# [ ETL processes using PySpark ] # Quick Summary

## 1. Environment Setup and SparkSession Creation

- **Install PySpark**: `pip install pyspark`
- **Start a SparkSession**: `from pyspark.sql import SparkSession; spark = SparkSession.builder.appName('ETL Process').getOrCreate()`

## 2. Data Extraction

- **Read Data from CSV**: `df = spark.read.csv('path/to/csv', inferSchema=True, header=True)`
- **Read Data from JSON**: `df = spark.read.json('path/to/json')`
- **Read Data from Parquet**: `df = spark.read.parquet('path/to/parquet')`
- **Read Data from a Database**: `df = spark.read.format("jdbc").option("url", jdbc_url).option("dbtable", "table_name").option("user", "username").option("password", "password").load()`

## 3. Data Transformation

- **Selecting Columns**: `df.select('column1', 'column2')`
- **Filtering Data**: `df.filter(df['column'] > value)`
- **Adding New Columns**: `df.withColumn('new_column', df['column'] + 10)`
- **Renaming Columns**: `df.withColumnRenamed('old_name', 'new_name')`
- **Grouping and Aggregating Data**: `df.groupBy('column').agg({'column2': 'sum'})`
- **Joining DataFrames**: `df1.join(df2, df1['id'] == df2['id'])`
- **Sorting Data**: `df.orderBy(df['column'].desc())`
- **Removing Duplicates**: `df.dropDuplicates()`

## 4. Handling Missing Values

- **Dropping Rows with Missing Values**: `df.na.drop()`
- **Filling Missing Values**: `df.na.fill(value)`
- **Replacing Values**: `df.na.replace(['old_value'], ['new_value'])`

By: Waleed Mousa

## 5. Data Type Conversion

- **Changing Column Types**: `df.withColumn('column', df['column'].cast('new_type'))`
- **Parsing Dates**: `from pyspark.sql.functions import to_date; df.withColumn('date', to_date(df['date_string']))`

## 6. Advanced Data Manipulations

- **Using SQL Queries**: `df.createOrReplaceTempView('table'); spark.sql('SELECT * FROM table WHERE column > value')`
- **Window Functions**: `from pyspark.sql.window import Window; from pyspark.sql.functions import row_number; df.withColumn('row', row_number().over(Window.partitionBy('column').orderBy('other_column')))`
- **Pivot Tables**: `df.groupBy('column').pivot('pivot_column').agg({'column2': 'sum'})`

## 7. Data Loading

- **Writing to CSV**: `df.write.csv('path/to/output')`
- **Writing to JSON**: `df.write.json('path/to/output')`
- **Writing to Parquet**: `df.write.parquet('path/to/output')`
- **Writing to a Database**: `df.write.format("jdbc").option("url", jdbc_url).option("dbtable", "table_name").option("user", "username").option("password", "password").save()`

## 8. Performance Tuning

- **Caching Data**: `df.cache()`
- **Broadcasting a DataFrame for Join Optimization**: `from pyspark.sql.functions import broadcast; df1.join(broadcast(df2), df1['id'] == df2['id'])`
- **Repartitioning Data**: `df.repartition(10)`
- **Coalescing Partitions**: `df.coalesce(1)`

## 9. Debugging and Error Handling

- **Showing Execution Plan**: `df.explain()`

By: Waleed Mousa

- **Catching Exceptions during Read**: Implement try-except blocks during data reading operations.

## 10. Working with Complex Data Types

- **Exploding Arrays**: `from pyspark.sql.functions import explode; df.select(explode(df['array_column']))`
- **Handling Struct Fields**: `df.select('struct_column.field1', 'struct_column.field2')`

## 11. Custom Transformations with UDFs

- **Defining a UDF**: `from pyspark.sql.functions import udf; @udf('return_type') def my_udf(column): return transformation`
- **Applying UDF on DataFrame**: `df.withColumn('new_column', my_udf(df['column']))`

## 12. Working with Large Text Data

- **Tokenizing Text Data**: `from pyspark.ml.feature import Tokenizer; Tokenizer(inputCol='text_column', outputCol='words').transform(df)`
- **TF-IDF on Text Data**: `from pyspark.ml.feature import HashingTF, IDF; HashingTF(inputCol='words', outputCol='rawFeatures').transform(df)`

## 13. Machine Learning Integration

- **Using MLlib for Predictive Modeling**: Building and training machine learning models using PySpark's MLlib.
- **Model Evaluation and Tuning**: `from pyspark.ml.evaluation import MulticlassClassificationEvaluator; MulticlassClassificationEvaluator().evaluate(predictions)`

## 14. Stream Processing

- **Reading from a Stream**: `dfStream = spark.readStream.format('source').load()`
- **Writing to a Stream**: `dfStream.writeStream.format('console').start()`

## 15. Advanced Data Extraction

- **Reading from Multiple Sources**: `df = spark.read.format('format').option('option', 'value').load(['path1', 'path2'])`
- **Incremental Data Loading**: Implementing logic to load data incrementally, based on timestamps or log tables.

## 16. Complex Data Transformations

- **Nested JSON Parsing**: `from pyspark.sql.functions import json_tuple; df.select(json_tuple('json_column', 'field1', 'field2'))`
- **Applying Map-Type Transformations**: Using `map` functions to transform key-value pair data.

## 17. Advanced Joins and Set Operations

- **Broadcast Join with Large and Small DataFrames**: Utilizing `broadcast` for efficient joins.
- **Set Operations (Union, Intersect, Except)**: Performing set operations like `df1.union(df2)`, `df1.intersect(df2)`, `df1.except(df2)`.

## 18. Data Aggregation and Summarization

- **Complex Aggregations**: `df.groupBy('group_col').agg({'num_col1': 'sum', 'num_col2': 'avg'})`
- **Rollup and Cube for Multi-Dimensional Aggregation**: `df.rollup('col1', 'col2').sum()`, `df.cube('col1', 'col2').mean()`

## 19. Advanced Data Filtering

- **Filtering with Complex Conditions**: `df.filter((df['col1'] > value) & (df['col2'] < other_value))`
- **Using Column Expressions**: `from pyspark.sql import functions as F; df.filter(F.col('col1').like('%pattern%'))`

## 20. Working with Dates and Times

- **Date Arithmetic**: df.withColumn('new_date', F.col('date_col') +
  F.expr('interval 1 day'))
- **Date Truncation and Formatting**: df.withColumn('month',
  F.trunc('month', 'date_col'))

## 21. Handling Nested and Complex Structures

- **Working with Arrays and Maps**: df.select(F.explode('array_col')),
  df.select(F.col('map_col')['key'])
- **Flattening Nested Structures**: df.selectExpr('struct_col.*')

## 22. Text Processing and Natural Language Processing

- **Regular Expressions for Text Data**: df.withColumn('extracted',
  F.regexp_extract('text_col', '(pattern)', 1))
- **Sentiment Analysis on Text Data**: Using NLP libraries to perform
  sentiment analysis on textual columns.

## 23. Advanced Window Functions

- **Window Functions for Running Totals and Moving Averages**: from
  pyspark.sql.window import Window; windowSpec =
  Window.partitionBy('group_col').orderBy('date_col');
  df.withColumn('cumulative_sum', F.sum('num_col').over(windowSpec))
- **Ranking and Row Numbering**: df.withColumn('rank',
  F.rank().over(windowSpec))

## 24. Data Quality and Consistency Checks

- **Data Profiling for Quality Assessment**: Generating statistics for
  each column to assess data quality.
- **Consistency Checks Across DataFrames**: Comparing schema and row
  counts between DataFrames for consistency.

## 25. ETL Pipeline Monitoring and Logging

By: Waleed Mousa

- **Implementing Logging in PySpark Jobs**: Using Python's logging module to log ETL process steps.
- **Monitoring Performance Metrics**: Tracking execution time and resource utilization of ETL jobs.

## 26. ETL Workflow Scheduling and Automation

- **Integration with Workflow Management Tools**: Automating PySpark ETL scripts using tools like Apache Airflow or Luigi.
- **Scheduling Periodic ETL Jobs**: Setting up cron jobs or using scheduler services for regular ETL tasks.

## 27. Data Partitioning and Bucketing

- **Partitioning Data for Efficient Storage**:
  ```
  df.write.partitionBy('date_col').parquet('path/to/output')
  ```
- **Bucketing Data for Optimized Query Performance**:
  ```
  df.write.bucketBy(42,
  'key_col').sortBy('sort_col').saveAsTable('bucketed_table')
  ```

## 28. Advanced Spark SQL Techniques

- **Using Temporary Views for SQL Queries**:
  ```
  df.createOrReplaceTempView('temp_view'); spark.sql('SELECT * FROM
  temp_view WHERE col > value')
  ```
- **Complex SQL Queries for Data Transformation**: Utilizing advanced SQL syntax for complex data transformations.

## 29. Machine Learning Pipelines

- **Creating and Tuning ML Pipelines**: Using PySpark's MLlib for building and tuning machine learning pipelines.
- **Feature Engineering in ML Pipelines**: Implementing feature transformers and selectors within ML pipelines.

## 30. Integration with Other Big Data Tools

- **Reading and Writing Data to HDFS**: Accessing Hadoop Distributed File System (HDFS) for data storage and retrieval.
- **Interfacing with Kafka for Real-Time Data Processing**: Connecting to Apache Kafka for stream processing tasks.

## 31. Cloud-Specific PySpark Operations

- **Utilizing Cloud-Specific Storage Options**: Leveraging AWS S3, Azure Blob Storage, or GCP Storage in PySpark.
- **Cloud-Based Data Processing Services Integration**: Using services like AWS Glue or Azure Synapse for ETL processes.

## 32. Security and Compliance in ETL

- **Implementing Data Encryption and Security**: Securing data at rest and in transit during ETL processes.
- **Compliance with Data Protection Regulations**: Adhering to GDPR, HIPAA, or other regulations in data processing.

## 33. Optimizing ETL Processes for Scalability

- **Dynamic Resource Allocation for ETL Jobs**: Adjusting Spark configurations for optimal resource usage.
- **Best Practices for Scaling ETL Processes**: Techniques for scaling ETL pipelines to handle growing data volumes.

By: Waleed Mousa