# 05 PySpark DataFrame Functions

## Comprehensive PySpark DataFrame Functions

PySpark DataFrames are a distributed collection of data organized into named columns. They are conceptually equivalent to a table in a relational database or a data frame in R/Python, but with optimizations for large-scale data processing provided by the Spark engine. This guide will explore essential PySpark DataFrame functions, providing detailed explanations, practical scenarios, and best practices.

## 1. Data Display & Inspection

These functions help you get a quick overview of your DataFrame's structure and content, crucial for initial data exploration and debugging.

### show()

**Purpose**: Displays the top `n` rows of the DataFrame in a tabular format. It's a convenient way to peek at your data without pulling all of it to the driver program.

**Syntax:** `df.show(n=20, truncate=True, vertical=False)`

```python
# Basic usage: Shows the first 20 rows, truncating long strings
df.show()

# Show first 5 rows without truncation: Useful for inspecting full string
values
df.show(5, truncate=False)

# Vertical display for wide DataFrames: Each row is displayed vertically,
# making it easier to read DataFrames with many columns
df.show(5, vertical=True)
```

**Explanation & Scenarios**:

- **Quick Data Exploration**: When you first load a dataset, `show()` is your go-to function to see if the data loaded correctly and to get a sense of its content.

- **Debugging Transformations**: After applying transformations, `show()` helps you verify if the changes were applied as expected. For instance, after adding a new column, you can use `show()` to see its values.

- **Presenting Results**: For small result sets or examples, `show()` can present data clearly in notebooks or command-line interfaces.

- **`truncate` Parameter**: When dealing with columns that might contain very long strings (like descriptions or URLs), `truncate=False` is invaluable to see the full content and avoid misinterpretations due to truncation.

- `vertical` **Parameter**: If your DataFrame has a large number of columns, the horizontal display can become unreadable. `vertical=True` rotates the display, showing one column per line for each row, making it much more readable.

**Best Practices**:

- Always use `truncate=False` when you need to inspect the full content of string columns during debugging or data quality checks.
- Limit the number of rows with the `n` parameter for large datasets to avoid overwhelming your console or notebook output.
- Utilize `vertical=True` for DataFrames with numerous columns to improve readability.

## collect()

**Purpose**: Retrieves all rows from the DataFrame and returns them as a list of `Row` objects to the Spark driver program.

**Syntax**: `df.collect()`

```python
# Collect all superhero data: This will bring all data to the driver
all_heroes = df.collect()
print(f"Total heroes: {len(all_heroes)}")
print(f"First hero: {all_heroes[0]}")
```

**Explanation & Scenarios**:

- **Small Result Sets**: Ideal for situations where you have already filtered or aggregated your data down to a very small size (e.g., a few hundred rows) and need to perform local processing or display it outside of Spark.
- **Lookup Tables**: If you have a small dimension table or a lookup table that fits comfortably in memory, `collect()` can be used to bring it to the driver for efficient lookups in subsequent local operations.
- **Testing with Small Datasets**: During development and testing phases with miniature datasets, `collect()` can be used to assert the correctness of transformations by comparing the collected data with expected outcomes.

**Best Practices**:

- ⚠️ **NEVER** use `collect()` on large datasets without prior filtering or aggregation. Doing so will inevitably lead to **Out-Of-Memory (OOM)** errors on the driver node, crashing your Spark application.
- For retrieving a limited number of rows, prefer `take(n)` or `first()` over `collect()`.
- If you need to iterate over a large DataFrame, consider using `foreach` for side effects, or re-evaluating your approach to leverage Spark's distributed processing capabilities instead of bringing data to the driver.

## select()

**Purpose:** Projects a new DataFrame by selecting a subset of columns or applying expressions to existing columns. It's fundamental for data projection and preparing data for further operations.

**Syntax:** `df.select(*cols)`

```python
# Select specific columns: Creates a new DataFrame with only 'name',
'universe', and 'debut_year'
df.select("name", "universe", "debut_year").show()

# Select with expressions: Renames 'universe' to 'UNIVERSE' and converts
its values to uppercase
df.select(col("name"), upper(col("universe")).alias("UNIVERSE")).show()

# Select all columns except some: Dynamically selects all columns
excluding 'attributes' and 'team'
df.select([c for c in df.columns if c not in ["attributes",
"team"]]).show()
```

**Explanation & Scenarios:**

- **Data Projection:** If you're only interested in a few columns from a wide DataFrame, `select()` allows you to project only those columns, reducing memory footprint and potentially speeding up subsequent operations.
- **Column-Specific Analysis:** Before performing aggregations or complex analyses on specific columns, `select()` can be used to isolate those columns, making your code clearer and more focused.
- **Preparing Data for Joins/Transformations:** Often, you'll need to select specific columns or derive new ones before joining DataFrames or applying complex transformations. `select()` facilitates this pre-processing step.
- **Renaming and Type Casting:** While `withColumnRenamed` is dedicated for renaming, `select()` with `alias()` can also rename columns as part of a projection. It can also be used for type casting by applying casting functions within the `select` statement.

**Best Practices:**

- Always select only the columns you need. This significantly reduces the amount of data processed and shuffled, leading to better performance.
- Utilize column expressions and built-in functions within `select()` for simple transformations and derivations.
- Combine `select()` with other operations like `filter()` or `withColumn()` for efficient data manipulation pipelines.

---

## Additional Inspection Functions

These functions provide further insights into the DataFrame's schema, statistics, and size, aiding in data understanding and validation.

## printSchema()

**Purpose**: Prints the schema of the DataFrame to the console in a tree-like format, showing column names, data types, and nullability.

```
df.printSchema()
```

**Scenario**: Indispensable for understanding the structure of a DataFrame, especially when working with complex nested types or external data sources where the schema might not be immediately obvious.

## describe()

**Purpose**: Computes basic descriptive statistics for numerical and string columns in the DataFrame.

```
df.describe().show()  # Statistical summary for numerical and string
columns
```

**Scenario**: Quickly get summary statistics like count, mean, stddev, min, and max for numerical columns, and count, unique, and frequency for string columns. Useful for initial data profiling and anomaly detection.

## dtypes

**Purpose**: Returns a list of tuples, where each tuple contains the column name and its data type as a string.

```
print(df.dtypes)  # List of (column_name, data_type) tuples
```

**Scenario**: Programmatically access column names and their types, useful for schema validation or dynamically processing columns based on their data type.

## count()

**Purpose**: Returns the number of rows in the DataFrame.

```
total_heroes = df.count()
print(f"Total number of heroes in the dataset: {total_heroes}")
```

**Scenario**: Get the total record count for data validation or to verify the size of a filtered dataset.

## first() and take()

**Purpose**: `first()` returns the first row as a `Row` object. `take(n)` returns the first `n` rows as a list of `Row` objects.

```
first_hero = df.first()
print(f"Details of the first hero: {first_hero}")

first_three = df.take(3)
print(f"Details of the first three heroes: {first_three}")
```

**Scenario**: Similar to `show(1)` or `show(n)` but returns the data as Python objects, making it easier for programmatic access in the driver program for quick checks or small-scale processing. Unlike `collect()`, these are safe to use on large DataFrames as they only retrieve a limited number of rows.

## Scenario-Based Questions & Answers (Data Display & Inspection)

1. **Scenario**: You've just loaded a new dataset of "Villains" and want to quickly see the first 10 entries and their full `origin_story` column, which you suspect might be very long.

   - **Question**: How would you achieve this using `show()`?
   - **Answer**:

     ```
     # Assuming 'villains_df' is your DataFrame
     # villains_df.show(10, truncate=False)
     df.show(5, truncate=False) # Using 'df' for demonstration
     ```

2. **Scenario**: Your team needs to confirm the exact number of unique "powers" present across all superheroes to build a categorization system. The `powers` column is an `ArrayType(StringType())`.

   - **Question**: How would you get this count without bringing the entire DataFrame to the driver, and then confirm how many rows are in the `df` DataFrame?
   - **Answer**:

```
# To get unique powers, we'd first explode the array, then count
distinct.
# This is a transformation, not a direct inspection function like
count().
# However, to answer the second part of the question (total rows in
df):
total_rows = df.count()
print(f"Total rows in df: {total_rows}")

from pyspark.sql.functions import explode, countDistinct

# Part 1: Get count of unique powers without bringing entire DataFrame
to driver
```

```python
unique_powers_count =
df.select(explode(col("powers")).alias("individual_power")) \

.select(countDistinct("individual_power").alias("unique_power_count")) \

                        .collect()[0]["unique_power_count"]

print(f"Number of unique powers: {unique_powers_count}")

# Part 2: Confirm total rows in DataFrame
total_rows = df.count()
print(f"Total rows in df: {total_rows}")

# Alternative approach to see the actual unique powers (if needed for
categorization):
unique_powers_list =
df.select(explode(col("powers")).alias("individual_power")) \
                        .distinct() \
                        .collect()

print("Unique powers are:")
for row in unique_powers_list:
    print(f"- {row['individual_power']}")
```

3. **Scenario:** You are debugging a complex transformation on your `df` DataFrame and want to inspect the schema after several `withColumn` operations to ensure the data types are as expected.

   - **Question:** Which function would you use to view the updated schema?
   - **Answer:**

   ```python
   df.printSchema()
   ```

4. **Scenario:** Your data scientist wants to get a quick summary of the `debut_year` column (min, max, average) and also see how many unique universes are present in the dataset.

   - **Question:** Which two inspection functions would be most efficient for this task?
   - **Answer:**

   ```python
   df.describe("debut_year", "universe").show()
   # Alternatively, for unique universe count specifically:
   # df.select("universe").distinct().count()
   ```

5. **Scenario:** You need to grab the very first superhero record to manually check its values, but you want to avoid `collect()` due to the potentially large dataset.

- **Question**: How would you retrieve just the first row?
- **Answer**:

```
first_record = df.first()
print(f"First record: {first_record}")
```

## 2. Column Operations

These functions allow you to manipulate, add, rename, or remove columns, which are crucial steps in feature engineering and data preparation.

### withColumn()

**Purpose**: Adds a new column to the DataFrame or replaces an existing one with the result of a given expression. This is a workhorse function for deriving new features.

**Syntax**: `df.withColumn(colName, col)`

```
# Add a new column based on condition: Categorizes heroes into 'Golden
Age', 'Silver Age', or 'Modern Age'
df_with_era = df.withColumn("era",
    when(col("debut_year") < 1950, "Golden Age")
    .when(col("debut_year") < 1980, "Silver Age")
    .otherwise("Modern Age")
)
df_with_era.select("name", "debut_year", "era").show()

# Add column with array size: Calculates the number of powers each hero
has
df_with_power_count = df.withColumn("power_count", size(col("powers")))
df_with_power_count.select("name", "powers",
"power_count").show(truncate=False)

# Add calculated column: Determines how many years a hero has been active
df_with_age = df.withColumn("years_active", year(current_date()) -
col("debut_year"))
df_with_age.select("name", "debut_year", "years_active").show()
```

**Explanation & Scenarios**:

- **Feature Engineering**: Create new features from existing ones. For example, calculating `age` from `birth_date`, or `total_price` from `quantity` and `unit_price`.
- **Data Categorization**: Apply conditional logic to categorize data into buckets, like the "era" example above, or assigning "high", "medium", "low" labels based on a numerical score.

- **Data Type Conversion**: Although `cast()` is more direct, `withColumn()` can be used to re-create a column with a different data type.
- **Cleaning and Standardization**: Apply string functions (e.g., `trim`, `lower`, `upper`) to clean textual data within a column.

**Best Practices**:

- **Chain `withColumn()` calls**: For multiple column additions or modifications, chaining `withColumn()` calls improves readability and can sometimes be more optimized than separate assignments.
- **Meaningful Column Names**: Always use descriptive and consistent naming conventions for new columns.
- **Performance Impact**: Be mindful of complex expressions. While Spark optimizes many operations, highly complex UDFs or very convoluted conditional logic can impact performance.

## withColumnRenamed()

**Purpose**: Renames an existing column in the DataFrame. This is useful for improving readability, adhering to naming conventions, or preparing data for downstream systems.

**Syntax**: `df.withColumnRenamed(existing, new)`

```python
# Rename single column: Changes 'debut_year' to 'first_appearance_year'
df_renamed = df.withColumnRenamed("debut_year", "first_appearance_year")
df_renamed.printSchema()

# Rename multiple columns: Renames 'name' to 'hero_name' and 'alias' to
'real_name'
df_multi_renamed = df.withColumnRenamed("name", "hero_name") \
                     .withColumnRenamed("alias", "real_name")
df_multi_renamed.printSchema()
```

**Explanation & Scenarios**:

- **Standardizing Naming Conventions**: Ensure all column names follow a consistent style (e.g., snake_case, camelCase).
- **Clarity and Readability**: Rename ambiguously named columns to be more descriptive.
- **Integration with Other Systems**: Sometimes, external systems or databases might have strict naming requirements, making renaming necessary before export.
- **Avoiding Name Collisions**: Before a join, you might rename columns in one DataFrame to avoid ambiguity if both DataFrames have columns with identical names but different meanings.

**Best Practices**:

- Use descriptive and unambiguous new column names.
- Follow your project's established naming conventions.

- When renaming multiple columns, chaining `withColumnRenamed()` is a clean approach.

## drop()

**Purpose**: Removes one or more columns from the DataFrame. This is crucial for reducing data size, removing sensitive information, or simplifying your DataFrame for specific analyses.

**Syntax:** `df.drop(*cols)`

```python
# Drop single column: Removes the 'attributes' column
df_no_attributes = df.drop("attributes")
df_no_attributes.printSchema()

# Drop multiple columns: Removes 'attributes', 'team', and
'first_appearance'
df_minimal = df.drop("attributes", "team", "first_appearance")
df_minimal.printSchema()

# Drop columns conditionally: Imagine a scenario where you want to drop
'alias'
# if a flag `drop_alias_column` is True.
drop_alias_column = True
columns_to_drop = ["alias"] if drop_alias_column else []
df_conditional_drop = df.drop(*columns_to_drop)
df_conditional_drop.printSchema()
```

**Explanation & Scenarios**:

- **Data Size Reduction**: Removing unnecessary columns reduces memory consumption and disk I/O, which can significantly improve performance for large datasets.
- **Privacy and Security**: Eliminate columns containing sensitive personal identifiable information (PII) or other confidential data.
- **Simplifying Data Models**: For specific analytical tasks or machine learning models, you might only need a subset of features, so dropping irrelevant columns simplifies the DataFrame.
- **Before Writing to Storage**: It's good practice to drop temporary or intermediate columns before writing the final DataFrame to a persistent storage.

**Best Practices:**

- Remove unused columns as early as possible in your data processing pipeline to minimize resource usage.
- For multiple columns, pass them as separate arguments or as a list to `drop()`.
- Consider using `select()` if you know exactly which columns you want to *keep*, as it can sometimes be more explicit than dropping many columns.

# pivot()

**Purpose**: Rotates a table-valued expression by turning the unique values from one column (the pivot column) into multiple new columns. This is essential for transforming "long" format data into "wide" format, often used for reporting or specific analytical needs.

**Syntax:** `df.groupBy(*cols).pivot(pivot_col).agg(*exprs)`

```python
# Pivot universe by debut decade: Counts heroes per universe per decade
df_pivot = df.withColumn("decade", floor(col("debut_year") / 10) * 10) \
             .groupBy("decade") \
             .pivot("universe") \
             .agg(count("name").alias("hero_count"))
df_pivot.show()
```

**Explanation & Scenarios**:

- **Creating Cross-Tabulations**: Generate summary tables where rows represent one category and columns represent another, with aggregated values at their intersection. For example, sales by product category and region.

- **Transforming Long to Wide Format**: If your data is in a "long" format (e.g., each power listed as a separate row for a hero), `pivot()` can transform it into a "wide" format where each power becomes a column.

- **Reporting and Analytics**: Often, business intelligence tools or dashboards prefer data in a wide format for easier visualization and reporting.

**Best Practices**:

- Always use `pivot()` in conjunction with `groupBy()` and an aggregation function (`agg()`). Without `agg()`, it won't perform any actual data transformation beyond creating the pivoted structure.

- **Limit Pivot Values**: Be cautious with pivot columns that have a very high cardinality (many unique values). Pivoting on such columns can lead to a DataFrame with an excessive number of columns, causing performance issues and making the data unmanageable. If you need to pivot on a high-cardinality column, consider pre-filtering the unique values.

- **Performance Implications**: Pivoting involves significant data shuffling across the cluster, which can be computationally expensive for very large DataFrames. Plan and test carefully.

---

# Scenario-Based Questions & Answers (Column Operations)

1. **Scenario:** You want to create a new column called `is_marvel_hero` which is `True` if the hero belongs to the "Marvel" universe and `False` otherwise.

   - **Question:** How would you achieve this using `withColumn()`?
   - **Answer:**

```
df_is_marvel = df.withColumn("is_marvel_hero", (col("universe") ==
"Marvel"))
df_is_marvel.select("name", "universe", "is_marvel_hero").show(5)
```

2. **Scenario**: Your data governance team has mandated that the `alias` column should now be named `secret_identity` for consistency across all datasets.
   - **Question**: How would you rename this column in the DataFrame?
   - **Answer**:

   ```
   df_renamed_alias = df.withColumnRenamed("alias", "secret_identity")
   df_renamed_alias.printSchema()
   ```

3. **Scenario**: For a simplified analysis, you only need the `name`, `universe`, and `powers` columns. All other columns (`alias`, `debut_year`, `attributes`, `team`, `first_appearance`) are not required.
   - **Question**: How would you remove all these unnecessary columns efficiently?
   - **Answer**:

   ```
   df_selected_columns = df.drop("alias", "debut_year", "attributes",
   "team", "first_appearance")
   df_selected_columns.printSchema()
   # Alternatively, using select for clarity if you know what to keep:
   # df_selected_columns = df.select("name", "universe", "powers")
   ```

4. **Scenario**: Your analytics team wants to see a cross-tabulation showing how many heroes from each `universe` have `Super Strength` versus those that don't.
   - **Question**: How would you use `pivot()` to show the count of heroes by universe based on whether they possess "Super Strength"?
   - **Answer**:

   ```
   df_strength_pivot = df.withColumn("has_super_strength",
   array_contains(col("powers"), "Super Strength")) \
                       .groupBy("universe") \
                       .pivot("has_super_strength") \
                       .agg(count("name").alias("hero_count"))
   df_strength_pivot.show()
   ```

5. **Scenario**: You want to add a column indicating the length of the `name` of each hero, and then immediately drop the `alias` column.
   - **Question**: How would you chain these two column operations together for readability and efficiency?

- **Answer:**

```python
df_chained_ops = df.withColumn("name_length", size(col("name"))) \
                .drop("alias")
df_chained_ops.select("name", "name_length",
"alias").show(truncate=False) # 'alias' won't be shown
df_chained_ops.printSchema()
```

# 3. Data Filtering & Querying

These functions allow you to select a subset of rows based on specified conditions, which is fundamental for focused analysis and data preparation.

## `where()` & `filter()`

**Purpose:** Both `where()` and `filter()` perform the same function: they filter rows from a DataFrame based on a given boolean condition. They are aliases for each other, so you can use whichever you find more readable or consistent with your background (e.g., `where` for SQL users, `filter` for R/Python users).

**Syntax:** `df.where(condition)` or `df.filter(condition)`

```python
# Filter Marvel heroes: Selects all heroes from the "Marvel" universe
marvel_heroes = df.filter(col("universe") == "Marvel")
marvel_heroes.show()

# Multiple conditions: Filters for Marvel heroes who debuted after 1960
modern_marvel = df.where((col("universe") == "Marvel") &
(col("debut_year") > 1960))
modern_marvel.show()

# String operations: Finds heroes whose names contain "Spider"
spider_heroes = df.filter(col("name").contains("Spider"))
spider_heroes.show()

# Array operations: Finds heroes who have "Super Strength" as one of their
powers
strong_heroes = df.filter(array_contains(col("powers"), "Super Strength"))
strong_heroes.show()

# Null filtering: Filters for heroes who have a team affiliation (team
column is not null)
heroes_with_teams = df.filter(col("team").isNotNull())
heroes_with_teams.show()
```

**Explanation & Scenarios:**

- **Data Subset Selection**: Extract specific subsets of data relevant to a particular analysis or report. For example, selecting only sales data from a particular region or for a specific product.
- **Data Validation**: Filter out erroneous or incomplete records based on certain criteria (e.g., records where a mandatory field is null, or numerical values are out of a valid range).
- **Pre-processing for ML**: Before training a machine learning model, you might filter out outliers or irrelevant data points.
- **Conditional Logic**: Apply complex boolean logic using `&` (AND), `|` (OR), and `~` (NOT) to combine multiple conditions.

**Best Practices:**

- **Filter Early**: Apply filters as early as possible in your data processing pipeline. This significantly reduces the amount of data that needs to be processed in subsequent steps, leading to better performance.
- **Use `col()` for Column References**: Always use `col()` or `df.column_name` when referring to columns within filter conditions to ensure Spark treats them as column expressions.
- **Parentheses for Clarity**: For complex conditions involving `&` and `|`, use parentheses to explicitly define the order of operations and improve readability.

## `sample()` vs `sampleBy()`

**Purpose**: These functions allow you to extract a random subset of rows from your DataFrame. This is especially useful for data exploration, prototyping, or creating smaller datasets for testing and machine learning.

### `sample()`

**Purpose**: Performs simple random sampling of rows from the DataFrame.

```python
# Random sample - 30% of data without replacement (default)
sample_df = df.sample(fraction=0.3, seed=42)
sample_df.show()
print(f"Sampled {sample_df.count()} rows out of {df.count()}")

# Sample with replacement - 50% of data with replacement
sample_with_replacement = df.sample(withReplacement=True, fraction=0.5,
seed=42)
sample_with_replacement.show()
print(f"Sampled with replacement: {sample_with_replacement.count()} rows")
```

**Explanation:**

- `fraction`: The approximate fraction of the rows to be sampled.
- `withReplacement`: If `True`, rows can be sampled multiple times (sampling with replacement). If `False` (default), each row can be sampled at most once.

- `seed` : An optional seed for the random number generator. Providing a seed ensures reproducible results across multiple runs.

## sampleBy()

**Purpose**: Performs stratified sampling, allowing you to sample different fractions of data for different categories within a specified column.

```python
# Stratified sampling by universe: Sample 30% of Marvel heroes and 50% of
DC heroes
stratified_sample = df.sampleBy("universe", fractions={"Marvel": 0.3,
"DC": 0.5}, seed=42)
stratified_sample.show()
print(f"Stratified Sample Counts (Marvel):
{stratified_sample.filter(col('universe') == 'Marvel').count()}")
print(f"Stratified Sample Counts (DC):
{stratified_sample.filter(col('universe') == 'DC').count()}")
```

**Explanation**:

- `col` : The column by which to stratify the sampling.
- `fractions` : A dictionary where keys are the categories in `col` and values are the desired sampling fractions for each category.
- `seed` : An optional seed for reproducibility.

**Explanation & Scenarios**:

- **Data Exploration**: Quickly create a smaller dataset to perform initial explorations or generate summary statistics without processing the entire large dataset.
- **Creating Training/Test Sets**: For machine learning, sampling can be used to split your data into training, validation, and test sets, especially `sampleBy()` for maintaining class distribution in stratified sampling.
- **Performance Testing**: Test the performance of your Spark jobs with a smaller, representative subset of data before running them on the full dataset.
- **Handling Imbalanced Data**: `sampleBy()` is particularly useful in machine learning contexts when dealing with imbalanced datasets, allowing you to undersample or oversample specific classes to balance the distribution.

**Best Practices**:

- Always set a `seed` when sampling if you need reproducible results. This is crucial for debugging and ensuring consistency in analyses or model training.
- Use `sampleBy()` when it's important to maintain the proportional representation of different categories within your sample, preventing bias.
- Be aware of potential **data skewness** when using `sample()`. If certain categories are very rare, a simple random sample might miss them entirely. `sampleBy()` can mitigate this.

## Suggested Readings:

1. https://www.scribbr.com/methodology/sampling-methods/
2. https://www.scribbr.com/methodology/stratified-sampling/

# Scenario-Based Questions & Answers (Data Filtering & Querying)

1. **Scenario**: You need a list of all DC universe heroes who debuted before 1940.

   - **Question**: How would you filter the DataFrame to get this specific subset of heroes?

   - **Answer**:

   ```python
   dc_early_heroes = df.filter((col("universe") == "DC") &
   (col("debut_year") < 1940))
   dc_early_heroes.show()
   ```

2. **Scenario**: The marketing team wants to identify all Marvel heroes who are also part of the "Defenders" team.

   - **Question**: How would you filter the DataFrame to find these heroes, given that `team` is an `ArrayType`?

   - **Answer**:

   ```python
   marvel_defenders = df.filter((col("universe") == "Marvel") &
   (array_contains(col("team"), "Defenders")))
   marvel_defenders.show()
   ```

3. **Scenario**: You are developing a new feature and need a quick, representative sample of 10% of your entire dataset for rapid prototyping, without worrying about specific categorical distributions.

   - **Question**: Which sampling function would you use, and how would you ensure that your sample is the same every time you run the code?

   - **Answer**:

   ```python
   proto_sample = df.sample(fraction=0.1, seed=123) # Using sample()
   for simple random sampling
   proto_sample.show()
   ```

4. **Scenario**: For a machine learning task, you need to create a training dataset where you have 70% of Marvel heroes and 80% of DC heroes, maintaining their respective proportions within the sample.

   - **Question**: How would you achieve this using a specific sampling method?

   - **Answer**:

```
ml_training_data = df.sampleBy("universe", fractions={"Marvel": 0.7,
"DC": 0.8}, seed=456)
ml_training_data.show()
```

5. **Scenario**: You've received a new batch of data, and you suspect some entries might have null values in the `name` column, which is critical for your analysis. You want to see if any such records exist.

   - **Question**: How would you filter the DataFrame to display only rows where the `name` column is null?

   - **Answer**:

```
heroes_with_null_name = df.filter(col("name").isNull())
heroes_with_null_name.show()
```

# 4. Data Transformation

These functions enable you to apply complex logic to your DataFrames, either by defining reusable functions or creating custom user-defined functions (UDFs).

## transform()

**Purpose**: The `transform()` method is a high-level function that allows you to chain multiple DataFrame transformations together in a clean and readable way. It takes a function that accepts a DataFrame and returns a DataFrame. This promotes code reusability and modularity.

**Syntax**: `df.transform(func)`

```
# Define reusable transformation functions
def add_hero_category(df_input):
    return df_input.withColumn("category",
        when(array_contains(col("powers"), "Magic"), "Mystic")
        .when(array_contains(col("powers"), "Genius Intellect"), "Tech")
        .otherwise("Physical")
    )

def add_power_level(df_input):
    return df_input.withColumn("power_level",
        when(size(col("powers")) ≥ 3, "High")
        .when(size(col("powers")) ≥ 2, "Medium")
        .otherwise("Low")
    )

# Apply transformations using transform():
# This applies add_hero_category first, then passes the result to
add_power_level
df_transformed =
```

```python
df.transform(add_hero_category).transform(add_power_level)
df_transformed.select("name", "powers", "category",
"power_level").show(truncate=False)

# Another scenario: Apply a series of cleaning steps
def clean_hero_names(df_input):
    return df_input.withColumn("name", upper(trim(col("name"))))

def normalize_universe(df_input):
    return df_input.withColumn("universe", upper(col("universe")))

df_cleaned = df.transform(clean_hero_names).transform(normalize_universe)
df_cleaned.select("name", "universe").show()
```

**Explanation & Scenarios:**

- **Reusable Transformation Pipelines**: Encapsulate common data cleaning, feature engineering, or enrichment steps into functions that can be easily applied to different DataFrames or at various stages of a pipeline.
- **Modular Code:** Break down complex data processing logic into smaller, manageable, and testable functions.
- **Readability**: Chaining `transform()` calls often makes the sequence of operations clearer than deeply nested method calls.
- **ETL/ELT Workflows**: Define and reuse standard transformations for different data sources or stages within an ETL/ELT process.

**Best Practices**:

- Ensure that the function passed to `transform()` always accepts a DataFrame as input and returns a DataFrame as output.
- Use meaningful function names that describe the transformation being performed.
- `transform()` itself is not a performance optimizer but a code organizer. The performance depends on the underlying Spark DataFrame operations within your functions.

## `map()` (for RDDs/Datasets) and `flatMap()` (for RDDs/Datasets)

**Purpose**: These are lower-level RDD transformations. While PySpark DataFrames are the preferred API for most tasks, understanding `map()` and `flatMap()` can be useful when you need more fine-grained control or are working with unstructured data that requires RDD-level processing before converting to a DataFrame. They operate on each element (row) of an RDD.

### `map()`

**Purpose**: Transforms each element of the RDD independently. It applies a function to each row and returns a new RDD with the results. There is a one-to-one mapping between input and output elements.

```
# Convert DataFrame to RDD, apply map, convert back
# Scenario: Create a new RDD with uppercase hero names and their universe
rdd_transformed = df.rdd.map(lambda row: (row.name.upper(), row.universe))
df_mapped = spark.createDataFrame(rdd_transformed, ["name_upper",
"universe"])
df_mapped.show()
```

## flatMap()

**Purpose**: Similar to `map()`, but each input element can be mapped to zero or more output elements. The results from all elements are then "flattened" into a single RDD.

```
# Flatten powers array
# Scenario: Create a new DataFrame where each row represents a (hero_name,
power) pair,
# effectively flattening the 'powers' array column.
powers_rdd = df.rdd.flatMap(lambda row: [(row.name, power) for power in
row.powers])
df_powers = spark.createDataFrame(powers_rdd, ["hero_name", "power"])
df_powers.show()
```

**Explanation & Scenarios**:

- **Data Parsing (RDD-level)**: If you're reading raw, unstructured text files (e.g., log files, CSVs with inconsistent formats), `map()` and `flatMap()` are often used in the initial RDD processing steps to parse each line into a structured format before converting to a DataFrame.

- **Complex Row-Level Logic**: When a transformation is too complex or custom for built-in DataFrame functions or even UDFs, you might resort to RDD operations. However, this is increasingly rare with the rich DataFrame API.

- **Flattening Nested Structures (complex scenarios)**: `flatMap()` is excellent for expanding collections within records into individual records, as shown with the `powers` array example.

**Best Practices:**

- **Prefer DataFrames**: Always try to use DataFrame operations (e.g., `withColumn` with built-in functions, `explode`) over RDD `map` / `flatMap` for performance and conciseness, especially for structured data. Spark's Catalyst optimizer cannot optimize RDD operations as effectively.

- **Conversion**: Remember to convert RDDs back to DataFrames using `spark.createDataFrame()` if you intend to continue with DataFrame operations.

## UDF (User Defined Functions)

**Purpose:** UDFs allow you to extend Spark's functionality by writing custom functions in Python (or Scala, Java) that can be applied to DataFrame columns. Use them when Spark's built-in functions don't cover your specific business logic.

**Syntax:**

1. Define a Python function.
2. Import `udf` from `pyspark.sql.functions` and a relevant `pyspark.sql.types`.
3. Decorate your Python function with `@udf(returnType= ... )` or register it using `spark.udf.register()`.
4. Apply the UDF using `withColumn()`.

```python
# Define UDF

@udf(returnType=StringType())
def get_power_category(powers):
    if powers is None: # Handle potential null input for robustness
        return "Unknown"
    if any("Magic" in power for power in powers):
        return "Mystic"
    elif any("Intellect" in power for power in powers):
        return "Tech"
    else:
        return "Physical"

# Apply UDF
df_with_category = df.withColumn("power_category",
get_power_category(col("powers")))
df_with_category.select("name", "powers",
"power_category").show(truncate=False)

# Another UDF scenario: Calculate a custom score based on attributes map
@udf(returnType=IntegerType())
def calculate_attribute_score(attributes_map):
    if attributes_map is None:
        return 0
    score = 0
    if "intelligence" in attributes_map and attributes_map["intelligence"]
== "high":
        score += 5
    if "wealth" in attributes_map and attributes_map["wealth"] ==
"billionaire":
        score += 3
    if "morality" in attributes_map and attributes_map["morality"] ==
"heroic":
        score += 2
    return score
```

```
df_with_score = df.withColumn("attribute_score",
calculate_attribute_score(col("attributes")))
df_with_score.select("name", "attributes",
"attribute_score").show(truncate=False)
```

**Explanation & Scenarios**:

- **Custom Business Logic**: Implement specific domain-driven rules that are not covered by standard SQL or built-in Spark functions.
- **Complex String Parsing**: When regular expressions or built-in string functions are insufficient for parsing complex text fields.
- **External Library Integration**: If your transformation requires functions from external Python libraries not natively supported by Spark SQL (e.g., custom NLP models, specific cryptographic functions).

**Best Practices**:

- **Prefer Built-in Functions**: Always prioritize Spark's built-in functions over UDFs. Built-in functions are highly optimized and executed natively on the JVM, offering significantly better performance. UDFs involve serialization/deserialization between Python and JVM, causing overhead.
- **Type Annotations**: Always specify the `returnType` for your UDF. This is crucial for Spark to optimize its execution plan and avoid runtime errors.
- **Vectorized UDFs (Pandas UDFs)**: For better performance with UDFs, especially when operating on Pandas Series or DataFrames, consider using PySpark's Pandas UDFs (also known as Vectorized UDFs). These can offer substantial performance improvements by executing the Python function on batches of data.
- **Testing**: Thoroughly test your UDFs with various edge cases, including null values and empty collections.

## Scenario-Based Questions & Answers (Data Transformation)

1. **Scenario:** You want to assign a `strength_tier` to each hero based on whether they have "Super Strength", "Enhanced Strength" (which doesn't exist in our current data but could be a future power), or no strength-related powers. If "Super Strength", tier is "High"; if "Enhanced Strength", tier is "Medium"; otherwise, "Low".

   - **Question:** How would you define a `transform()` function to add this `strength_tier` column?
   - **Answer:**

   ```
   def add_strength_tier(df_input):
       return df_input.withColumn("strength_tier",
           when(array_contains(col("powers"), "Super Strength"),
   ```

```
"High")
        .when(array_contains(col("powers"), "Enhanced Strength"),
"Medium") # Example for potential future power
        .otherwise("Low")
    )
df_with_strength_tier = df.transform(add_strength_tier)
df_with_strength_tier.select("name", "powers",
"strength_tier").show(truncate=False)
```

2. **Scenario:** You need to calculate a `power_density` score for each hero, defined as the number of powers divided by their `debut_year` (as a rough inverse relationship for older heroes potentially having fewer explicitly listed powers). This requires a custom calculation.

   - **Question:** How would you use a UDF to calculate this `power_density`?
   - **Answer:**

```
@udf(returnType=DoubleType()) # Assuming
pyspark.sql.types.DoubleType is imported
def calculate_power_density(powers_list, debut_year):
    if powers_list is None or debut_year is None or debut_year == 0:
        return 0.0
    return len(powers_list) / debut_year

from pyspark.sql.types import DoubleType # Import required type

df_power_density = df.withColumn("power_density",
calculate_power_density(col("powers"), col("debut_year")))
df_power_density.select("name", "powers", "debut_year",
"power_density").show(truncate=False)
```

3. **Scenario:** You have hero records and need to list each hero's individual power as a separate row. For example, Spider-Man would appear three times, once for each of his powers.

   - **Question:** While `explode()` is ideal for DataFrames, how would you achieve this flattening using `flatMap()` on an RDD, and then convert it back to a DataFrame?
   - **Answer:**

```
# Already demonstrated in the flatMap() section, reiterating for
clarity.
powers_expanded_rdd = df.rdd.flatMap(lambda row: [(row.name, p) for
p in row.powers])
df_powers_expanded = spark.createDataFrame(powers_expanded_rdd,
```

```
["hero_name", "individual_power"])
df_powers_expanded.show()
```

4. **Scenario:** You want to standardize the `universe` column by ensuring all entries are in uppercase and remove any leading/trailing whitespace from the `name` column. You'd like to bundle these cleaning steps into a reusable function for `transform()`.

   - **Question:** How would you create and apply a `transform()` function that performs these two cleaning operations?

   - **Answer:**

```
def standardize_hero_data(df_input):
    return df_input.withColumn("universe", upper(col("universe"))) \
                   .withColumn("name", trim(col("name")))

df_standardized = df.transform(standardize_hero_data)
df_standardized.select("name", "universe").show()
```

5. **Scenario:** You need to create a UDF that checks if a hero's `debut_year` is considered "classic" (before 1960). The UDF should return "Classic" or "Modern".

   - **Question:** Implement this UDF and apply it to the DataFrame.

   - **Answer:**

```
@udf(returnType=StringType())
def get_debut_era(debut_year):
    if debut_year is None:
        return "Unknown"
    return "Classic" if debut_year < 1960 else "Modern"

df_with_era_udf = df.withColumn("debut_era_category",
get_debut_era(col("debut_year")))
df_with_era_udf.select("name", "debut_year",
"debut_era_category").show()
```

# 5. Data Cleaning & Deduplication

These functions are crucial for maintaining data quality by handling missing values and removing redundant records.

## dropDuplicates()

**Purpose:** Removes duplicate rows from a DataFrame. You can either remove exact duplicate rows across all columns or specify a subset of columns to consider for duplication.

**Syntax:** `df.dropDuplicates(subset=None)`

```python
# Create some duplicate data for demonstration
duplicate_data = [
    ("Spider-Man", "Marvel", "Peter Parker", 1962, ["Wall-Crawling"],
{"int": "high"}, ["Avengers"], date(1962, 8, 1)),
    ("Spider-Man", "Marvel", "Peter Parker", 1962, ["Wall-Crawling"],
{"int": "high"}, ["Avengers"], date(1962, 8, 1)), # Duplicate
    ("Iron Man", "Marvel", "Tony Stark", 1963, ["Genius Intellect"],
{"wealth": "billionaire"}, ["Avengers"], date(1963, 3, 1)),
    ("Iron Man", "Marvel", "Tony Stark", 1963, ["Genius Intellect"],
{"wealth": "billionaire"}, ["Avengers"], date(1963, 3, 2)) # Different
first_appearance
]
df_dup = spark.createDataFrame(duplicate_data, schema)
print("Original DataFrame with duplicates:")
df_dup.show(truncate=False)

# Remove complete duplicates: Drops rows that are identical across all
columns
df_no_duplicates = df_dup.dropDuplicates()
print("DataFrame after dropping all duplicates:")
df_no_duplicates.show(truncate=False)

# Remove duplicates based on specific columns: Keeps one record for each
unique (name, universe) pair
df_unique_names = df_dup.dropDuplicates(["name", "universe"])
print("DataFrame after dropping duplicates on 'name' and 'universe':")
df_unique_names.show(truncate=False)

# Remove duplicates considering only name: If multiple entries exist for
the same hero name, keep only one.
# Note: Spark does not guarantee which of the duplicate rows is kept,
typically it's the first one encountered during processing.
df_unique_heroes = df_dup.dropDuplicates(["name"])
print("DataFrame after dropping duplicates on 'name':")
df_unique_heroes.show(truncate=False)
```

**Explanation & Scenarios**:

- **Data Integrity**: Ensure that each record is unique based on a primary key or a combination of identifying columns.

- **Preventing Double Counting**: In aggregation scenarios, duplicates can lead to inflated counts or sums. Removing them ensures accurate metrics.

- **Merging Data**: After combining data from multiple sources (e.g., using `union()`), `dropDuplicates()` is often used to consolidate redundant entries.

- **"Golden Record" Selection**: While `dropDuplicates()` doesn't explicitly allow you to choose *which* duplicate to keep (like the latest one), it implicitly keeps one. For more

control, you'd typically use `Window Functions` with `row_number()` and then filter.

**Best Practices:**

- **Define `subset`**: Always specify the `subset` of columns that uniquely identify a record. If not specified, Spark considers all columns, which might be too strict if minor non-key differences exist.

- **Understand Data Grain**: Before deduplicating, clearly define what constitutes a "duplicate" in your dataset. Is it a duplicate if all columns are the same, or just a few key identifiers?

- **Performance**: Deduplication can be a costly operation as it involves shuffling and sorting data across partitions. Apply it judiciously.

---

## `fillna()` & `fill()`

**Purpose**: These are aliases for handling null (or NaN in the case of numerical columns) values in a DataFrame. They allow you to replace nulls with a specified value.

**Syntax**: `df.fillna(value, subset=None)` or `df.fill(value, subset=None)`

```python
# Create a DataFrame with some null values for demonstration
data_with_nulls = [
    ("Spider-Man", "Marvel", "Peter Parker", 1962, ["Wall-Crawling"],
{"intelligence": "high"}, ["Avengers"], date(1962, 8, 1)),
    ("Unknown Hero", None, "Alias Unknown", None, [], {}, [], None), #
Contains multiple nulls
    ("Iron Man", "Marvel", "Tony Stark", 1963, ["Genius Intellect"],
{"wealth": "billionaire"}, ["Avengers"], date(1963, 3, 1))
]
schema_with_nulls = StructType([
    StructField("name", StringType(), True),
    StructField("universe", StringType(), True),
    StructField("alias", StringType(), True),
    StructField("debut_year", IntegerType(), True),
    StructField("powers", ArrayType(StringType()), True),
    StructField("attributes", MapType(StringType(), StringType()), True),
    StructField("team", ArrayType(StringType()), True),
    StructField("first_appearance", DateType(), True)
])
df_nulls = spark.createDataFrame(data_with_nulls, schema_with_nulls)
print("Original DataFrame with nulls:")
df_nulls.show(truncate=False)

# Fill all null values with a specific value (e.g., "Unknown" for strings,
0 for integers)
# This applies to all columns that match the type of the fill value
df_filled_global = df_nulls.fillna("Unknown", subset=["universe",
"alias"])
```

```python
df_filled_global = df_filled_global.fillna(0, subset=["debut_year"])
print("DataFrame after global fillna (string and int):")
df_filled_global.show(truncate=False)

# Fill specific columns using a dictionary: Recommended for heterogeneous
null filling
df_filled_selective = df_nulls.fillna({"universe": "Unspecified",
"debut_year": 1900})
print("DataFrame after selective fillna with dictionary:")
df_filled_selective.show(truncate=False)

# Fill with different values by column type (more robust approach)
df_filled_typed = df_nulls.fillna({
    "name": "Unnamed Hero",
    "debut_year": 1900,
    "powers": [], # Fill ArrayType nulls with empty array
    "team": [],   # Fill ArrayType nulls with empty array
    "first_appearance": date(1900, 1, 1) # Fill DateType nulls
})
print("DataFrame after typed fillna:")
df_filled_typed.show(truncate=False)
```

**Explanation & Scenarios**:

- **Data Quality Improvement**: Replace missing data with sensible defaults (e.g., 0, "N/A", mean/median) to make the dataset more complete for analysis.
- **Preparing Data for Analysis/ML**: Many analytical tools and machine learning algorithms cannot handle null values. `fillna()` is a common pre-processing step.
- **Preventing Errors**: Null values can cause unexpected errors or incorrect results in calculations or joins.
- **Consistency**: Ensure consistent representation of missing data across your dataset.

**Best Practices**:

- **Column-Specific Fill Values**: Use a dictionary to specify different fill values for different columns, especially if columns have different data types (e.g., 0 for integers, "N/A" for strings, average for numerical columns).
- **Business Logic**: The choice of fill value should always be guided by business logic and the nature of the data. Replacing nulls with 0 might be appropriate for a `quantity` but not for an `age`.
- **Document Fill Strategies**: Clearly document how null values are handled, as this can significantly impact downstream analysis.
- **Validate After Filling**: After `fillna()`, run checks (e.g., `df.filter(col("column").isNull()).count()`) to ensure nulls were handled as expected.

# Scenario-Based Questions & Answers (Data Cleaning & Deduplication)

1. **Scenario**: You've ingested hero data from multiple sources, and sometimes the same hero appears with slightly different `first_appearance` dates, but their `name` and `universe` are always consistent. You want to keep only one record per unique (name, universe) pair, ignoring the `first_appearance` differences.
   - **Question**: How would you deduplicate the DataFrame to achieve this?
   - **Answer**:

   ```
   df_unique_heroes_by_name_universe = df.dropDuplicates(["name",
   "universe"])
   df_unique_heroes_by_name_universe.show()
   ```

2. **Scenario**: A new batch of hero data has arrived, and it's missing `debut_year` for some entries (represented as `None`). You want to replace these missing `debut_year` values with a placeholder year, `1900`, to avoid issues in numerical calculations.
   - **Question**: How would you fill these null `debut_year` values?
   - **Answer**:

   ```
   df_filled_debut_year = df.fillna({"debut_year": 1900})
   df_filled_debut_year.filter(col("name") == "Unknown
   Hero").show(truncate=False) # Assuming "Unknown Hero" might have
   null debut_year from a broader dataset
   ```

3. **Scenario**: You notice that some hero entries might have duplicate `powers` listed within their `powers` array (e.g., `["Flight", "Flight", "Super Strength"]`). Before further analysis, you want to ensure each hero's `powers` array contains only distinct powers.
   - **Question**: How would you remove duplicates within the `powers` array for each hero? (Note: This requires an array function, not `dropDuplicates` on the DataFrame itself).
   - **Answer**:

   ```
   from pyspark.sql.functions import array_distinct
   df_distinct_powers = df.withColumn("powers",
   array_distinct(col("powers")))
   df_distinct_powers.select("name", "powers").show(truncate=False)
   ```

4. **Scenario**: You have a `team` column which is an `ArrayType(StringType())`. If a hero has no team, the value might be `None` (null). You want to replace these null `team` entries with an empty array `[]` so that downstream array operations don't fail.

- **Question:** How would you use `fillna()` to replace `None` in the `team` column with an empty array?
- **Answer:**

```python
# Assuming df_nulls as created in the fillna example with a null
team
df_filled_team_nulls = df_nulls.fillna({"team": []})
df_filled_team_nulls.show(truncate=False)
```

5. **Scenario:** You've merged data from two sources, and now some heroes have duplicate rows, but you only want to keep the record with the earliest `first_appearance` date if there are duplicates based on `name`.

- **Question:** Explain how you would approach this using a combination of window functions and `dropDuplicates()` concepts, even though `dropDuplicates()` itself doesn't offer this control. (No code answer needed, just description).
- **Answer:** You would use a **Window Function** partitioned by `name` and ordered by `first_appearance` in ascending order. Then, you'd apply `row_number()` to assign a rank to each duplicate. Finally, you'd filter the DataFrame to keep only rows where `row_number` is 1, effectively selecting the earliest appearance. This is a common pattern for "deduplication with preference."

## 6. Sorting & Ordering

These functions arrange the rows of your DataFrame based on the values in one or more columns, enabling ordered analysis and presentation.

## orderBy() & sort()

**Purpose:** Both `orderBy()` and `sort()` are aliases for the same functionality: sorting the DataFrame rows by one or more columns. The default sort order is ascending.

**Syntax:** `df.orderBy(*cols)` or `df.sort(*cols)`

```python
# Sort by single column (ascending by default): Sorts heroes by their
debut year, earliest first
df_sorted = df.orderBy("debut_year")
df_sorted.select("name", "debut_year").show()

# Sort by multiple columns: Sorts by universe (ascending), then by debut
year (ascending)
df_multi_sorted = df.orderBy("universe", "debut_year")
df_multi_sorted.select("name", "universe", "debut_year").show()

# Sort with direction: Sorts by debut year in descending order
df_desc_sorted = df.orderBy(desc("debut_year"))
df_desc_sorted.select("name", "debut_year").show()
```

```
# Complex sorting: Sorts by universe (asc), then by debut year (desc),
then by name (asc)
df_complex_sorted = df.orderBy(
    asc("universe"),
    desc("debut_year"),
    asc("name")
)
df_complex_sorted.select("name", "universe", "debut_year").show()
```

**Explanation & Scenarios**:

- **Preparing Data for Presentation**: When displaying results, sorting ensures a logical and readable order, making it easier for users to consume the information.
- **Finding Top/Bottom Records**: After aggregation or calculation, sorting helps identify the highest or lowest values easily (e.g., top 10 best-selling products, oldest heroes).
- **Ensuring Consistent Ordering**: In some distributed operations or when writing data to files that rely on order, sorting can provide deterministic output.
- **Window Functions Pre-requisite**: Sorting is a fundamental part of defining a window for window functions (e.g., calculating running totals, ranks).

**Best Practices**:

- Use `asc()` and `desc()` explicitly for clear indication of sorting direction, even for ascending where it's the default.
- **Performance Impact**: Sorting involves a global shuffle of data across the cluster, which can be a very expensive operation for large datasets. Use it only when necessary and try to filter data before sorting if possible.
- For small results or debugging, sort before `show()` to get consistent output.

## Scenario-Based Questions & Answers (Sorting & Ordering)

1. **Scenario**: Your team wants to see a list of all heroes, sorted from the most recently debuted to the oldest.
   - **Question**: How would you sort the DataFrame to achieve this?
   - **Answer**:

   ```
   df_latest_debut_first = df.orderBy(desc("debut_year"))
   df_latest_debut_first.select("name", "debut_year").show()
   ```

2. **Scenario**: You need to organize the heroes first by their `universe` (alphabetically) and then, within each universe, by the `power_count` (number of powers) in descending order (heroes with more powers first).
   - **Question**: How would you sort the DataFrame for this complex requirement?

- **Answer:**

```
df_sorted_by_universe_and_powers = df.withColumn("power_count",
size(col("powers"))) \
                                        .orderBy(asc("universe"),
desc("power_count"))
df_sorted_by_universe_and_powers.select("name", "universe",
"powers", "power_count").show(truncate=False)
```

3. **Scenario**: You're preparing a report and need to list all heroes, but for heroes who debuted in the same year, you want them sorted alphabetically by their `name`.
   - **Question**: How would you sort to achieve this, prioritizing `debut_year` then `name`?
   - **Answer:**

```
df_report_sorted = df.orderBy(asc("debut_year"), asc("name"))
df_report_sorted.select("name", "debut_year").show()
```

4. **Scenario**: After performing a `groupBy` operation, you get a summary table and want to order the results by the aggregated count from highest to lowest.
   - **Question**: If you had a `universe_hero_counts` DataFrame from a `groupBy` and `count` operation, how would you sort it by the `hero_count` column in descending order?
   - **Answer:**

```
universe_hero_counts =
df.groupBy("universe").agg(count("name").alias("hero_count"))
sorted_universe_counts =
universe_hero_counts.orderBy(desc("hero_count"))
sorted_universe_counts.show()
```

5. **Scenario**: You have a very large DataFrame and need to check the first few rows to ensure a previous transformation worked correctly. You want to see these first rows in a consistent order every time you run the code, regardless of Spark's internal partitioning.
   - **Question**: How would you combine `orderBy()` and `show()` for this debugging purpose?
   - **Answer:**

```
df.orderBy("name").show(5) # Sort by a consistent column like 'name'
before showing
```

# 7. Aggregation & Grouping

These functions are used to summarize data, typically by grouping rows that share common characteristics and then applying aggregate functions (like sum, count, average) to each group.

## groupBy()

**Purpose**: Groups rows in the DataFrame based on one or more columns, allowing you to apply aggregate functions to each group. This is one of the most fundamental operations for data summarization and analytics.

**Syntax:** `df.groupBy(*cols).agg(*exprs)`

```python
# Basic grouping: Counts heroes, calculates average debut year, and finds
min/max debut year per universe
universe_stats = df.groupBy("universe").agg(
    count("name").alias("hero_count"),
    avg("debut_year").alias("avg_debut_year"),
    min("debut_year").alias("first_debut"),
    max("debut_year").alias("latest_debut")
)
universe_stats.show()

# Complex aggregations: Group by universe, then collect all hero names
into a list
# and calculate the average number of powers per hero in that universe.
power_stats = df.groupBy("universe").agg(
    count("name").alias("total_heroes"),
    collect_list("name").alias("hero_names"),
    avg(size(col("powers"))).alias("avg_power_count")
)
power_stats.show(truncate=False)

# Grouping by multiple columns
team_universe_stats = df.withColumn("single_team", explode(col("team"))) \
                        .groupBy("universe", "single_team") \
                        .agg(count("name").alias("heroes_in_team"))
team_universe_stats.show()
```

**Explanation & Scenarios**:

- **Summary Reports**: Generate reports like "total sales per region," "average customer age per demographic," or "number of products per category."

- **Data Aggregation for Dashboards**: Prepare aggregated metrics that can be easily visualized in business intelligence dashboards.

- **Feature Engineering for ML**: Create aggregated features (e.g., total purchases for a user, average rating for a product) that can be used in machine learning models.

- **Analyzing Group Characteristics**: Understand the characteristics of different groups within your data (e.g., average power count in Marvel vs. DC).

**Best Practices**:

- Use meaningful aliases for aggregated columns using `.alias()` for clarity.
- **Performance**: `groupBy()` is a wide transformation, meaning it involves shuffling data across the network (moving data from one partition to another). This can be expensive for very large datasets and many distinct groups.
- Choose appropriate aggregation functions based on your analytical needs (e.g., `sum`, `avg`, `min`, `max`, `count`, `collect_list`, `collect_set`).

---

## union() & unionByName()

**Purpose**: These functions combine two or more DataFrames vertically (appending rows from one DataFrame to another). They are essential for consolidating data from different sources or for combining partitioned data.

### union()

**Purpose**: Combines two DataFrames by position. The schemas of the DataFrames must be identical (same number of columns, same column names, and same data types in the same order). If schemas differ, it will result in an error.

```
# Create sample DataFrames
marvel_df = df.filter(col("universe") == "Marvel")
dc_df = df.filter(col("universe") == "DC")

# Union DataFrames: Combines all Marvel and DC heroes back into one
DataFrame
combined_df = marvel_df.union(dc_df)
combined_df.show(5)
print(f"Combined DataFrame count: {combined_df.count()}")
```

**Explanation**:

- Strict Schema Match: Requires precise column order and type matching.

**Best Practices**:

- **Ensure Identical Schemas**: Before using `union()`, verify that both DataFrames have the exact same schema. You might need to use `select()` and `withColumn()` to align columns.
- **Deduplication**: After a `union()`, you might need to use `dropDuplicates()` if the source DataFrames could contain overlapping records.

### unionByName()

**Purpose**: Combines two DataFrames by matching columns by name, rather than by position. This is more flexible as it allows for different column orders and can handle missing columns.

**Syntax:** `df1.unionByName(df2, allowMissingColumns=False)`

```python
# Create DataFrames with different column orders
df1_partial = df.select("name", "universe", "debut_year")
df2_partial = df.select("universe", "debut_year", "name") # Different
order

# Union by name (default allowMissingColumns=False, will work here as all
columns exist)
df_union_by_name = df1_partial.unionByName(df2_partial)
df_union_by_name.show(5)

# Allow missing columns: If one DataFrame has columns the other doesn't,
# the missing columns will be filled with nulls in the combined DataFrame.
df_only_name_universe = df.select("name", "universe")
df_union_missing = df1_partial.unionByName(df_only_name_universe,
allowMissingColumns=True)
df_union_missing.show(5)
```

**Explanation:**

- `allowMissingColumns=False` (default): All columns present in one DataFrame must also be present in the other.
- `allowMissingColumns=True` : If a column exists in one DataFrame but not the other, it will be filled with `null` values in the DataFrame where it was missing.

**Explanation & Scenarios**:

- **Consolidating Heterogeneous Data**: Combine datasets from different periods or sources that might have slightly different schemas (e.g., new columns added over time).
- **Appending Incremental Data**: Efficiently add new batches of data to an existing dataset where schema evolution might have occurred.
- **Schema Evolution**: When dealing with evolving schemas, `unionByName(allowMissingColumns=True)` is incredibly useful for gracefully handling new or deprecated columns.

**Best Practices:**

- Use `unionByName()` whenever column order is not guaranteed or when dealing with schema evolution.
- Carefully consider `allowMissingColumns` . If set to `True` , ensure that `null` values in the newly introduced columns are acceptable for downstream processing.
- Similar to `union()` , consider `dropDuplicates()` after `unionByName()` if overlap is possible.

## Scenario-Based Questions & Answers (Aggregation & Grouping)

1. **Scenario:** You want to find out the total number of heroes in each universe and the earliest debut year for a hero in that universe.

   - **Question:** How would you achieve this using `groupBy()` and appropriate aggregation functions?

   - **Answer:**

   ```python
   universe_summary = df.groupBy("universe").agg(
       count("name").alias("total_heroes"),
       min("debut_year").alias("earliest_debut_year")
   )
   universe_summary.show()
   ```

2. **Scenario:** Your data analytics team needs a list of all hero names for each `team` they belong to, along with the total count of heroes in that team. (Note: A hero can be in multiple teams, so you'll need to flatten the `team` array first).

   - **Question:** How would you use `explode()` and `groupBy()` to get this team-level aggregation?

   - **Answer:**

   ```python
   df_exploded_teams = df.withColumn("single_team",
   explode(col("team")))
   team_hero_counts = df_exploded_teams.groupBy("single_team").agg(
       count("name").alias("hero_count"),
       collect_list("name").alias("hero_names")
   )
   team_hero_counts.show(truncate=False)
   ```

3. **Scenario:** You have two DataFrames: `df_marvel_new_heroes` (new Marvel heroes) and `df_dc_new_heroes` (new DC heroes). Both have the same columns (`name`, `universe`, `debut_year`) in the same order. You want to combine them into a single DataFrame.

   - **Question:** Which `union` function would you use, and why?

   - **Answer:**

   ```python
   # Let's create dummy new hero DataFrames
   new_marvel_data = [("Captain Marvel", "Marvel", "Carol Danvers",
   1968, ["Flight"], {"strong": "yes"}, ["Avengers"], date(1968, 3,
   1))]
   new_dc_data = [("Cyborg", "DC", "Victor Stone", 1980, ["Tech
   Manipulation"], {"cybernetic": "yes"}, ["Justice League"],
   date(1980, 7, 1))]
   df_marvel_new_heroes = spark.createDataFrame(new_marvel_data,
   schema) # Reusing schema for simplicity
   df_dc_new_heroes = spark.createDataFrame(new_dc_data, schema)
   ```

```
# Since schemas are identical, union() is sufficient and slightly
more performant.
combined_new_heroes = df_marvel_new_heroes.union(df_dc_new_heroes)
combined_new_heroes.show()
```

**Explanation**: You'd use `union()` because the schemas are guaranteed to be identical in terms of column names, order, and types. `unionByName()` would also work, but `union()` is slightly more efficient when strict schema matching is met.

4. **Scenario**: You're combining historical data with current data. The historical data (`df_old_heroes`) might have a `city_of_origin` column that isn't present in the `df` DataFrame (current data), and the column order might be different. You want to merge them, keeping all columns, filling missing ones with `null`.

- **Question**: Which `union` function is appropriate here, and what parameter would you set?

- **Answer**:

```
# Dummy old heroes DataFrame with a new column and different order
old_data = [
    ("Ghost Rider", "Marvel", 1972, "Hell"),
    ("Plastic Man", "DC", 1941, "Gotham")
]
old_schema = StructType([
    StructField("name", StringType(), True),
    StructField("universe", StringType(), True),
    StructField("debut_year", IntegerType(), True),
    StructField("city_of_origin", StringType(), True)
])
df_old_heroes = spark.createDataFrame(old_data, old_schema)

# To combine, use unionByName with allowMissingColumns=True
# Need to select matching columns in df and fill others with null
# For a practical union, ensure matching core columns or select a
subset
df_subset_for_union = df.select("name", "universe", "debut_year") #
Example, matching some columns for demo
combined_heroes_flexible =
df_subset_for_union.unionByName(df_old_heroes,
allowMissingColumns=True)
combined_heroes_flexible.show()
```

**Explanation**: You would use `unionByName(allowMissingColumns=True)`. This allows DataFrames with differing schemas and column orders to be combined. Columns present in one but not the other will be filled with `null` values in the resulting DataFrame where they were missing.

5. **Scenario:** You want to calculate the average number of `powers` per hero for each `debut_year`.

- **Question:** How would you group the DataFrame and perform this aggregation?
- **Answer:**

```python
avg_powers_per_year = df.groupBy("debut_year").agg(
    avg(size(col("powers"))).alias("average_powers")
)
avg_powers_per_year.orderBy("debut_year").show()
```

# 8. Data Joins & Relationships

Joining DataFrames is a cornerstone of data processing in Spark, allowing you to combine data from different sources based on common columns.

## join()

**Purpose:** Combines two DataFrames horizontally based on one or more common columns (join keys). This is equivalent to SQL JOIN operations.

**Syntax:** `df1.join(df2, on=None, how='inner')`

```python
# Create a powers reference DataFrame (example: mapping powers to types)
powers_ref = spark.createDataFrame([
    ("Super Strength", "Physical"),
    ("Flight", "Physical"),
    ("Magic", "Mystic"),
    ("Genius Intellect", "Mental"),
    ("Time Manipulation", "Mystic"),
    ("Wall-Crawling", "Physical")
], ["power_name", "power_type"])
powers_ref.show()

# Explode powers for joining: Each hero's power becomes a separate row
# This allows joining on individual power names.
df_powers_exploded = df.select("name", "universe",
explode("powers").alias("power_name"))
df_powers_exploded.show(5)

# Inner join: Combines heroes with their power types, only keeping records
# where a match is found in both DFs
df_with_power_types_inner = df_powers_exploded.join(powers_ref,
on="power_name", how="inner")
df_with_power_types_inner.show()

# Left join (left_outer): Keeps all records from the left DataFrame
(df_powers_exploded)
```

```python
# and matches with powers_ref. If no match, power_type will be null.
df_left_join = df_powers_exploded.join(powers_ref, on="power_name",
how="left")
df_left_join.show()

# Right join (right_outer): Keeps all records from the right DataFrame
(powers_ref)
# and matches with df_powers_exploded. If no match, hero details will be
null.
df_right_join = df_powers_exploded.join(powers_ref, on="power_name",
how="right")
df_right_join.show()

# Full outer join: Keeps all records from both DataFrames, matching where
possible.
# Nulls will appear for unmatched records on either side.
df_full_outer_join = df_powers_exploded.join(powers_ref, on="power_name",
how="full_outer")
df_full_outer_join.show()

# Join with different column names: Explicitly specify join conditions
# Drop the duplicate power_name column from the right side after join
df_join_different = df_powers_exploded.join(
    powers_ref,
    df_powers_exploded.power_name == powers_ref.power_name,
    "inner"
).drop(powers_ref.power_name) # Drop the redundant column after join
df_join_different.show()

# Self-join scenario: Find heroes from the same universe who debuted in
the same year
df_self_join = df.alias("hero1").join(
    df.alias("hero2"),
    (col("hero1.universe") == col("hero2.universe")) &
    (col("hero1.debut_year") == col("hero2.debut_year")) &
    (col("hero1.name") ≠ col("hero2.name")), # Exclude joining a hero
with themselves
    "inner"
).select(
    col("hero1.name").alias("hero1_name"),
    col("hero2.name").alias("hero2_name"),
    col("hero1.universe"),
    col("hero1.debut_year")
)
df_self_join.show()
```

**Join Types:**

- `inner` : Returns only the rows where the join condition is met in *both* DataFrames. This is the default.
- `left_outer` (or `left` ): Returns all rows from the left DataFrame, and the matched rows from the right DataFrame. If no match, columns from the right DataFrame will be `null` .
- `right_outer` (or `right` ): Returns all rows from the right DataFrame, and the matched rows from the left DataFrame. If no match, columns from the left DataFrame will be `null` .
- `full_outer` (or `full` ): Returns all rows when there is a match in one of the DataFrames. If no match, columns from the unmatched side will be `null` .
- `cross` : Returns the Cartesian product of rows from both DataFrames. Each row from `df1` is combined with every row from `df2` . **Use with extreme caution on large datasets**, as it can result in a massive DataFrame. Implicit `crossJoin` is often disallowed for safety, requiring explicit `crossJoin()` method call.
- `left_semi` : Returns rows from the left DataFrame that have a match in the right DataFrame. Columns from the right DataFrame are not included. Similar to `WHERE EXISTS` in SQL.
- `left_anti` : Returns rows from the left DataFrame that do *not* have a match in the right DataFrame. Columns from the right DataFrame are not included. Similar to `WHERE NOT EXISTS` in SQL.

**Explanation & Scenarios**:

- **Data Enrichment**: Add information from a lookup table or a dimension table to a fact table (e.g., add customer demographics to sales transactions).
- **Relational Database Equivalence**: Mimic common SQL JOIN operations to combine data that is logically related but stored in separate DataFrames.
- **Fact and Dimension Tables**: A common use case in data warehousing, joining fact tables with dimension tables to provide context.
- **Merging Data for Analysis**: Combine various datasets (e.g., hero details with their comic appearances, or attributes) for a comprehensive view.

**Best Practices**:

- **Choose Appropriate Join Type**: Carefully select the join type ( `inner` , `left` , `right` , `full` , `semi` , `anti` , `cross` ) based on your desired output and how you want to handle unmatched rows.
- **Broadcast Joins**: For joining a large DataFrame with a small DataFrame (typically less than `10GB` ), use `broadcast()` the smaller DataFrame. Spark will send the smaller DataFrame to all worker nodes, avoiding a shuffle of the large DataFrame and significantly improving performance.

```
from pyspark.sql.functions import broadcast
df_with_power_types_broadcast =
df_powers_exploded.join(broadcast(powers_ref), on="power_name",
how="inner")
```

- **Handle Duplicate Column Names**: After a join, if both DataFrames have columns with the same name (and they are not join keys), Spark will keep both, differentiating them by their origin (e.g., `df1.column_name`, `df2.column_name`). Use `drop()` or `select()` to manage these. When using explicit join conditions (`df1.col_a == df2.col_b`), you often need to `drop()` one of the join key columns if they are identical.
- **Performance Implications**: Joins, especially wide transformations like `inner`, `left`, `right`, `full_outer`, involve significant data shuffling across the network. Optimize join keys and consider data partitioning. `left_semi` and `left_anti` are generally more efficient as they only need to check for existence, not retrieve data from the right side.

## Scenario-Based Questions & Answers (Data Joins & Relationships)

1. **Scenario**: You have a separate DataFrame, `villain_data`, with `villain_name`, `nemesis_hero_name`, and `threat_level`. You want to find out the `threat_level` of nemeses for each superhero, joining on the superhero's `name`. Only show heroes that *have* a nemesis in the `villain_data`.
   - **Question**: How would you perform this join?
   - **Answer**:

   ```
   villain_data = spark.createDataFrame([
       ("Joker", "Batman", "High"),
       ("Lex Luthor", "Superman", "High"),
       ("Green Goblin", "Spider-Man", "Medium"),
       ("Doctor Doom", "Fantastic Four", "High") # No direct match in
   df for Fantastic Four hero
   ], ["villain_name", "nemesis_hero_name", "threat_level"])

   df_heroes_with_threat = df.join(
       villain_data,
       df.name == villain_data.nemesis_hero_name,
       "inner"
   ).select(df.name, df.universe, villain_data.villain_name,
   villain_data.threat_level)
   df_heroes_with_threat.show()
   ```

2. **Scenario**: You want to list all superheroes and, if available, their associated `power_type` from the `powers_ref` DataFrame. Heroes without a defined `power_type` should still be included, with `null` in the `power_type` column.
   - **Question**: Which join type would you use for this, and how would you implement it?
   - **Answer**:

```
# Re-using df_powers_exploded and powers_ref from previous examples
df_left_join_all_heroes_powers = df_powers_exploded.join(powers_ref,
on="power_name", how="left_outer")
df_left_join_all_heroes_powers.show()
```

3. **Scenario:** You're specifically interested in the `power_type` of each unique power listed in the `powers_ref` DataFrame. You want to see which heroes (if any) possess each `power_type`, but also include power types that no current hero has (like "Telepathy", if it existed in `powers_ref` but not in `df`).

   - **Question:** Which join type would you use to achieve this, ensuring all `power_type` values from `powers_ref` are present?

   - **Answer:**

```
# Assuming we add "Telepathy" to powers_ref for demonstration
powers_ref_extended = spark.createDataFrame([
    ("Super Strength", "Physical"),
    ("Flight", "Physical"),
    ("Magic", "Mystic"),
    ("Genius Intellect", "Mental"),
    ("Time Manipulation", "Mystic"),
    ("Wall-Crawling", "Physical"),
    ("Telepathy", "Mental") # A power type not in current hero data
], ["power_name", "power_type"])

df_right_join_all_power_types =
df_powers_exploded.join(powers_ref_extended, on="power_name",
how="right_outer")
df_right_join_all_power_types.show()
```

4. **Scenario:** You need to find pairs of heroes from *different* universes who debuted in the same year.

   - **Question:** How would you use a self-join to identify such pairs?

   - **Answer:**

```
df_cross_universe_debut = df.alias("hero1").join(
    df.alias("hero2"),
    (col("hero1.debut_year") == col("hero2.debut_year")) &
    (col("hero1.universe") ≠ col("hero2.universe")), # Different
universes
    "inner"
).select(
    col("hero1.name").alias("hero1_name"),
    col("hero1.universe").alias("hero1_universe"),
    col("hero2.name").alias("hero2_name"),
```

```
        col("hero2.universe").alias("hero2_universe"),
        col("hero1.debut_year")
    )
    df_cross_universe_debut.show()
```

5. **Scenario**: You have a small DataFrame `special_powers_list` with a list of unique, critical powers (e.g., "Magic", "Time Travel"). You want to identify *only* the heroes from your `df` DataFrame who possess any of these `special_powers`. You don't need any additional details from `special_powers_list`, just the heroes who have them.

   - **Question**: Which efficient join type would you use for this "existence check"?
   - **Answer**:

```
special_powers_list = spark.createDataFrame([
    ("Magic",),
    ("Time Travel",)
], ["power_name"])

# First, explode the powers of the heroes to make them joinable
df_hero_powers_exploded = df.select("name",
explode("powers").alias("power_name"))

# Use left_semi join to get heroes that have any of the special
powers
heroes_with_special_powers = df_hero_powers_exploded.join(
    special_powers_list,
    on="power_name",
    how="left_semi"
).distinct() # Use distinct to get unique hero names
heroes_with_special_powers.show()
```

# 9. Data Iteration & Processing

While Spark's strength lies in parallel, set-based operations, sometimes you need to process data row by row, or perform actions that have side effects.

## foreach()

**Purpose**: Applies a function to each row of the DataFrame. It's primarily used for operations that have *side effects*, such as writing to an external database, sending data to a message queue, or printing to the console (for debugging small datasets). **It does not return a new DataFrame.**

**Syntax**: `df.foreach(func)`

```
# Define a function to process each row
def print_hero_details(row):
```

```python
    # This print statement will execute on the worker nodes.
    # Its output might not appear in the driver's console immediately or
in order.
    print(f"Processing Hero: {row.name} from {row.universe}, debuted
{row.debut_year}")

# Apply to a small sample for demonstration (not recommended for large
datasets due to print overhead)
print("\n--- Using foreach() on a limited DataFrame (output might be
scattered) ---")
df.limit(3).foreach(print_hero_details)

# Better approach for local iteration/debugging: collect and iterate
print("\n--- Using collect() and iterating locally (recommended for small
results) ---")
for row in df.limit(3).collect():
    print(f"Hero (collected): {row.name} from {row.universe}, debuted
{row.debut_year}")

# Scenario: Write each hero's name to a log file (illustrative)
# This would involve actual file I/O on workers.
# Example:
# from datetime import datetime
# def log_hero_to_file(row):
#     with
open(f"/tmp/hero_log_{datetime.now().strftime('%Y%m%d%H')}.log", "a") as
f:
#         f.write(f"Logged: {row.name} - {row.universe}\n")
#
# df.foreach(log_hero_file) # Will write one line per hero to a
distributed log file
```

**Explanation & Scenarios**:

- **External System Integration**: Sending individual records to external systems (e.g., writing to a NoSQL database row by row, pushing events to Kafka, calling a REST API for each record).

- **Logging and Monitoring**: Performing custom logging logic for each record, though usually, Spark's built-in logging mechanisms are preferred.

- **Debugging Small Datasets**: For very small datasets, `foreach` can be used with print statements for quick inspection, but be aware of its distributed nature (output might be interleaved or appear on worker logs).

**Best Practices**:

- **Use Sparingly**: Avoid `foreach()` for transformations or aggregations. These should always be done with DataFrame operations for performance and scalability.

- **Performance Implications**: `foreach()` can be very slow if the action performed for each row is expensive (e.g., frequent I/O to a non-distributed sink) or if it involves a lot of

communication back to the driver.

- **No Return Value**: Remember that `foreach()` does not return a DataFrame. It is purely for side effects.
- **Consider `collect()` for Local Processing**: If you truly need to process data row by row in Python on the driver (e.g., for simple display or very small lookups), `collect()` the data first and then iterate over the resulting Python list. This is safer than `foreach` for local iteration.
- **Distributed Logging**: If logging each record to a file, ensure the file system is accessible to all workers and consider potential file contention issues.

## Scenario-Based Questions & Answers (Data Iteration & Processing)

1. **Scenario**: After a complex set of transformations, you want to perform a final data validation step where for each hero, you check if their `debut_year` is greater than their `first_appearance` year. If it is, you want to log a warning message for that hero.

   - **Question**: How would you use `foreach()` to perform this check and print a warning for each problematic hero?

   - **Answer**:

   ```python
   def validate_debut_year(row):
       if row.debut_year is not None and row.first_appearance is not None:
           if row.debut_year > row.first_appearance.year:
               print(f"WARNING: Hero {row.name} has a debut_year ({row.debut_year}) later than first_appearance year ({row.first_appearance.year})")

   df.foreach(validate_debut_year)
   ```

2. **Scenario**: You need to send hero names and their universes to an external API endpoint one by one. You want to simulate this process for the first 5 heroes.

   - **Question**: How would you use `collect()` to safely retrieve these 5 heroes and then iterate through them to simulate the API call?

   - **Answer**:

   ```python
   def send_to_api(hero_name, hero_universe):
       print(f"Simulating API call: Sending {hero_name} from {hero_universe}")
       # In a real scenario, this would be an actual API call
       # import requests
       # requests.post("your_api_endpoint", json={"name": hero_name,
   ```

```
    "universe": hero_universe})

    for row in df.limit(5).collect():
        send_to_api(row.name, row.universe)
```

3. **Scenario**: You have a small subset of heroes and need to perform a very specific, computationally intensive Python operation on each hero's `powers` array that isn't easily expressed with Spark's built-in functions or even simple UDFs (e.g., complex graph traversal on a subset of powers). You don't need to return a new DataFrame, just perform a side effect.

   - **Question**: While acknowledging performance implications, how would you set up `foreach()` to apply a Python function that performs this complex operation on the `powers` list of each hero?

   - **Answer**:

   ```
   def complex_power_analysis(row):
       hero_name = row.name
       powers = row.powers
       # Imagine a very complex graph traversal or external call here
       # For demonstration, just print
       print(f"Performing complex analysis for {hero_name} with powers:
   {powers}")
       # result = some_complex_function(powers)
       # print(f"Analysis result for {hero_name}: {result}")

   # Limit to a very small subset for such intensive operations
   df.limit(2).foreach(complex_power_analysis)
   ```

4. **Scenario**: You are running a Spark job on a cluster, and you notice that `print()` statements within a `foreach` are not showing up consistently in your driver logs.
   - **Question**: Explain why this might be happening and what the recommended approach for debugging is if you need to see intermediate row data. (No code required).
   - **Answer**: Print statements within `foreach` are executed on the **worker nodes**, not the driver. Their output often goes to the worker logs, which might not be directly streamed to your console or central driver logs. The recommended approach for debugging intermediate row data is to **filter the DataFrame down to a very small subset** (e.g., `df.filter( ... ).limit(10)`), then `collect()` that small subset to the driver, and then iterate and print those rows locally. This ensures all output is visible on the driver.

5. **Scenario**: You have successfully transformed a large DataFrame and want to save each record into a separate file in a specific directory on a distributed file system, with the filename derived from the hero's name.

- **Question**: While this is better handled by `df.write`, how would `foreach()` *theoretically* enable this for a single hero for demonstration purposes (ignoring best practices for distributed writes)?
  - **Answer**:

```
# This is a hypothetical example and not recommended for distributed
writing.
# Spark's df.write.format("json").save() or similar is the correct
way.
# This is just to illustrate foreach for a side effect (file
writing).

# from pyspark.sql.functions import concat_ws, lit
#
# def save_hero_to_file(row):
#     file_name = row.name.replace(" ", "_").lower() + ".json"
#     # In a real cluster, this would write to HDFS/S3 etc.
#     # This simple open() is for local demonstration
#     with open(f"/tmp/heroes_individual/{file_name}", "w") as f:
#         f.write(row.json()) # Converts Row object to JSON string
#
# # Create a directory first: os.makedirs("/tmp/heroes_individual",
exist_ok=True)
# df.limit(1).foreach(save_hero_to_file) # Process just one hero
```

## 10. Data Partitioning

Data partitioning is a crucial concept in Spark for optimizing performance. It determines how data is physically distributed across the cluster's nodes.

## `partitionBy()` (for writing data)

**Purpose**: When writing a DataFrame to disk (e.g., Parquet, ORC), `partitionBy()` physically organizes the data into subdirectories based on the values of specified columns. This greatly improves read performance for queries that filter on these partition columns.

**Syntax**: `df.write.partitionBy(*cols).format("parquet").save(path)`

```
# Scenario: Save heroes data, partitioned by universe.
# This creates directories like 'heroes_partitioned/universe=Marvel/',
'heroes_partitioned/universe=DC/'
print("Saving partitioned data by universe ... ")
df.write.partitionBy("universe").mode("overwrite").parquet("heroes_partiti
oned")
print("Data partitioned and saved to 'heroes_partitioned' directory.")
# To verify (not executable in current context but conceptually):
# !ls -R heroes_partitioned/
```

```python
# Multiple partition columns: Partition by universe, then by debut decade
print("\nSaving data partitioned by universe and decade ... ")
df.withColumn("decade", floor(col("debut_year") / 10) * 10) \
  .write.partitionBy("universe", "decade") \
  .mode("overwrite") \
  .parquet("heroes_multi_partitioned")
print("Data partitioned and saved to 'heroes_multi_partitioned' directory.")
# !ls -R heroes_multi_partitioned/
```

**Explanation & Scenarios**:

- **Query Optimization**: When you query a partitioned table and filter on the partition columns (e.g., `SELECT * FROM heroes WHERE universe = 'Marvel'`), Spark can perform "partition pruning." It only reads the relevant subdirectories, skipping unnecessary data, leading to significantly faster queries.
- **Data Organization**: Provides a logical and physical organization of your data on storage (e.g., HDFS, S3).
- **Performance for Filtered Reads**: Ideal for very large datasets where specific subsets are frequently queried.

**Best Practices**:

- **Choose Partition Columns Wisely**: Select columns with a reasonable cardinality (number of unique values).
  - **Too few partitions (low cardinality)**: Leads to large, unmanageable files within partitions and less parallelism.
  - **Too many partitions (high cardinality)**: Creates an excessive number of small files (the "small file problem"), which can lead to inefficient read performance and overhead in distributed file systems.
- **Consider Query Patterns**: Partition by columns that are frequently used in `WHERE` clauses for filtering.
- The partition columns will not appear as regular columns in the DataFrame when read back; instead, their values are inferred from the directory structure.

## `repartition()` & `coalesce()`

**Purpose:** These functions control the number of partitions of a DataFrame in memory. This influences the parallelism of subsequent operations and can be critical for performance tuning.

```python
# Check initial number of partitions (might vary based on Spark config/data size)
print(f"Initial number of partitions: {df.rdd.getNumPartitions()}")

# Repartition to specific number of partitions: Increases parallelism or
```

```python
balances data
df_repartitioned = df.repartition(4) # Recreate DataFrame with 4
partitions
print(f"Number of partitions after repartition(4):
{df_repartitioned.rdd.getNumPartitions()}")

# Repartition by column: Distributes data based on hash of the column(s)
# Ensures rows with the same universe value go to the same partition
df_repartitioned_by_col = df.repartition(col("universe"))
print(f"Number of partitions after repartition by 'universe':
{df_repartitioned_by_col.rdd.getNumPartitions()}")

# Coalesce to reduce partitions: Efficiently reduces partitions without a
full shuffle if possible
df_coalesced = df_repartitioned.coalesce(2) # Reduces from 4 partitions to
2
print(f"Number of partitions after coalesce(2) from 4 partitions:
{df_coalesced.rdd.getNumPartitions()}")
```

**Explanation**:

- `repartition(numPartitions)` : Creates exactly `numPartitions` partitions. It always involves a **full shuffle** of data across the network, even if reducing the number of partitions. Useful for increasing parallelism if current partitions are too few, or for balancing skewed data.

- `repartition(*cols)` : Repartitions the DataFrame by hash-partitioning based on the specified columns. Rows with the same values in these columns will end up in the same partition. This is often used before joins or aggregations on those columns to minimize shuffle.

- `coalesce(numPartitions)` : Reduces the number of partitions to `numPartitions` . It tries to avoid a full shuffle by moving data only when necessary. It can only **decrease** the number of partitions. It's more efficient than `repartition()` for reducing partitions, as it minimizes data movement.

**Explanation & Scenarios**:

- **Performance Tuning**: Optimizing the number of partitions is crucial for Spark performance.
    - **Too few partitions**: Underutilizes cluster resources, leading to bottlenecks.
    - **Too many partitions (especially tiny ones)**: Introduces overhead in task scheduling and management, leading to the "small task problem."

- **Preparing for Aggregations/Joins**: Repartitioning by join keys or group-by keys before a `join()` or `groupBy()` operation can significantly reduce the shuffle cost during those operations.

- **Before Writing Data**: Coalescing partitions before writing to a distributed file system can reduce the number of output files, preventing the "small file problem" (which affects read performance and metadata management).

**Best Practices:**

- Use `repartition()` to **increase** or set a specific number of partitions (always involves a shuffle).
- Use `coalesce()` to **efficiently reduce** the number of partitions (minimizes shuffle).
- The optimal number of partitions often depends on your cluster size, data size, and the number of cores available. A common heuristic is 2-4 partitions per CPU core.
- Monitor Spark UI to identify skewed partitions or bottlenecks related to partitioning.

## `MapType` (Map/Dict Operations)

**Purpose:** PySpark provides functions to interact with `MapType` columns (which are equivalent to dictionaries or hash maps), allowing you to extract keys, values, or specific entries.

```python
# Extract map keys and values: Convert the 'attributes' map into an array
of keys and an array of values
df_with_attr_keys = df.withColumn("attribute_keys",
map_keys(col("attributes")))
df_with_attr_values = df.withColumn("attribute_values",
map_values(col("attributes")))
df_with_attr_keys.select("name", "attributes",
"attribute_keys").show(truncate=False)
df_with_attr_values.select("name", "attributes",
"attribute_values").show(truncate=False)

# Access specific map values: Get the 'intelligence' level from the
'attributes' map
df_with_intelligence = df.withColumn("intelligence_level",
    col("attributes").getItem("intelligence")) # Or col("attributes")
["intelligence"]
df_with_intelligence.select("name", "attributes",
"intelligence_level").show(truncate=False)

# Create map from arrays: Construct a new map column from two arrays (keys
and values)
df_with_new_map = df.withColumn("new_map",
    map_from_arrays(array(lit("key1"), lit("key2")),
                    array(lit("value1"), lit("value2"))))
df_with_new_map.select("name", "new_map").show(truncate=False)

# Filtering based on map key existence
df_filtered_by_attr_key =
df.filter(map_keys(col("attributes")).contains("wealth"))
df_filtered_by_attr_key.select("name", "attributes").show(truncate=False)
```

**Explanation & Scenarios**:

- **Parsing Semi-structured Data**: When you have JSON-like data or flexible key-value pairs stored in a `MapType` column, these functions are essential for extracting and manipulating the individual elements.
- **Feature Extraction**: Extract specific attributes (e.g., `intelligence`, `wealth`) from a generic `attributes` map column to create dedicated features.
- **Data Transformation**: Transform existing lists into a map, or vice versa, for easier consumption by downstream applications.
- **Filtering**: Filter rows based on the existence of certain keys or values within the map.

## Scenario-Based Questions & Answers (Data Partitioning & MapType Operations)

1. **Scenario**: You have a very large DataFrame of hero events and users frequently query this data by `universe` and then by `debut_year`. To optimize query performance when saving this data, you want to partition it by these two columns.

   - **Question**: How would you save the `df` DataFrame as Parquet files, partitioned by `universe` and `debut_year`?
   - **Answer**:

   ```
   # Assuming you want to add the decade as a partition, otherwise just
   use debut_year
   # df.write.partitionBy("universe",
   "debut_year").mode("overwrite").parquet("heroes_events_partitioned")
   # For demonstration, using existing df
   df.write.partitionBy("universe",
   "debut_year").mode("overwrite").parquet("heroes_events_by_universe_y
   ear")
   print("Hero events data saved, partitioned by universe and
   debut_year.")
   ```

2. **Scenario**: Your current DataFrame has 200 partitions, but your cluster has only 10 cores available. You want to reduce the number of partitions to 10 to better match your cluster's parallelism and avoid excessive task overhead without triggering a full shuffle if possible.

   - **Question**: Which function would you use to efficiently reduce the number of partitions to 10?
   - **Answer**:

   ```
   # First, let's artificially repartition df to 200 for this scenario
   df_large_partitions = df.repartition(200)
   print(f"DataFrame artificially repartitioned to
   {df_large_partitions.rdd.getNumPartitions()} partitions.")
   ```

```
df_optimized_partitions = df_large_partitions.coalesce(10)
print(f"DataFrame coalesced to
{df_optimized_partitions.rdd.getNumPartitions()} partitions.")
```

3. **Scenario**: You want to extract the `morality` attribute from the `attributes` map for each hero. If a hero doesn't have a `morality` attribute, it should simply appear as `null`.

   - **Question**: How would you add a new column `morality_level` that contains this information?
   - **Answer**:

```
df_with_morality = df.withColumn("morality_level",
col("attributes").getItem("morality"))
df_with_morality.select("name", "attributes",
"morality_level").show(truncate=False)
```

4. **Scenario**: You need to create a new `MapType` column called `contact_info` for each hero, where the keys are "email" and "phone", and the values are "hero_email@example.com" and "555-1234" respectively.

   - **Question**: How would you create this `MapType` column using PySpark functions?
   - **Answer**:

```
# from pyspark.sql.functions import map_from_arrays, array, lit
df_with_contact = df.withColumn("contact_info",
    map_from_arrays(array(lit("email"), lit("phone")),
                    array(lit(col("name") + "@example.com"),
lit("555-1234"))))
df_with_contact.select("name", "contact_info").show(truncate=False)
```

5. **Scenario**: You want to filter the DataFrame to only include heroes who have the attribute "wealth" defined in their `attributes` map, regardless of its value.

   - **Question**: How would you filter based on the existence of a key in a `MapType` column?
   - **Answer**:

```
df_wealthy_heroes =
df.filter(map_keys(col("attributes")).contains("wealth"))
df_wealthy_heroes.select("name", "attributes").show(truncate=False)
```

## 11. Additional Essential Functions

PySpark's DataFrame API is rich with functions for various data manipulation tasks. Here are some commonly used categories.

# Window Functions

**Purpose:** Perform calculations across a set of DataFrame rows that are related to the current row. Unlike aggregate functions that reduce the number of rows, window functions return a value for each row. They are powerful for calculating moving averages, rankings, and lead/lag values.

```python
# Define window specification: Partition by 'universe' and order by
'debut_year'
window_spec = Window.partitionBy("universe").orderBy("debut_year")

# Add row numbers within each universe: Assigns a unique, sequential
number within each universe
df_with_row_num = df.withColumn("row_num", row_number().over(window_spec))
df_with_row_num.select("name", "universe", "debut_year",
"row_num").orderBy("universe", "debut_year").show()

# Add ranking: Assigns a rank based on debut_year within each universe
(skips ranks for ties)
df_with_rank = df.withColumn("debut_rank", rank().over(window_spec))
df_with_rank.select("name", "universe", "debut_year",
"debut_rank").orderBy("universe", "debut_year").show()

# Add dense ranking: Assigns a rank without skipping numbers for ties
df_with_dense_rank = df.withColumn("debut_dense_rank",
dense_rank().over(window_spec))
df_with_dense_rank.select("name", "universe", "debut_year",
"debut_dense_rank").orderBy("universe", "debut_year").show()

# Add lag/lead values: Get the debut_year of the previous/next hero in the
same universe by debut_year order
df_with_lag = df.withColumn("prev_debut_year", lag("debut_year",
1).over(window_spec))
df_with_lag.select("name", "universe", "debut_year",
"prev_debut_year").orderBy("universe", "debut_year").show()

# Calculate running sum of debut years within each universe
window_sum_spec =
Window.partitionBy("universe").orderBy("debut_year").rowsBetween(Window.un
boundedPreceding, Window.currentRow)
df_with_running_sum = df.withColumn("running_debut_sum",
sum("debut_year").over(window_sum_spec))
df_with_running_sum.select("name", "universe", "debut_year",
"running_debut_sum").orderBy("universe", "debut_year").show()
```

**Scenarios:**

- **Ranking**: Finding the top N items within a group (e.g., top 5 performing products per region).
- **Running Totals/Averages**: Calculating cumulative sums, moving averages, or running counts (e.g., cumulative sales over time).
- **Lead/Lag Analysis**: Comparing a row's value to the value of a preceding or succeeding row (e.g., difference in sales from previous month).
- **Deduplication with Preference**: When used with `row_number()`, allows selection of a "golden record" among duplicates based on specific criteria.

## String Functions

**Purpose**: Manipulate and transform string columns. PySpark provides a rich set of string functions, mirroring those found in SQL.

```python
# String operations
df_string_ops = df.withColumn("name_upper", upper(col("name"))) \
                  .withColumn("name_lower", lower(col("name"))) \
                  .withColumn("name_length", length(col("name"))) \
                  .withColumn("alias_trimmed", trim(col("alias"))) \
                  .withColumn("universe_replaced",
regexp_replace(col("universe"), "Marvel", "MARVEL_COMICS"))
df_string_ops.select("name", "name_upper", "name_lower", "name_length",
"alias", "alias_trimmed", "universe", "universe_replaced").show(5,
truncate=False)

# Scenario: Extracting a substring (e.g., first 5 characters of name)
from pyspark.sql.functions import substring
df_substring = df.withColumn("name_prefix", substring(col("name"), 1, 5))
df_substring.select("name", "name_prefix").show()

# Scenario: Concatenating strings
from pyspark.sql.functions import concat_ws
df_full_title = df.withColumn("full_title", concat_ws(" (", col("name"),
col("alias"), lit(")"))).drop("alias")
df_full_title.select("name", "full_title").show(truncate=False)
```

**Scenarios**:

- **Data Cleaning**: Removing whitespace, changing case, or replacing characters.
- **Pattern Matching**: Using regular expressions to extract or replace patterns.
- **Feature Engineering**: Creating new features from text data, like string length or specific substrings.
- **Data Masking**: Replacing sensitive parts of a string.

## Array Functions

**Purpose**: Operate on columns of `ArrayType`, allowing you to inspect, transform, or filter based on array elements.

```python
# Array operations
df_array_ops = df.withColumn("has_super_strength",
array_contains(col("powers"), "Super Strength")) \
                .withColumn("unique_powers",
array_distinct(col("powers"))) \
                .withColumn("power_count", size(col("powers")))

df_array_ops.select("name", "powers", "has_super_strength",
"unique_powers", "power_count").show(truncate=False)

# Scenario: Filtering heroes who have AT LEAST TWO teams
from pyspark.sql.functions import array_size
df_multi_team_heroes = df.filter(array_size(col("team")) ≥ 2)
df_multi_team_heroes.select("name", "team").show(truncate=False)

# Scenario: Checking if an array is empty
from pyspark.sql.functions import is_empty
df_no_powers_heroes = df.withColumn("no_powers", is_empty(col("powers")))
df_no_powers_heroes.select("name", "powers",
"no_powers").show(truncate=False)

# Scenario: Intersecting arrays (find common powers between two heroes,
hypothetical)
from pyspark.sql.functions import array_intersect
df_intersect_powers = df.alias("hero1").join(
    df.alias("hero2"),
    (col("hero1.name") == "Spider-Man") & (col("hero2.name") == "Hulk"),
    "cross" # Only one pair needed for demonstration
).select(
    col("hero1.name").alias("hero1_name"),
    col("hero2.name").alias("hero2_name"),
    array_intersect(col("hero1.powers"),
col("hero2.powers")).alias("common_powers")
)
df_intersect_powers.show(truncate=False)
```

**Scenarios**:

- **Data Filtering**: Identifying records where an array contains a specific element or has a certain size.
- **Data Cleaning**: Removing duplicate elements within an array.
- **Feature Engineering**: Creating new features based on array properties (e.g., number of

elements, existence of specific elements).

- **Set Operations**: Performing set-like operations on arrays (intersection, union, difference).

---

## Date/Time Functions

**Purpose:** Manipulate and extract information from `DateType` and `TimestampType` columns.

```python
# Date operations
# df_date_ops = df.withColumn("debut_formatted",
date_format(col("first_appearance"), "yyyy-MM-dd")) \
#                 .withColumn("days_since_debut", datediff(current_date(),
col("first_appearance"))) \
#                 .withColumn("month_of_debut",
month(col("first_appearance"))) \
#                 .withColumn("year_of_debut",
year(col("first_appearance")))
# df_date_ops.select("name", "first_appearance", "debut_formatted",
"days_since_debut", "month_of_debut", "year_of_debut").show()

# Corrected for the current dataset, which has `DateType` for
`first_appearance`
from pyspark.sql.functions import date_format, datediff, current_date,
month, year, dayofmonth, dayofweek, dayofyear, weekofyear

df_date_ops = df.withColumn("debut_formatted",
date_format(col("first_appearance"), "yyyy-MM-dd")) \
                .withColumn("days_since_debut", datediff(current_date(),
col("first_appearance"))) \
                .withColumn("month_of_debut",
month(col("first_appearance"))) \
                .withColumn("year_of_debut",
year(col("first_appearance"))) \
                .withColumn("day_of_month_debut",
dayofmonth(col("first_appearance"))) \
                .withColumn("day_of_week_debut",
dayofweek(col("first_appearance"))) \
                .withColumn("day_of_year_debut",
dayofyear(col("first_appearance"))) \
                .withColumn("week_of_year_debut",
weekofyear(col("first_appearance")))

df_date_ops.select("name", "first_appearance", "debut_formatted",
"days_since_debut",
                   "month_of_debut", "year_of_debut",
"day_of_month_debut",
                   "day_of_week_debut", "day_of_year_debut",
"week_of_year_debut").show()
```

```python
# Scenario: Calculate age of hero (in years) based on current date
from pyspark.sql.functions import datediff, current_date, expr
df_hero_age = df.withColumn("hero_age_years", datediff(current_date(),
col("first_appearance")) / 365.25)
df_hero_age.select("name", "first_appearance", "hero_age_years").show()

# Scenario: Filter heroes who debuted in a specific month (e.g., August)
df_august_debuts = df.filter(month(col("first_appearance")) == 8)
df_august_debuts.select("name", "first_appearance").show()
```

**Scenarios**:

- **Time-Series Analysis**: Extracting components like year, month, day, hour for time-based aggregations.

- **Age/Duration Calculation**: Computing differences between dates (e.g., employee tenure, product shelf life).

- **Filtering by Date Range**: Selecting data within specific date or time windows.

- **Data Standardization**: Formatting date/time columns into consistent strings.

```python
from pyspark.sql.functions import date_format

df_date_ops = df.withColumn("debut_formatted",
date_format(col("first_appearance"), "yyyy-MM-dd")) \
                .withColumn("days_since_debut", datediff(current_date(),
col("first_appearance"))) \
                .withColumn("month_of_debut",
month(col("first_appearance"))) \
                .withColumn("year_of_debut",
year(col("first_appearance"))) \
                .withColumn("day_of_month_debut",
dayofmonth(col("first_appearance"))) \
                .withColumn("day_of_week_debut",
dayofweek(col("first_appearance"))) \
                .withColumn("day_of_year_debut",
dayofyear(col("first_appearance"))) \
                .withColumn("week_of_year_debut",
weekofyear(col("first_appearance"))) \
                .withColumn("day_name_debut",
date_format(col("first_appearance"), "EEEE")) \
                .withColumn("month_name_debut",
date_format(col("first_appearance"), "MMMM"))

df_date_ops.select(
    "name", "first_appearance", "debut_formatted",
    "days_since_debut", "month_of_debut", "month_name_debut",
    "year_of_debut", "day_of_month_debut", "day_of_week_debut",
```

```
        "day_name_debut", "day_of_year_debut", "week_of_year_debut"
).show(truncate=False)
```

## 📅 Format Patterns You Can Use:

| Pattern | Meaning | Example |
|---------|---------|---------|
| EEEE | Full name of the day | Thursday |
| EE | Abbreviated day | Thu |
| MMMM | Full month name | August |
| MMM | Abbreviated month | Aug |
| MM | Numeric month | 08 |

## Scenario-Based Questions & Answers (Additional Essential Functions)

1. **Scenario**: You want to rank heroes within each `universe` based on their `debut_year`, with the earliest debut having rank 1. If multiple heroes debuted in the same year, they should receive the same rank, and no ranks should be skipped.

   ```
   * **Question**: Which window function and ranking type would you use?
   * **Answer**:
     ```python
   ```

   from pyspark.sql import Window
   from pyspark.sql.functions import col, dense_rank
   window_spec_universe_debut = Window.partitionBy("universe").orderBy("debut_year")
   df_ranked_heroes = df.withColumn("debut_dense_rank",
   dense_rank().over(window_spec_universe_debut))
   df_ranked_heroes.select("name", "universe", "debut_year",
   "debut_dense_rank").orderBy("universe", "debut_dense_rank").show()
   ```

2. **Scenario**: For data presentation, you need to display hero names in lowercase and replace all spaces in their names with underscores.
   - **Question**: How would you achieve this using string functions?
   - **Answer**:

     ```
     df_cleaned_names = df.withColumn("name_lower_underscore",
     lower(regexp_replace(col("name"), " ", "_")))
     df_cleaned_names.select("name", "name_lower_underscore").show()
     ```

3. **Scenario**: You want to find all heroes who have "Flight" as one of their powers and also belong to at least one `team`.

- **Question**: How would you filter the DataFrame using array functions and a combined condition?
- **Answer**:

```
df_flying_team_heroes = df.filter(array_contains(col("powers"),
"Flight") & (size(col("team")) > 0))
df_flying_team_heroes.select("name", "powers",
"team").show(truncate=False)
```

4. **Scenario**: Your analysts need to categorize heroes based on which quarter of the year they first appeared.

- **Question**: How would you add a `debut_quarter` column using date functions?
- **Answer**:

```
from pyspark.sql.functions import quarter
df_with_quarter = df.withColumn("debut_quarter",
quarter(col("first_appearance")))
df_with_quarter.select("name", "first_appearance",
"debut_quarter").show()
```

5. **Scenario**: You want to identify the longest `name` among all superheroes and display that hero's record.

- **Question**: How would you combine string functions with sorting/ordering to find this hero?
- **Answer**:

```
df_longest_name = df.withColumn("name_length", length(col("name"))) \
                    .orderBy(desc("name_length")) \
                    .limit(1)
df_longest_name.select("name", "name_length").show()
```

## Performance Tips and Best Practices

Mastering PySpark DataFrames goes beyond knowing individual functions; it's about understanding how to combine them efficiently for large-scale data processing.

## 1. Optimization Strategies

- `select()` **Early to Reduce Data Size**: Always project only the necessary columns as early as possible in your pipeline. Fewer columns mean less data to process, shuffle, and store,

which directly translates to faster execution and lower memory footprint.

- **Filter Data as Early as Possible**: Similar to `select()`, apply `filter()` (or `where()`) transformations early. Reducing the number of rows before complex transformations (like joins or aggregations) minimizes the amount of data that needs to be processed.

- **Use Appropriate Join Types and Broadcast for Small Tables**:
  - Choose the correct join type (`inner`, `left`, `right`, `full`, `semi`, `anti`) for your specific use case to avoid unnecessary data transfer.
  - When joining a large DataFrame with a small DataFrame (typically under 10GB, but check Spark documentation/configuration for exact thresholds), use `broadcast(smaller_df)` on the smaller DataFrame. This sends the entire smaller DataFrame to all worker nodes, avoiding a costly shuffle of the larger DataFrame.

- **Cache Frequently Accessed DataFrames (`df.cache()`)**: If you're going to reuse a DataFrame multiple times in your application, especially after expensive transformations, `cache()` it. This stores the DataFrame (or its partitions) in memory (and optionally on disk) across different operations, preventing re-computation. Remember to `unpersist()` it when no longer needed to free up resources.

## 2. Memory Management

- **Use `coalesce()` to Reduce Partitions Before Writing**: Before writing a DataFrame to external storage, use `coalesce()` to reduce the number of partitions. Too many small files (the "small file problem") can degrade read performance and overload metadata services in distributed file systems. `coalesce()` is efficient as it minimizes data shuffling.

- **Avoid `collect()` on Large Datasets**: As emphasized repeatedly, `collect()` brings all data to the driver program, which will cause OutOfMemory errors for large datasets. Use `show()`, `take()`, `first()`, or write to external storage instead.

- **Use `limit()` for Testing Transformations**: When developing or debugging, use `.limit(N)` after transformations to inspect a small subset of the resulting data without triggering a full computation on the entire dataset.

## 3. Code Organization

- **Create Reusable Transformation Functions**: Encapsulate common sequences of transformations into Python functions that accept and return a DataFrame. This promotes modularity, testability, and reduces code duplication, especially when using `transform()`.

- **Use Meaningful Variable Names**: Clear and descriptive variable names improve code readability and maintainability.

- **Document Complex Transformations**: Add comments or docstrings to explain the purpose and logic of non-trivial transformations, especially for complex UDFs or window functions.

- **Test with Small Datasets First**: Always develop and test your Spark code on small, representative samples of your data before deploying to full production datasets. This speeds up the development cycle and saves compute resources.

## 4. Common Pitfalls to Avoid

- **Don't use `collect()` on large datasets**: This is the most common mistake for beginners and leads to driver OOM errors.
- **Avoid Unnecessary Shuffling Operations**: Operations like `groupBy()`, `orderBy()`, `repartition()`, and `join()` (without broadcast) involve data shuffling, which is the most expensive operation in Spark. Minimize these, and when unavoidable, ensure they are optimized (e.g., repartitioning by join keys before a join).
- **Don't Create Too Many Small Partitions**: As mentioned with `partitionBy()` and `repartition()`, too many tiny partitions lead to overhead in task scheduling and can overwhelm the file system.
- **Be Cautious with UDFs (Performance Impact)**: While UDFs are powerful, they are less performant than built-in Spark functions because they break Spark's Catalyst optimizer's ability to optimize the query plan. Use them only when absolutely necessary and consider Pandas UDFs for better performance if possible.

## 5. Testing and Debugging

- **Use `explain()` to Understand Execution Plans**: `df.explain()` (or `df.explain(True)` for detailed output) shows the physical and logical execution plan of your DataFrame operations. This is invaluable for identifying performance bottlenecks, understanding data shuffles, and verifying that Spark is optimizing your queries as expected.
- **Test with Subset of Data First**: Reduces iteration time during development.
- **Use `show()` and `printSchema()` for Debugging**: Quick checks to see the data content and schema after each transformation step.
- **Monitor Spark UI for Performance Issues**: The Spark Web UI (typically accessible on port 4040 of your driver node) provides detailed information about jobs, stages, tasks, and executors, including shuffle reads/writes, memory usage, and execution times. This is the ultimate tool for diagnosing performance problems.

## Conclusion

This comprehensive guide has covered a wide array of essential PySpark DataFrame functions, from basic data inspection to complex transformations, aggregations, and joins. We've explored how to manipulate columns, handle missing values, deduplicate data, and organize your DataFrame efficiently.

Remember that **DataFrame operations in Spark are lazy**: they are only executed when an **action** (like `show()`, `collect()`, `count()`, or `write()`) is called. This lazy evaluation, combined with Spark's Catalyst optimizer, allows for highly optimized execution plans.

The key to mastering PySpark DataFrames is not just memorizing functions, but understanding **when and how to use each function effectively**, always keeping **performance, scalability, and code readability** in mind. By following these best practices, you can build robust and efficient data processing pipelines for large-scale datasets.