

```

1 Lab. AWS S3 Bucket 생성 및 데이터 저장
2
3 1. S3 Bucket 생성
4 1)[서비스] > [스토리지] > S3
5 2)[버킷 만들기] 버튼 클릭
6 3)[버킷 만들기] 페이지에서, [버킷 이름] : {계정}-datalake-bucket
7 4)[AWS 리전] : 아시아 태평양(서울) ap-northeast-2
8 5)[버킷 만들기] 버튼 클릭
9
10
11 2. AWS CLI를 사용하여 Bucket List 출력하기
12 1)Windows Command 창 또는 macOS Terminal에서
13 2)AWS Access Key ID와 AWS Secret Access Key 입력
14 $ aws configure
15 AWS Access Key ID [*****GYVQ]:
16 AWS Secret Access Key [*****jWiz]:
17 Default region name [ap-northeast-2]:
18 Default output format [json]:
19
20 3)S3 Bucket List 출력
21 $ aws s3 ls /
22 2023-03-08 09:56:42 henry-datalake-bucket
23
24 4)해당 Bucket 내용 출력
25 $ aws s3 ls s3://{bucket name}
26 <--- Bucket 내에 어떤 Object도 없기 때문에 아무 것도 출력되지 않음.
27
28
29 3. Lab에서 사용할 Public DataSet 확인
30 1)Google에서 "aws public datasets"로 검색
31 2)검색 결과에서 [Open Data on AWS] 링크 클릭,
32 https://aws.amazon.com/ko/opendata/?wwps-cards.sort-by=item.additionalFields.sortDate&wwps-cards.sort-order=desc
33 3)페이지에서 [Find publicly available data on AWS] 버튼 클릭
34 4)검색창에 "taxi" 입력하여 "New York City Taxi and Limousine Commission(TLC) Trip Record Data" 클릭
35 -https://aws.amazon.com/marketplace/pp/prodview-okyonroqg5b2u?sr=0-1&ref\_=beagle&applicationId=AWSMPContessa
36 5)[New York City Taxi and Limousine Commission (TLC) Trip Record Data] 페이지에서, [Description] 탭에서 [Documentation]의 링크
37 클릭
38 6)[TLC Trip Record Data] 페이지에서, [Data Dictionaries and MetaData] 섹션의 "Yellow Trips Data Dictionary" 클릭하여 문서의 내용 파악
39 7)다시 [New York City Taxi and Limousine Commission (TLC) Trip Record Data] 페이지로 돌아와서, [Resources on AWS] 탭으로 이동
40 8)[AWS CLI Access]의 값 확인
41 aws s3 ls s3://nyc-tlc/
42 9)Windows Command 창 또는 macOS의 Terminal에서,
43 $ aws s3 ls s3://nyc-tlc/
44 PRE csv_backup/
45 PRE misc/
46 PRE trip data/
47 10)Object들 중에서 "trip data" 검색
48 $ aws s3 ls s3://nyc-tlc/"trip data"/
49
50 11)검색 결과 중 "2022-10" 필터하기
51 -macOS
52 $ aws s3 ls s3://nyc-tlc/"trip data"/ | grep 2022-10
53 -Windows
54 >aws s3 ls s3://nyc-tlc/"trip data"/ | findstr "2022-10"
55
56 2022-12-20 06:42:16 12051434 fhv_tripdata_2022-10.parquet
57 2022-12-20 06:42:12 495083481 fhvhv_tripdata_2022-10.parquet
58 2022-12-20 06:42:14 1444642 green_tripdata_2022-10.parquet
59 2022-12-20 06:42:12 57061938 yellow_tripdata_2022-10.parquet
60
61 4. "trip-data"의 데이터를 위에서 생성한 나의 Bucket으로 복사하기
62 1)"trip-data"의 green_tripdata_2022-10.parquet를 위에서 생성한 나의 Bucket으로 복사하기
63 $ aws s3 cp s3://nyc-tlc/"trip data"/green_tripdata_2022-10.parquet
64 s3://henry-datalake-bucket/input/green_tripdata_2022-10.parquet
65 copy: s3://nyc-tlc/trip data/green_tripdata_2022-10.parquet to
66 s3://henry-datalake-bucket/input/green_tripdata_2022-10.parquet
67
68 2)"trip-data"의 yellow_tripdata_2022-10.parquet를 위에서 생성한 나의 Bucket으로 복사하기
69 $ aws s3 cp s3://nyc-tlc/"trip data"/yellow_tripdata_2022-10.parquet
70 s3://henry-datalake-bucket/input/yellow_tripdata_2022-10.parquet
71 copy: s3://nyc-tlc/trip data/yellow_tripdata_2022-10.parquet to
72 s3://henry-datalake-bucket/input/yellow_tripdata_2022-10.parquet
73
74 3)해당 파일들 복사되었는지 확인하기
75 $ aws s3 ls s3://henry-datalake-bucket/input/
76 2023-03-08 10:29:22 1444642 green_tripdata_2022-10.parquet
77 2023-03-08 10:31:50 57061938 yellow_tripdata_2022-10.parquet
78
79 5. Local Machine에 CSV 파일 다운로드하여 Head 확인하기
80 1)[New York City Taxi and Limousine Commission (TLC) Trip Record Data]의 CSV 파일 목록 확인
81 $ aws s3 ls s3://nyc-tlc/csv_backup/

```

```
79 2)특정 CSV 파일 다운로드
80 $ aws s3 cp s3://nyc-tlc/csv_backup/yellow_tripdata_2020-04.csv .
81
```

```
82 3)CSV 파일 앞 부분 확인
83 -Windows
84 >more yellow_tripdata_2020-04.csv
85 -macOS
86 $ head yellow_tripdata_2020-04.csv
87
88
```

VendorID, tpep_pickup_datetime, tpep_dropoff_datetime, passenger_count, trip_distance, RatecodeID, store_and_fwd_flag, PULocationID, DOLocationID, payment_type, fare_amount, extra, mta_tax, tip_amount, tolls_amount, improvement_surcharge, total_amount, congestion_surcharge

```
89 1,2020-04-01 00:41:22,2020-04-01 01:01:53,1,1.20,1,N,41,24,2,5.5,0.5,0.5,0,0,0.3,6.8,0
90 1,2020-04-01 00:56:00,2020-04-01 01:09:25,1,3.40,1,N,95,197,1,12.5,0.5,0.5,2.75,0,0.3,16.55,0
91 1,2020-04-01 00:00:26,2020-04-01 00:09:25,1,2.80,1,N,237,137,1,10,3,0.5,1,0,0.3,14.8,2.5
92 1,2020-04-01 00:24:38,2020-04-01 00:34:38,0,2.60,1,N,68,142,1,10,3,0.5,1,0,0.3,14.8,2.5
93 2,2020-04-01 00:13:24,2020-04-01 00:18:26,1,1.44,1,Y,263,74,1,6.5,0.5,0.5,3,0,0.3,13.3,2.5
94 2,2020-04-01 00:24:36,2020-04-01 00:33:09,1,2.93,1,N,75,170,2,10.5,0.5,0.5,0,0,0.3,14.3,2.5
```