# Data Engineering Core Skills – Exercises
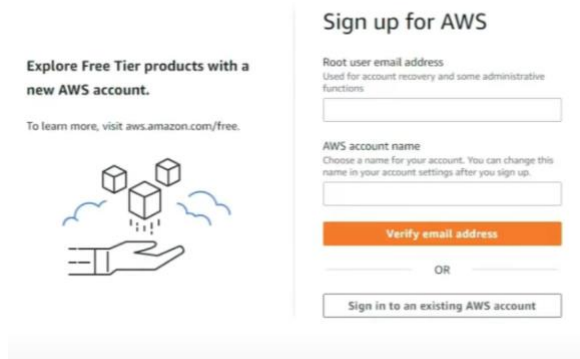
## 1. Cracking the Case

### 1.1.    Configuration

- **AWS Account Setup –** You'll need an AWS Account for this tutorial. To create an account, click <u>HERE</u>:
  - Click Free Account



  - Enter your email as root user
  - Enter an account name and credit card details (you will not be charged)



  - Enter a region – close to your country



  - Select support plan - Free

- Login – as root user



- **IAM Role Setup for AWS Glue**
    - Navigate to the IAM console
    - Click on **Roles** and **Create Role**.
    - Select Glue service from the section
    - **Choose the service that will use this role**.
    - Choose Glue for *Select your use case.*



    - Select **AWSGlueServiceRole** from **Attach Permissions Policies**.

## Add permissions

**Permissions policies** (Selected 1/857)
Choose one or more policies to attach to your new role.

🔍 Filter policies by property or policy name and press enter

"awsglueser" ✕ | **Clear filters**

| ☐ | Policy name 🔗 |
|---|---|
| ☐ | ⊞ AWSGlueServiceRole-aws-audio-analysis |
| ☐ | ⊞ AWSGlueServiceRole-aws-audio-analysis-iam |
| ☐ | ⊞ AWSGlueServiceRole-demo |
| ☐ | ⊞ AWSGlueServiceRole-hh |
| ☐ | ⊞ AWSGlueServiceRole-xxx |
| ☐ | ⊞ 📦 AWSGlueServiceNotebookRole |
| ☑ | ⊞ 📦 AWSGlueServiceRole |

- Now, write S3Full in the filter box above to add **S3FullAccess**
- Leave the tag section blank and then create the role.
- **S3 Bucket Setup**
  - Navigate to the S3 console.
  - Click **Create bucket**. Give it a name: revenue-loss-raw-data-lake

### Create bucket Info
Buckets are containers for data stored in S3.

**General configuration**

AWS Region
Europe (Stockholm) eu-north-1

Bucket type    Info

○ General purpose
Recommended for most use cases and access patterns.
General purpose buckets are the original S3 bucket type.
They allow a mix of storage classes that redundantly
store objects across multiple Availability Zones.

○ Directory
Recommended for low-latency use cases. These buckets
use only the S3 Express One Zone storage class, which
provides faster processing of data within a single
Availability Zone.

Bucket name    Info

revenue-loss-raw-data-lake

Bucket name must be unique within the global namespace and follow the bucket naming rules. See rules for bucket naming 🔗

Copy settings from existing bucket - *optional*
Only the bucket settings in the following configuration are copied.

Choose bucket

Format: s3://bucket/prefix

**Object Ownership**    Info
Control ownership of objects written to this bucket from other AWS accounts and the use of access control lists (ACLs). Object ownership
determines who can specify access to objects.

○ ACLs disabled (recommended)
All objects in this bucket are owned by this account.
Access to this bucket and its objects is specified using
only policies.

○ ACLs enabled
Objects in this bucket can be owned by other AWS
accounts. Access to this bucket and its objects can be
specified using ACLs.

Object Ownership
Bucket owner enforced

**Block Public Access settings for this bucket**
Public access is granted to buckets and objects through access control lists (ACLs), bucket policies, access point policies, or all. In order to
ensure that public access to this bucket and its objects is blocked, turn on Block all public access. These settings apply only to this bucket
and its access points. AWS recommends that you turn on Block all public access, but before applying any of these settings, ensure that your
applications will work correctly without public access. If you require some level of public access to this bucket or objects within, you can
customize the individual settings below to suit your specific storage use cases. Learn more 🔗

☑ Block *all* public access

- Select your region
- Finish creating bucket.


- **AWS Glue Setup**

- Navigate to the [Glue console](#)
- In getting started, **Click on Setup roles and users from**
- Click on 'Choose Roles' and Confirm.



- Click on Next, Next and Apply Changes.

# 1.2.   Create Glue Script

- Navigate to the [Glue console](#)
- Click ETL Jobs
- Click on Script editor
- Enter the following

```python
import requests
from awsglue.context import GlueContext
from pyspark.context import SparkContext
from awsglue.dynamicframe import DynamicFrame  # Import the correct module

# Initialize GlueContext and SparkContext
sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)

# Make the GET request to the API
api_url = "https://github.com/nextgendata-academy/revenue-loss/blob/main/data/clickstream_data.csv"
response = requests.get(api_url)

# Check if the request was successful
if response.status_code == 200:

    # Read the CSV data directly from the API response
    csv_data = response.content.decode('utf-8')

    # API response is a string, split it into rows by breaking it at (\n)
    # Convert the CSV data into an RDD format that Spark understands
    rdd = sc.parallelize(csv_data.splitlines())
```

```
# Convert RDD into Spark DataFrame, with a schema
df = glueContext.spark_session.read.csv(rdd, header=True, inferSchema=True)

# Convert the Spark DataFrame to a Glue DynamicFrame
dynamic_frame = DynamicFrame.fromDF(df, glueContext)  # Corrected method usage

# Write the data to S3 in CSV format
glueContext.write_dynamic_frame.from_options(
    frame=dynamic_frame,
    connection_type="s3",
    connection_options={"path": "s3://revenue-loss-raw-data-lake"},
    format="csv"
)

else:
    print(f"Failed to fetch data. Status code: {response.status_code}")
```

- Save Script and Run
- Click on "Runs" Tab to see the if the job succeeded.
- Navigate to the S3 console.
- Select the bucket to see if the data is ingested.

# 2. Section 2: Hands on Exercise
## 2.1. Create Data Brew Job for Data Cleaning

## 2.2. Create Glue Job for Data Cleaning

```
import sys
from awsglue.transforms import *
from awsglue.utils import getResolvedOptions
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.job import Job
from pyspark.sql.functions import col, when

# Initialize Glue context
glueContext = GlueContext(SparkContext.getOrCreate())
spark = glueContext.spark_session

# Load data from S3 (assuming it's already in the Glue Data Catalog)
clickstream_data = glueContext.create_dynamic_frame.from_catalog(database="your_database",
table_name="your_table")

# Convert to DataFrame for easier manipulation
df = clickstream_data.toDF()

# Data cleaning and transformation
# Handle missing ProductID by removing rows
df_cleaned = df.filter(df.ProductID.isNotNull())

# Standardize date formats (assuming SessionDate is in multiple formats)
df_cleaned = df_cleaned.withColumn("SessionDate", date_format(col("SessionDate"), "yyyy-MM-dd"))

# Handle outliers in TimeSpent by capping values above a threshold
df_cleaned = df_cleaned.withColumn("TimeSpent", when(col("TimeSpent") > 3600,
3600).otherwise(col("TimeSpent")))

# Remove duplicate entries based on SessionID
df_cleaned = df_cleaned.dropDuplicates(["SessionID"])

# Convert DataFrame back to DynamicFrame for writing
dynamic_frame_cleaned = DynamicFrame.fromDF(df_cleaned, glueContext, "dynamic_frame_cleaned")

# Save the cleaned data back to S3
glueContext.write_dynamic_frame.from_options(
    frame = dynamic_frame_cleaned,
    connection_type = "s3",
    connection_options = {"path": "s3://your-bucket/cleaned_data/"},
    format = "parquet"
)

# End Glue job
job.commit()
```

# 3. Section 2: Hands on Exercise
## 3.1. Create Glue Job to solve mystery: Lost customers

```
import sys
from awsglue.transforms import *
from awsglue.utils import getResolvedOptions
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.job import Job
from pyspark.sql.functions import col, count, when

# Initialize Glue context and job parameters
args = getResolvedOptions(sys.argv, ['JOB_NAME'])
sc = SparkContext()
glueContext = GlueContext(sc)
spark = glueContext.spark_session
job = Job(glueContext)
job.init(args['JOB_NAME'], args)

# Input S3 path and output paths
input_s3_path = "s3://your-bucket-name/clickstream_data.csv"
output_s3_path_cleaned = "s3://your-bucket-name/cleaned_clickstream_data/"
output_s3_path_disappearing = "s3://your-bucket-name/disappearing_visitors/"
output_s3_path_checkout = "s3://your-bucket-name/checkout_dropoff/"

# Load clickstream data from S3
clickstream_df = spark.read.csv(input_s3_path, header=True, inferSchema=True)

# === 1. The Disappearing Visitors: Identify users who only visited pages without engagement ===
disappearing_visitors = clickstream_df.groupBy("customerId", "sessionId").agg(
    count(when(col("eventType") == "pageView", 1)).alias("pageViewCount"),
    count(when(col("eventType") != "pageView", 1)).alias("otherEventCount")
).filter((col("pageViewCount") > 0) & (col("otherEventCount") == 0))

# Write disappearing visitors data to S3
disappearing_visitors.write.mode("overwrite").csv(output_s3_path_disappearing)

# === 2. The Missing Product Insights: Data Cleaning ===
cleaned_clickstream = clickstream_df.fillna({'productId': -1}) \
    .fillna({'customerId': 'Unknown'}) \
    .dropna(subset=['sessionId', 'eventType'])

# Drop duplicates
cleaned_clickstream = cleaned_clickstream.dropDuplicates()

# Write cleaned data to S3
cleaned_clickstream.write.mode("overwrite").csv(output_s3_path_cleaned)

# === 3. The Checkout Drop-Off: Identify users who added items to cart but didn't checkout ===
checkout_drop_off = cleaned_clickstream.groupBy("customerId", "sessionId").agg(
    count(when(col("eventType") == "addToCart", 1)).alias("addToCartCount"),
    count(when(col("eventType") == "checkout", 1)).alias("checkoutCount")
).filter((col("addToCartCount") > 0) & (col("checkoutCount") == 0))

# Write checkout drop-off data to S3
checkout_drop_off.write.mode("overwrite").csv(output_s3_path_checkout)

# Commit job
job.commit()
```

# 4. Section 4: Hands on Exercise

```
import sys
from awsglue.transforms import *
```

```python
from awsglue.utils import getResolvedOptions
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.job import Job
from pyspark.sql.functions import col, count, when

# Initialize Glue context and job parameters
args = getResolvedOptions(sys.argv, ['JOB_NAME'])
sc = SparkContext()
glueContext = GlueContext(sc)
spark = glueContext.spark_session
job = Job(glueContext)
job.init(args['JOB_NAME'], args)

from pyspark.sql import functions as F

# Load the cleaned clickstream data
clickstream_df = spark.read.csv("s3://your-bucket-name/cleaned_clickstream_data/", header=True,
inferSchema=True)

# === Fact Table: Customer Sessions ===
fact_customer_sessions = clickstream_df.groupBy("customerId", "sessionId").agg(
    F.count(when(col("eventType") == "pageView", 1)).alias("pageViewCount"),
    F.count(when(col("eventType") == "view_product", 1)).alias("productViewCount"),
    F.count(when(col("eventType") == "addToCart", 1)).alias("addToCartCount"),
    F.count(when(col("eventType") == "checkout", 1)).alias("checkoutCount"),
    F.max("sessionDate").alias("sessionDate"),  # Optional to store the latest session date
    F.max("timeSpent").alias("timeSpent")  # Total time spent in session
)

# Create a sessionType column to label the session based on conditions
fact_customer_sessions = fact_customer_sessions.withColumn(
    "sessionType",
    F.when((F.col("pageViewCount") > 0) & (F.col("productViewCount") == 0), "ViewingButNoProduct")
    .when((F.col("addToCartCount") > 0) & (F.col("checkoutCount") == 0), "CheckoutDropOff")
    .when((F.col("pageViewCount") > 0) & (F.col("productViewCount") == 0) & (F.col("addToCartCount") == 0),
"Disappearing")
    .otherwise("Other")
)

# Write the fact table to S3
fact_customer_sessions.write.mode("overwrite").csv("s3://your-bucket-name/fact_customer_sessions/")

# === Dimension Table: Customer ===
dim_customer = clickstream_df.select("customerId", "customerName", "email", "age", "city").distinct()

# Write the customer dimension table to S3
dim_customer.write.mode("overwrite").csv("s3://your-bucket-name/dim_customer/")

# === Dimension Table: Product ===
dim_product = clickstream_df.select("productId").distinct()

# Optionally, join product metadata if available
# product_metadata = spark.read.csv("s3://your-bucket-name/product_metadata.csv", header=True,
inferSchema=True)
# dim_product = dim_product.join(product_metadata, "productId", "left")

# Write the product dimension table to S3
dim_product.write.mode("overwrite").csv("s3://your-bucket-name/dim_product/")


# Commit job
job.commit()
```