```python
from pyspark.sql import SparkSession
from pyspark.sql.functions import col, when, count, max, min
from pyspark.sql.window import Window
import numpy as np

# Initialize Spark session
spark = SparkSession.builder.appName("FillMissingValues").getOrCreate()

# Sample Data
data = [
    (1, 25, 'North', 'M', '2025-01-01', 150),
    (2, None, 'East', None, '2025-01-02', None),
    (3, 30, 'South', 'F', None, 200),
    (4, 22, None, 'M', '2025-01-03', 180),
    (5, 28, 'West', 'F', None, None),
]

# Column names
columns = ['Customer_ID', 'Age', 'Region', 'Gender', 'Last_Visit', 'Purchase_Amou

# Create DataFrame
df = spark.createDataFrame(data, columns)

# Show the original DataFrame
df.show()

# Define a function to fill missing values dynamically
def fill_missing_values(df):
    # Get column types
    column_types = df.dtypes

    # Loop through each column based on type
```

```python
24   # Show the original DataFrame
25   df.show()
26
27   # Define a function to fill missing values dynamically
28   def fill_missing_values(df):
29       # Get column types
30       column_types = df.dtypes
31
32       # Loop through each column based on type
33       for column, dtype in column_types:
34           if dtype == 'int' or dtype == 'double':
35               # For numeric columns, fill with median
36               median_value = df.approxQuantile(column, [0.5], 0)[0]
37               df = df.fillna({column: median_value})
38           elif dtype == 'string':
39               # For string columns, fill with 'Unknown'
40               df = df.fillna({column: 'Unknown'})
41           else:
42               # For other columns (like date, etc.), fill with 'Unknown'
43               df = df.fillna({column: 'Unknown'})
44
45       return df
46
47   # Apply the function to fill missing values
48   filled_df = fill_missing_values(df)
49
50   # Show the updated DataFrame
51   filled_df.show()
52
53   # Stop the Spark session
54   spark.stop()
```

Capgemini

Karthik Kondpak

```
+----------+----+------+------+----------+---------------+
|Customer_ID| Age|Region|Gender|Last_Visit|Purchase_Amount|
+----------+----+------+------+----------+---------------+
|         1|  25| North|     M|2025-01-01|            150|
|         2|NULL|  East|  NULL|2025-01-02|           NULL|
|         3|  30| South|     F|      NULL|            200|
|         4|  22|  NULL|     M|2025-01-03|            180|
|         5|  28|  West|     F|      NULL|           NULL|
+----------+----+------+------+----------+---------------+


+----------+-------+-------+-------+----------+---------------+
|Customer_ID|    Age| Region| Gender|Last_Visit|Purchase_Amount|
+----------+-------+-------+-------+----------+---------------+
|         1|     25|  North|      M|2025-01-01|            150|
|         2|   NULL|   East|Unknown|2025-01-02|           NULL|
|         3|     30|  South|      F|   Unknown|            200|
|         4|     22|Unknown|      M|2025-01-03|            180|
|         5|     28|   West|      F|   Unknown|           NULL|
+----------+-------+-------+-------+----------+---------------+
```