|Movie|1983|1|
|TV Show|1981|1|
|TV Show|1971|5|
|TV Show|1988|2|
|TV Show|1972|1|
|TV Show|Nse Iype-Elimi|1|
|Movie|1956|2|
|Movie|Charles Rocket|1|
|Movie|1997|33|
|Movie|2015|397|
|Movie|1969|2|
|s1|TV Show|2010|153|
|Movie|1993|24|
|TV Show|1977|6|
|TV Show|2020|436|
|TV Show|1997|4|
|Movie|2016|657|
|Movie|1992|29|
|TV Show|1945|1|
+-----+
only showing top 28 rows
```

GroupBy type, release\_year and count of showld

```
In [28]: from pyspark.sql.functions import when
df = data.withColumn("duration", when(data.duration == "90 min", "90") \
    .when(data.duration == "2 Seasons", "2") \
    .otherwise(data.duration))
df.show()

+-----+
|show_id|type|title|director|cast|country|date_added|release_year|rating|duration|listed_in|
+-----+
|s1|Movie|Dick Johnson Is Dead|Kirsten Johnson|null|United States|September 25, 2021|2020|PG-13|90|Documentaries|As her fat
her nears...
|s2|TV Show|Blood & Water|null|Ama Qamata, Khosi...|South Africa|September 24, 2021|2021|TV-MA|2|International TV ...|After cros
sing pa...
|s3|TV Show|Ganglands|Julien Leclercq|Samir Bouajila, Tr...|NaN|September 24, 2021|2021|TV-MA|1|Crime TV Shows, I...|To protec
t his fa...
|s4|TV Show|Jailbirds New Orleans|null|null|NaN|September 24, 2021|2021|TV-MA|1|Docuseries, Reality TV|Feuds, fil
trations...
|s5|TV Show|Kota Factory|null|Mayur More, Jiten...|India|September 24, 2021|2021|TV-MA|2|International TV ...|In a city
of coac...
|s6|TV Show|Midnight Mass|Mike Flanagan|Kate Siegel, Zach...|null|September 24, 2021|2021|TV-MA|1|TV Dramas, TV Hor...|The arriva
l of a ...
|s7|Movie|My Little Pony: A...|Robert Cullen, Jo...|Vanessa Hudgens, ...|null|September 24, 2021|2021|PG|91|Children & Family...|Equestri
a's divid...
|s8|Movie|Sankofa|Haile Gerima|Kofi Ghanaba, Oya...|United States, Gh...|September 24, 2021|1993|TV-MA|125 min|Dramas, Independe...|On a phot
o shoot ...
|s9|TV Show|The Great British...|Andy Devonshire|Mel Giedroyc, Sue...|United Kingdom|September 24, 2021|2021|TV-14|9|British TV Shows, I...|A talente
d batch ...
|s10|Movie|The Starling|Theodore Melfi|Melissa McCarthy,...|United States|September 24, 2021|2021|PG-13|184|Comedies, Dramas|A woman a
djusting...
|s11|TV Show|Vendetta: Truth, ...|null|null|null|September 24, 2021|2021|TV-MA|1|Crime TV Shows, D...|Sicily bo
asts a ...
|s12|TV Show|Bangkok Breaking|Kongkiat Komesiri|Sukollawat Kanaro...|null|September 24, 2021|2021|TV-MA|1|Crime TV Shows, I...|Strugglin
g to ear...
|s13|Movie|Je Suis Karl|Christian Schwochow|Luna Wedler, Jann...|Germany, Czech Re...|September 23, 2021|2021|TV-MA|127 min|Dramas, Internati...|After mos
t of her...
|s14|Movie|Confessions of an...|Bruno Garotti|Klara Castanho, L...|null|September 22, 2021|2021|TV-PG|91|Children & Family...|When the c
lever b...
|s15|TV Show|Crime Stories: In...|null|null|null|September 22, 2021|2021|TV-MA|1|British TV Shows, ...|Cameras fo
llowing...
|s16|TV Show|Dear White People|null|Logan Browning, B...|United States|September 22, 2021|2021|TV-MA|4|TV Comedies, TV D...|Students
of colo...
|s17|Movie|Europe's Most Dan...|Pedro de Echave G...|null|null|September 22, 2021|2020|TV-MA|67 min|Documentaries, In...|Declassifi
ed docu...
|s18|TV Show|Falsa identidad|null|Luis Ernesto Fran...|Mexico|September 22, 2021|2020|TV-MA|2|Crime TV Shows, S...|Strangers
Diego a...
|s19|Movie|Intrusion|Adam Salky|Freida Pinto, Log...|null|September 22, 2021|2021|TV-14|94|Thrillers|After a d
eadly ho...
|s20|TV Show|Jaguar|null|Blanca Suárez, Iv...|null|September 22, 2021|2021|TV-MA|1|International TV ...|In the 19
60s, a H...
+-----+
only showing top 28 rows
```

Groupby type and avg durations

```
In [31]: from pyspark.sql.functions import col, sum, avg, max
df_groupby("type") \
    .agg(avg("duration").alias("avg durations")) \
    .show(truncate=False)

+-----+
|type|avg durations|
+-----+
|TV Show|[1.7654320987854322]|
|Movie|[199.3898768862828]|
|William Wyler|1|
|null|1|
+-----+
```

Saving analysis report- save as tables - partition by type

```
In [34]: df.write.option("header", True) \
        .partitionBy("type") \
        .mode("overwrite") \
        .csv("D:/ADATA/netflix")
```