

```

1 Lab2. 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] : {계정명}-datalake-bucket
7   4)[AWS 리전] : 아시아 태평양(서울) ap-northeast-2
8   5)[버킷 만들기] 버튼 클릭
9
10
11
12 2. AWS CLI를 사용하여 Bucket List 출력하기
13   1)Windows Command 창 또는 macOS Terminal에서
14   2)AWS Access Key ID와 AWS Secret Access Key 입력
15     $ aws configure
16     AWS Access Key ID [None]: <---- Access Key 붙여넣기
17     AWS Secret Access Key [None]: <---- Secret Key 붙여넣기
18     Default region name [None]: ap-northeast-2
19     Default output format [None]: json
20
21   3)S3 Bucket List 출력
22     $ aws s3 ls /
23     2023-03-08 09:56:42 {계정명}-datalake-bucket
24
25   4)해당 Bucket 내용 출력
26     $ aws s3 ls s3://{bucket name}
27     <--- Bucket 내에 어떤 Object도 없기 때문에 아무 것도 출력되지 않음.
28
29
30
31 3. Lab에서 사용할 Public DataSet 확인
32   1)Google에서 "aws public datasets"로 검색
33
34   2)검색 결과에서 [Open Data on AWS] 링크 클릭,
35     https://aws.amazon.com/ko/opendata/?wwps-cards.sort-by=item.additionalFields.sortDate&wwps-cards.sort-order=desc
36
37   3)페이지에서 [Find publicly available data on AWS] 버튼 클릭
38
39   4)검색창에 "taxi" 입력하여 "New York City Taxi and Limousine Commission(TLC) Trip Record Data" 클릭
40     -https://aws.amazon.com/marketplace/pp/prodview-okyonroqg5b2u?sr=0-1&ref_=beagle&applicationId=AWSMPContessa
41
42   5)[New York City Taxi and Limousine Commission (TLC) Trip Record Data] 페이지에서, [Description] 탭에서 [Documentation]의 링크
43     클릭
44
45   6)[TLC Trip Record Data] 페이지에서, [Data Dictionaries and MetaData] 섹션의 "Yellow Trips Data Dictionary" 클릭하여 문서의 내용 파악
46
47   7)다시 [New York City Taxi and Limousine Commission (TLC) Trip Record Data] 페이지로 돌아와서, [Resources on AWS] 탭으로 이동
48
49   8)[AWS CLI Access]의 값 확인
50     aws s3 ls s3://nyc-tlc/
51
52   9)Windows Command 창 또는 macOS의 Terminal에서,
53     $ aws s3 ls s3://nyc-tlc/
54     PRE csv_backup/
55     PRE misc/
56     PRE trip data/
57
58   10)Object들 중에서 "trip data" 검색
59     $ aws s3 ls s3://nyc-tlc/"trip data"/
60
61   11)검색 결과 중 "2022-10" 필터하기
62     -macOS
63       $ aws s3 ls s3://nyc-tlc/"trip data"/ | grep 2022-10
64
65     -Windows
66       >aws s3 ls s3://nyc-tlc/"trip data"/ | findstr "2022-10"
67
68     2022-12-20 06:42:16 12051434 fhv_tripdata_2022-10.parquet
69     2022-12-20 06:42:12 495083481 fhvhv_tripdata_2022-10.parquet
70     2022-12-20 06:42:14 1444642 green_tripdata_2022-10.parquet
71     2022-12-20 06:42:12 57061938 yellow_tripdata_2022-10.parquet
72
73 4. "trip-data"의 데이터를 위에서 생성한 나의 Bucket으로 복사하기
74   1)"trip-data"의 green_tripdata_2022-10.parquet를 위에서 생성한 나의 Bucket으로 복사하기
75     $ aws s3 cp s3://nyc-tlc/"trip data"/green_tripdata_2022-10.parquet
76     s3://{계정명}-datalake-bucket/input/green_tripdata_2022-10.parquet
77     copy: s3://nyc-tlc/trip data/green_tripdata_2022-10.parquet to
78     s3://{계정명}-datalake-bucket/input/green_tripdata_2022-10.parquet
79
80   2)"trip-data"의 yellow_tripdata_2022-10.parquet를 위에서 생성한 나의 Bucket으로 복사하기
81     $ aws s3 cp s3://nyc-tlc/"trip data"/yellow_tripdata_2022-10.parquet
82     s3://{계정명}-datalake-bucket/input/yellow_tripdata_2022-10.parquet

```

```

80 copy: s3://nyc-tlc/trip data/yellow_tripdata_2022-10.parquet to
81 s3://{계정명}-datalake-bucket/input/yellow_tripdata_2022-10.parquet
82
83 3) 해당 파일들 복사되었는지 확인하기
84 $ aws s3 ls s3://{계정명}-datalake-bucket/input/
85 2023-03-08 10:29:22 1444642 green_tripdata_2022-10.parquet
86 2023-03-08 10:31:50 57061938 yellow_tripdata_2022-10.parquet
87
88
89 5. Local Machine에 CSV 파일 다운로드하여 Head 확인하기
90 1)[New York City Taxi and Limousine Commission (TLC) Trip Record Data]의 CSV 파일 목록 확인
91 $ aws s3 ls s3://nyc-tlc/csv_backup/
92
93 2) 특정 CSV 파일 다운로드
94 $ aws s3 cp s3://nyc-tlc/csv_backup/yellow_tripdata_2020-04.csv . <---마지막 '.' 주의
95
96 3) CSV 파일 앞 부분 확인
97 -Windows
98 >more yellow_tripdata_2020-04.csv
99
100 -macOS
101 $ head yellow_tripdata_2020-04.csv
102
103 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
104 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
105 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
106 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
107 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
108 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
109 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

```