

# AWS기반 COVID 19 데이터 분석과 시각화

---

발표 날짜	2021.09.27
팀명	N조
팀원 이름	김현수 나혜주 레이첼 간



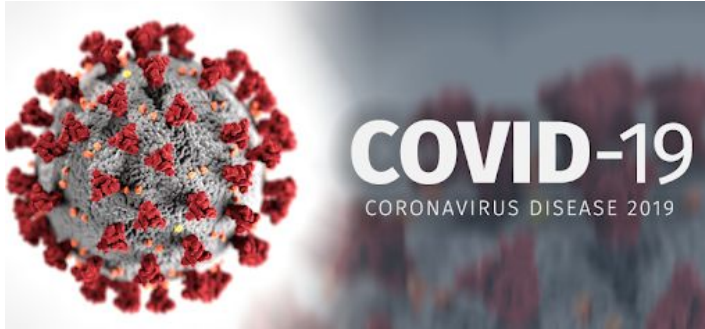
# Table of Contents

---

- Overview
- Goal/Problem & Requirement
- Approach
- Basic Spec
- Specifications
- Development Environment
- Architecture
- Current Status
- Further Plan
- Division & Assignment of Work
- Schedule

# Overview

---



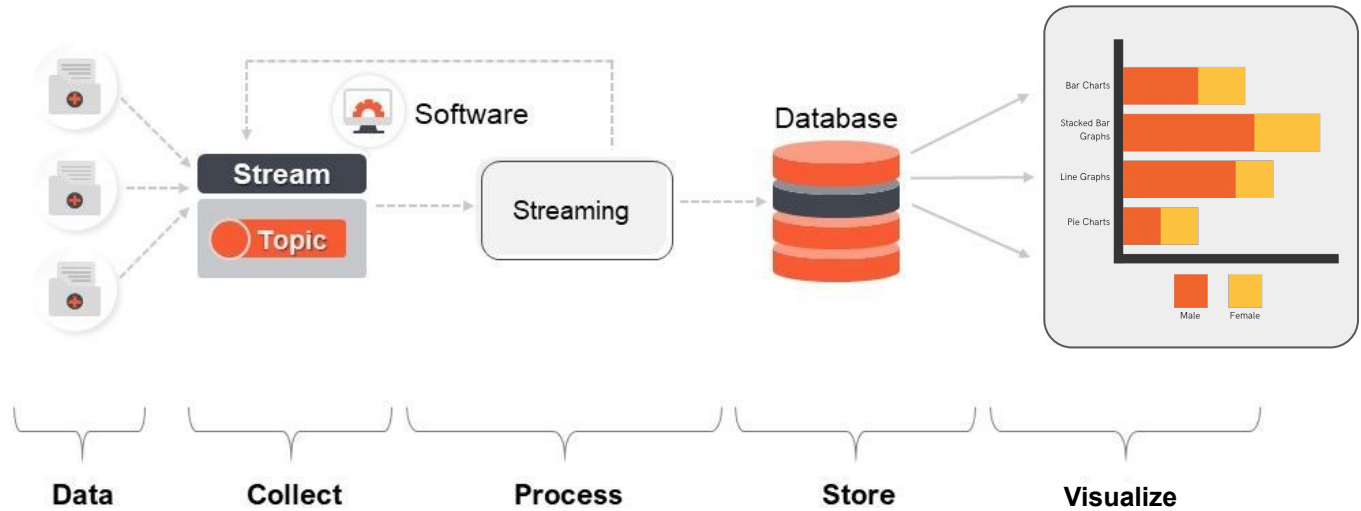
## 프로젝트 내용

- COVID 19 관련 다양한 공개 데이터와 소셜 미디어 자료를 활용해 데이터 수집과 저장,
- 데이터 전처리와 분석, 그리고 시각화까지 데이터 분석 플랫폼을 구축한다.

## 교육/훈련 효과

- 데이터 분석의 전반적인 과정과 그에 필요한 지식 습득, 데이터 분석 오픈 소스 활용 경험,
- 데이터 시각화를 통한 인사이트 도출

# Goal



COVID-19의 다양한 데이터 셋을 활용해 데이터 수집과 저장, 데이터 전처리와 분석, 그리고 시각화까지 데이터 분석 플랫폼을 구축한다.

- 필요에 따라 AWS의 다양한 데이터 분석 서비스를 사용
- 사용자가 원하는 형태의 데이터 시각화 구현

## 구현 내용

- 데이터 수집기 : COVID-19 및 관련 자료 수집기
- 데이터 처리기: 원하는 정보로 가공 가능하도록 만들기
- 대시 보드: 다양한 COVID-19 트렌드 및 관련 소셜 미디어 데이터 연동

# Requirement

- AWS 서비스와 기타 API, CLI를 사용해 데이터를 수집 및 이관하는 아키텍처를 3개 이상 구현
- 4개 이상의 데이터셋을 통해 적절한 데이터 전처리 과정을 구현
- 서로 다른 4개 이상의 AWS 서비스를 사용해 데이터 분석 아키텍처를 구현
- 원하는 형태의 대시보드 기능 구현이 가능하며 분석에 대한 기능 항목 3개 이상을 설정하고 구현



# Approach 주제탐색

---

## 주제 선정 시 고려 사항

데이터 셋이

- (1) 코로나 전과 후 각각에 대해 동일 조건으로 있어야 함
- (2) 최근 데이터까지도 있어야 함
- (3) 데이터의 주기가 1년 미만

## 데이터로 해보고 싶은 것

- (1) 코로나가 경제, 환경, 생활 등 다양한 분야에 어떤 영향을 미쳤는지  
알아보자
- (2) 코로나 사태와 이전 전염병 사태, 나아가 경제 위기 등등을 비교해보자.

# Approach 주제선정

(1)주가



(2)실업률

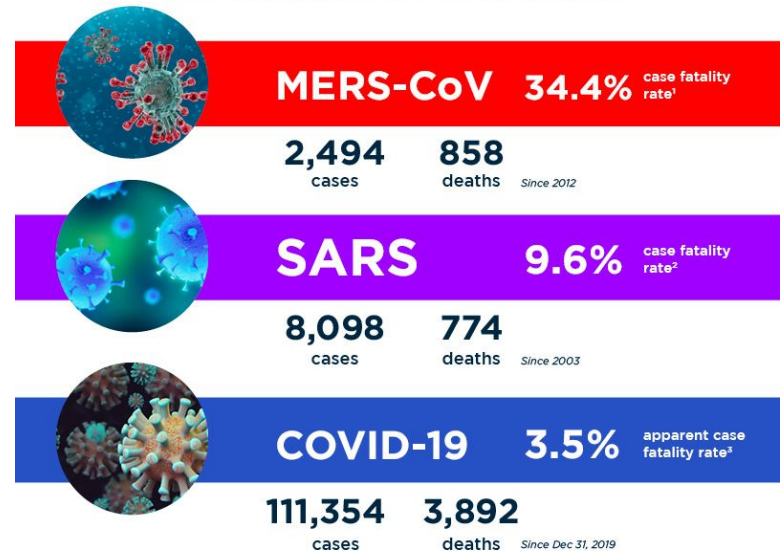


(3)대기오염도



## • 다른 질병, 경제적 위기 vs 코로나

### How COVID-19 Compares to Other Coronaviruses as of 3/9/20



Source:  
<sup>1</sup>WHO [www.who.int/emergencies/mers-cov/en](http://www.who.int/emergencies/mers-cov/en)  
<sup>2</sup>WHO [www.who.int/ith/diseases/sars/en](http://www.who.int/ith/diseases/sars/en)  
<sup>3</sup>Johns Hopkins University [gisanddata.maps.arcgis.com/apps/opsdashboard/index.html#/bda7594740fd40299423467b48e9ecf6](https://gisanddata.maps.arcgis.com/apps/opsdashboard/index.html#/bda7594740fd40299423467b48e9ecf6)

# Basic Specifications

---

## Lake Formation

- 3개의 주제(대기오염, 실업률, 주가)를 기반으로 데이터 수집
- 데이터 레이크를 생성하고 메타데이터를 저장하기

## 데이터 시각화

- AWS 툴을 활용해 데이터를 분석해 다른 세계적 질병과 경제적 위기를 COVID-19 데이터랑 비교하기
- AWS 서비스의 시각화 기능을 통해 Interactive 대시보드를 구축하고 Frontend 애플리케이션과 원활하게 결합해 웹사이트를 배포하기

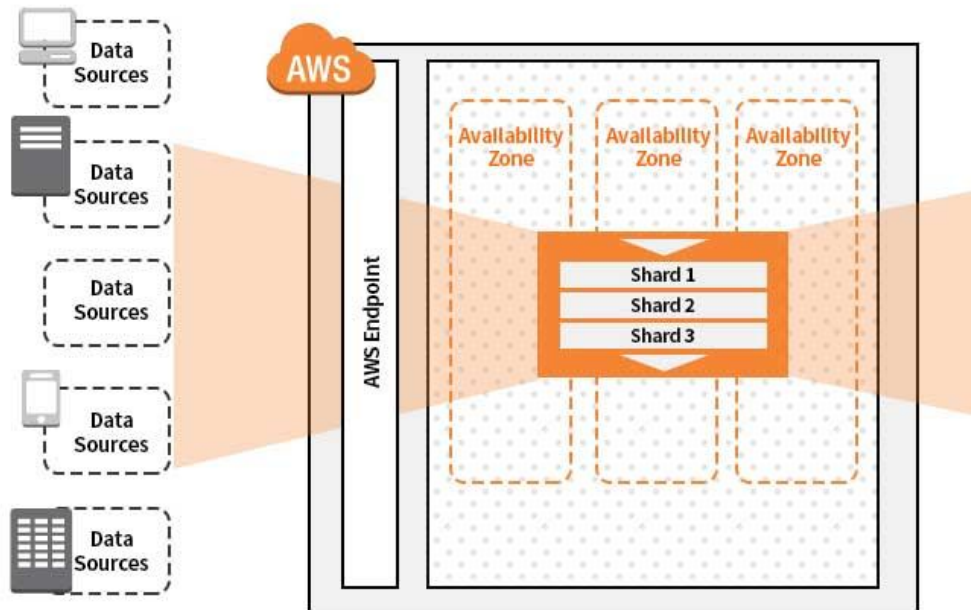


# Specifications (1)데이터 수집 및 이동

---

## Challenge

- Data Source에 따른 효율적인 데이터 수집 및 이동 방법 확인
- Static Data vs Streaming Data의 수집 및 이동
- AWS CLI, SDK로 데이터 수집 및 이동을 자동화 시키는 방법 고안



## Data Source

- 시기 : 2000년~현재
- NASDAQ / S&P 500 / 다우 / Nikkei225 / Shanghai / CSI 300 / FTSE 100 / KOSPI 200

- 시기 : 1955~2021년 8월
- OECD회원국

- 시기: 2013년~현재 (매일 업데이트)
- scope: 전세계 (132 나라)

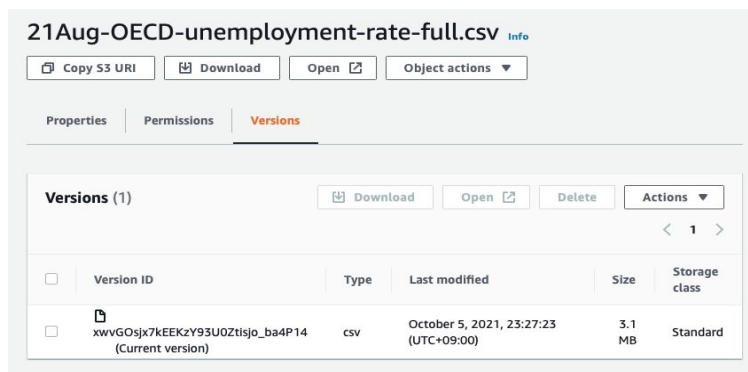
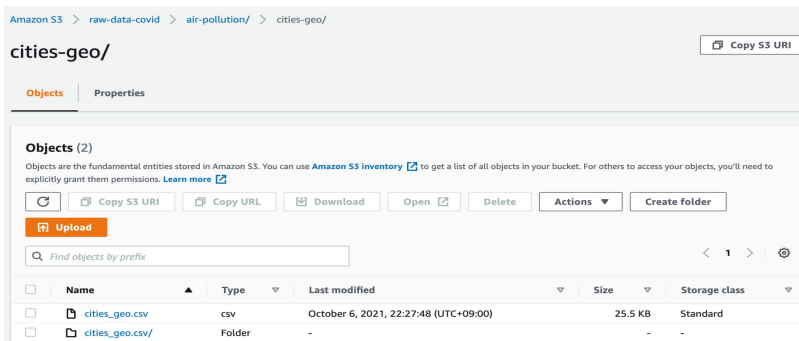


# Specifications (1)데이터 수집 및 이동

## Static Data Source → S3

AWS console – S3 raw bucket

- 버킷에 저장 후 S3 데이터 레이크 스토리지 위치 등록



boto3 library (python의 AWS sdk)

```
2 import boto3
3 import requests
4 import json
5
6 # Retrieve data from source (api) and save to local file
7 def get_file_api_local_save(api, file=None):
8     """
9     Example:
10     """
11     response = requests.get("https://aqicn.org/data-platform/covid19/airquality-covid19-cities.json")
12     if file is None:
13         file = open('save_local.csv')
14     data = response.text
15     parsed = json.loads(data)
16     file.write(parsed)
17     file.close()
18     return file
19
20 # return saved file pointer
21 pass
22
23 # Send file to s3 bucket
24 def send_file_to_s3(filePath, bucket):
25     """
26     """
27     client = boto3.client(
28         's3',
29         aws_access_key_id = 'AK*****XC',
30         aws_secret_access_key = 'RONA2y0Uws7QZ*****0Lv+DnHy',
31         region_name = 'ap-northeast-2'
32     )
33     clientResponse = client.Bucket(bucket).upload_file(filePath)
34     """
35     # return true if data is file in filePath is safely saved in s3Path else false
36     pass
```

# Specifications (2)데이터 처리

## Challenge

- File format, Table Schema 등 다양한 형태의 데이터를 SQL, Python을 이용해 여러 형태로 가공
- 필요없는 Column 삭제, 중복 데이터 제거의 자동화 등 데이터 처리를 위한 Script 작성  
→ vanilla python, spark (python or scala) 활용

## Tools

- Python으로 데이터 가공 후 S3로 이동 → 속도 문제 때문에 작은 데이터 프로세싱에 사용
- Pyspark → 큰 데이터 프로세싱에 사용
- AWS Glue의 Crawling을 이용한 Metadata 추출

```
1 # Format cities_geo
2
3 import requests
4 import json
5 import pandas as pd
6
7 response = requests.get("https://aqicn.org/data-platform/covid19/airquality-covid19-cities.json")
8 text = json.loads(response.text)
9 data = text['data']
10
11 df = pd.DataFrame(data)
12 places = df['Place']
13 places_list = places.values.tolist()
14
15 df = pd.DataFrame(places_list)
16 df = df.rename(columns={"name": "city"})
17 df = df.apply(lambda col: [col[0][0], col[0][1], col[1], col[2], col[3], col[4]], axis=1, result_type="expand")
18 df = df.rename(mapper={0: "latitude", 1: "longitude", 2: "feature", 3: "city", 4: "country", 5: "population"}, axis=1)
19
20 df.to_csv('cities_geo.csv', header=False, index=False)
```

Jobs A job is your business logic required to perform extract, transform and load (ETL) work. Job runs are initiated by triggers which can be scheduled or driven by events.

User preferences

Add job Action Filter by tags and attributes Showing: 1 - 3

Name	Type	ETL language	Script location	Last modified	Job bookmark
<input type="checkbox"/> clean-ap-2020-python	Python shell		s3://aws-glue-...	6 October 2021 2:01 P...	Disable
<input type="checkbox"/> clean-ap-2020-spark	Spark	python	s3://aws-glue-...	6 October 2021 2:19 P...	Disable
<input checked="" type="checkbox"/> clean-cities-geo-csv	Spark	python	s3://aws-glue-...	6 October 2021 4:39 P...	Disable

History Details Script Metrics

```
1 import sys
2 from awslogs.transforms import *
3 from awslogs.utils import getResolvedOptions
4 from pyspark.context import SparkContext
5 from awslogs.context import GlueContext
6 from awslogs.job import Job
7
8 ## @params: [JOB_NAME]
9 args = getResolvedOptions(sys.argv, ['JOB_NAME'])
10
11 sc = SparkContext()
12 glueContext = GlueContext(sc)
13 spark = glueContext.spark_session
14 job = Job(glueContext)
15 job.init(args['JOB_NAME'], args)
16 ## @type: DataSource
17 ## @name: DataSource = "aws-logs" table name = "aws-logs" schema = "aws-logs"
```

Edit script



# Specifications (2)데이터 처리

## Glue Crawler

- AWS Lake Formation → Crawlers → Ingest database & Crawl Data
- Crawlers를 만들어 메타 데이터 캡처 자동화 → 카탈로그 엔터티 생성

### AWS Glue

#### Data catalog

##### Databases

##### Tables

##### Connections

##### Crawlers

##### Classifiers

##### Schema registries

##### Schemas

##### Settings

##### ETL

##### AWS Glue Studio New

##### Blueprints

##### Workflows

##### Jobs

##### ML Transforms

##### Triggers

##### Dev endpoints

##### Notebooks

##### Security

##### Security configurations

##### Tutorials

##### Add crawler

Jobs A job is your business logic required to perform extract, transform and load (ETL) work. Job runs are initiated by triggers which can be scheduled or driven by events.

[User preferences](#)

<input checked="" type="checkbox"/>	Name	Type	ETL language	Script location	Last modified	Job bookmark
<input checked="" type="checkbox"/>	ingest-Parquet	Spark	python	s3://aws-glue...	6 October 2021 2:29 A...	Disable

##### History

##### Details

##### Script

##### Metrics

```
1 import sys
2 from awsglue.transforms import *
3 from awsglue.utils import getResolvedOptions
4 from pyspark.context import SparkContext
5 from awsglue.context import GlueContext
6 from awsglue.job import Job
7
8 # @params: [JOB_NAME]
9 args = getResolvedOptions(sys.argv, ['JOB_NAME'])
10
11 sc = SparkContext()
12 glueContext = GlueContext(sc)
13 spark = glueContext.spark_session
14 job = Job(glueContext)
15 job.init(args['JOB_NAME'], args)
16
17 # @type: DataSource
18 # @args: [database = "ingest", table_name = "batch_person", transformation_ctx = "datasource0"]
19 # @return: datasource0
20 # @inputs: []
21 datasource0 = glueContext.create_dynamic_frame.from_catalog(database = "ingest", table_name = "batch_person", transformation_ctx = "datasource0")
22
23 # @type: ApplyMapping
24 # @args: [mapping = [{"location", "string", "location", "string"}, {"indicator", "string", "indicator", "string"}, {"subject", "string", "subject", "string"}]]
25 # @return: applymapping1
26 # @inputs: [frame = datasource0]
27 applymapping1 = ApplyMapping.apply(frame = datasource0, mappings = [{"location", "string", "location", "string"}, {"indicator", "string", "indicator", "string"}, {"subject", "string", "subject", "string"}])
28
29 # @type: DataSink
30 # @args: [connection_type = "s3", connection_options = {"path": "s3://rachel-datalake-cleaned/Parquet-verified"}, format = "json", transformation_ctx = "datasink2"]
31 # @return: datasink2
32 # @inputs: [frame = applymapping1]
```

[Edit script](#)

# Specifications (2)데이터 처리

## AWS Glue

- AWS Glue용 IAM 역할 생성

## Glue ETL Job

- raw data를 파싱 → 효율적인 columnar data인 'parquet' 포맷으로 변환 → cleaned S3 bucket에 저장
- 데이터베이스 테이블 생성 확인

Job: ingest-Parquet

Action Save Run job

Generate diagram

Insert template at cursor Source Target Target Location Transform Spigot

Database Name ingest  
Table Name batch\_person

Transform Name ApplyMapping

Path s3://rachel-datalake-cleaned/Parquet-verified

```
1 import sys
2 from aws glue.transforms import *
3 from aws glue.utils import getResolvedOptions
4 from pyspark.context import SparkContext
5 from aws glue.context import GlueContext
6 from aws glue.job import Job
7
8 ## @param: [JOB_NAME]
9 args = getResolvedOptions(sys.argv, ['JOB_NAME'])
10
11 sc = SparkContext()
12 glueContext = GlueContext(sc)
13 spark = glueContext.spark_session
14 job = Job(glueContext)
15 job.init(args['JOB_NAME'], args)
16
17 ## @type: DataSource
18 ## @args: [database = 'ingest', table_name = 'batch_person', transformation_ctx = 'datasource0']
19 ## @return: DataSource
20 ## @inputs: []
21
22 datasource0 = glueContext.create_dynamic_frame_from_catalog(database = 'ingest', table_name = 'batch_person', transformation_ctx = 'datasource0')
23
24 ## @type: ApplyMapping
25 ## @args: [mapping = [{"location": "string", "location": "string"}, {"indicator": "string", "indicator": "string"}], frame = datasource0]
26 ## @return: ApplyMapping
27
28 applymapping1 = ApplyMapping.apply(frame = datasource0, mappings = [{"location": "string", "location": "string"}, {"indicator": "string", "indicator": "string"}])
29
30 ## @type: DataSink
31 ## @args: [path = 's3://rachel-datalake-cleaned/Parquet-verified', format = 'parquet', partition_keys = []]
32 ## @return: DataSink
33
34 sink = DataSink.apply(path = 's3://rachel-datalake-cleaned/Parquet-verified', format = 'parquet', partition_keys = [])
35 sink.write(applymapping1)
```

Logs Schema

## Schema

	Column name	Data type	Partition key	Comment
1	location	string		
2	indicator	string		
3	subject	string		
4	measure	string		
5	frequency	string		
6	time	string		
7	value	double		
8	flag codes	string		

# Specifications (3)데이터 분석

## Challenge

- AWS EMR (Hadoop, Spark 등)을 이용해서 대규모 데이터 처리 → 서로 다른 결과 도출
- AI 툴인 AWS Sagemaker를 활용한 데이터 상관관계 및 미래 예측
- SQL 사용 → Amazon Athena를 이용해 S3에 저장된 데이터를 query해 원하는 종류와 기준의 데이터만 추출
- 자주 사용되는 작업을 효과적으로 처리할 수 있는 쿼리 생성 & 저장

The screenshot displays the Amazon Athena console interface. On the left, the 'Data' sidebar shows the 'Data Source' as 'AwsDataCatalog' and the 'Database' as 'ap-cities-geo'. Below this, 'Tables and views' are listed, including 'cities\_geo\_csv' and 'cities\_geo\_json'. The main panel shows a SQL query being executed:

```
1 SELECT country, city, latitude, longitude, feature, population
2 FROM "ap-cities-geo"."cities_geo_csv"
3 where col0 >= 0
4 limit 10;
```

The query status is 'Completed'. The execution details show: 'Time in queue: 0.142 sec', 'Run time: 0.552 sec', and 'Data scanned: 28.07 KB'. The results are displayed in a table with 4 columns: 'country', 'city', 'latitude', and 'longitude'.

country	city	latitude	longitude
MX	"Ecatepec de Morelos"	19.60492	-99.06064
MX	Guadalajara	20.66682	-103.39182
MX	Cuernavaca	18.9261	-99.23075

# Specifications (3)데이터 분석

## Amazon Athena SQL Query

- SQL문을 사용하여 쿼리 진행
- 쿼리 결과를 S3 curated 버킷 AthenaResults/ 폴더에서 저장

SQL 함수를 이용해 특정 테이블 열을 합계

Query 2

```
1 select
2   unemployment.location,
3   -- aggregate fields and get average value
4   avg(unemployment.value) as "Average Index"
5 from "verified"."parquet_verified" as unemployment
6 group by
7   unemployment.location
```

Ln 4, Col 47

Run Cancel Save as Clear Create ▾

Completed Time in queue: 0.119 sec Run time: 1.216 sec Data scanned: 8.48 MB

Results (42) Download results

Search rows

location	Average Index
LTU	10.789457116161602
FRA	9.912914962906521
JPN	2.720287243533484
NLD	5.957355172764229
EST	8.969936221690585

SQL ORDER BY 데이터 테이블 조회 정렬

Query 2

```
1 select
2   unemployment.time as Date,
3   unemployment.value as "Unemployment Index"
4 from "verified"."parquet_verified" as unemployment
5   where unemployment.location='KOR'
6   -- arrange in ascending time
7 order by time asc
```

Ln 7, Col 18

Run Cancel Save as Clear Create ▾

Completed Time in queue: 0.119 sec Run time: 1.596 sec Data scanned: 8.48 MB

Results (500+) Download results

Search rows

Date	Unemployment Index
1989	1.866859
1989	3.068215
1989	2.583333
1989-01	3.171043
1989-01	1.0662



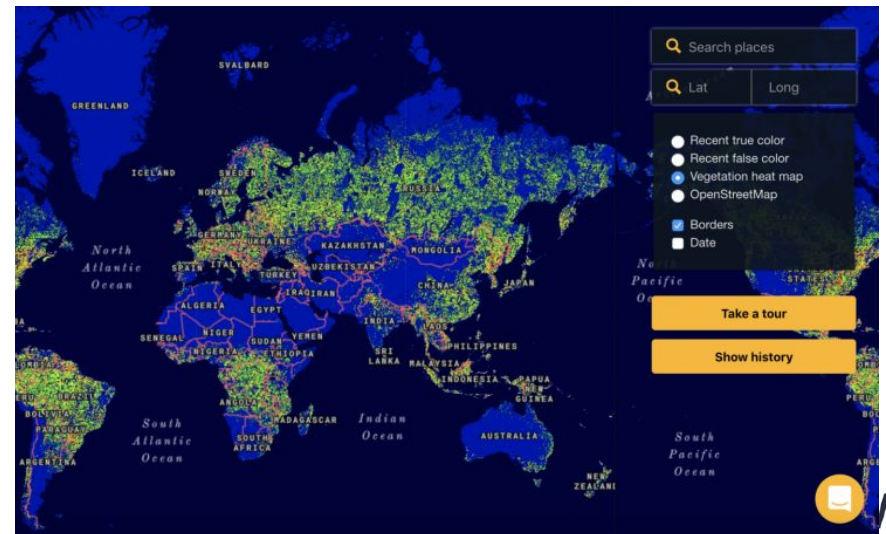
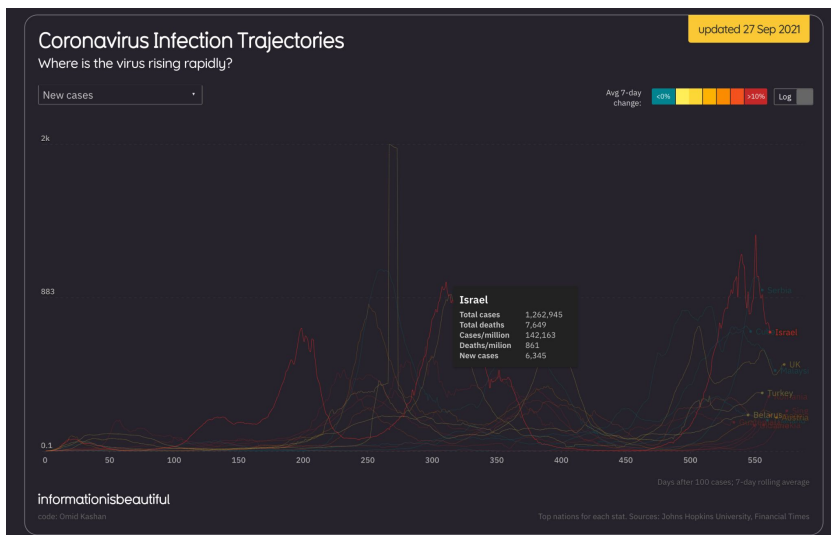
# Specifications (4)데이터 시각화

## Challenges

- Amazon QuickSight를 이용해 최종 결과를 도출, QuickSight 내에서 여러 형태로 데이터를 볼 수 있도록 customize
- Interactive하게 여러 데이터의 관계를 보여주기 위해 python libraries (folium, bokeh) 활용

## 결과물 : Interactive Website

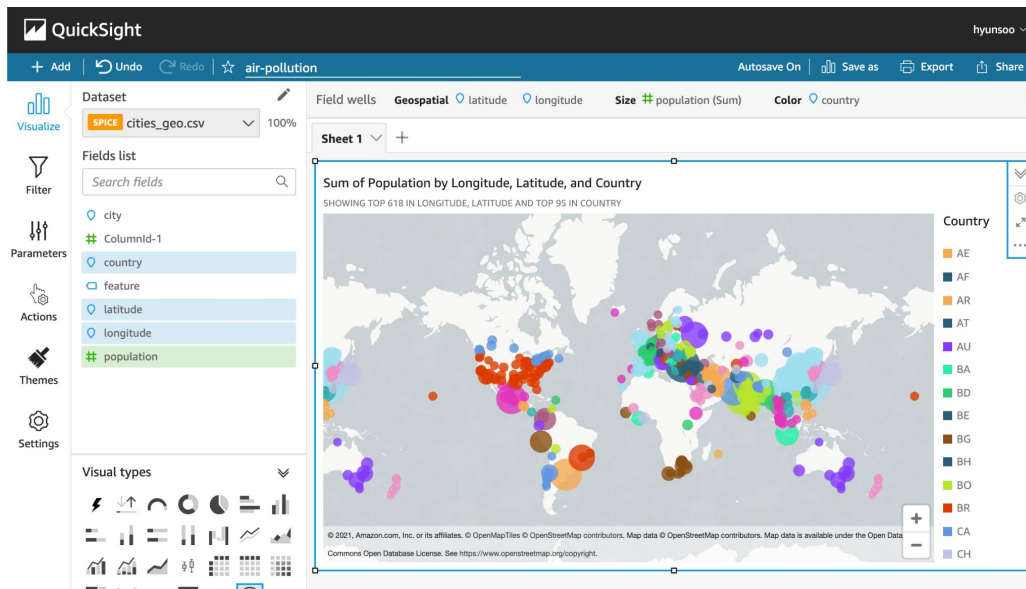
→ 여러가지 그래프가 보여지고, 그래프를 클릭하면 세부 정보를 더 볼 수 있는 형태



# Specifications (4)데이터 시각화

## 1. Heat Map

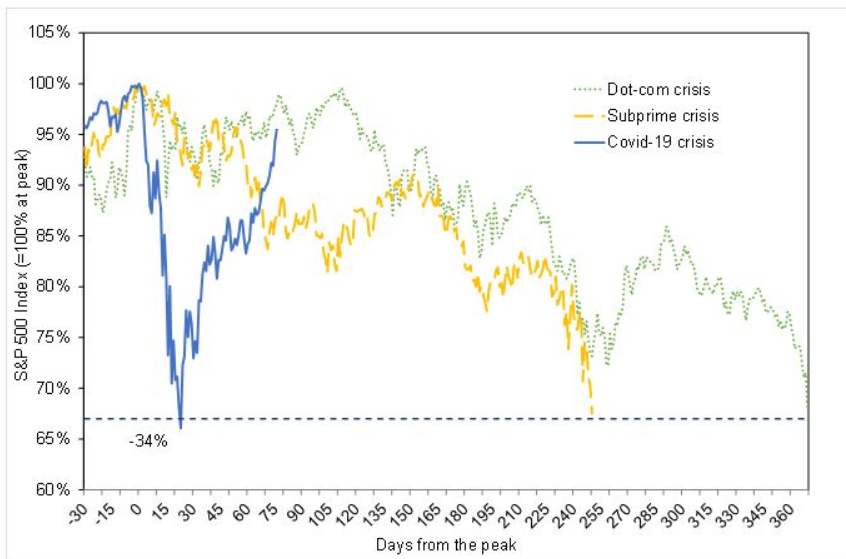
- Interactive Feature
  - heatmap에 air quality 변화를 시간에 따라 변화 시켜서 코로나 전후 비교
  - Geo heatmap, slider로 시간 조정하면 시간에 맞는 air quality를 세계 지도에 표시
- Prediction
  - 코로나가 없었을 경우의 air quality 추정



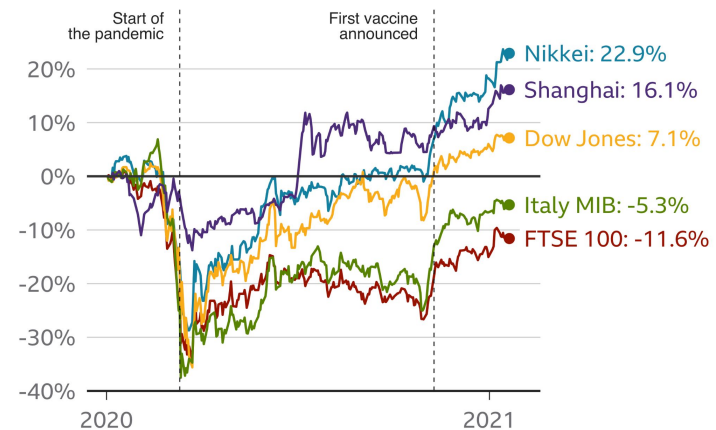
# Specifications (4)데이터 시각화

## 2. Line Graph

- Interactive Feature
  - 주가 지수 그래프를 클릭하면 sector 별 변동량 그래프 노출
  - 그래프 위 hover 시 major event 관련 정보 및 그래프 상세보기 가능
  - 선택한 국가별 지수 그래프 비교
- Data Relation
  - major event의 지속일과 주가 지수 관계



### The impact of coronavirus on stock markets since the start of the outbreak



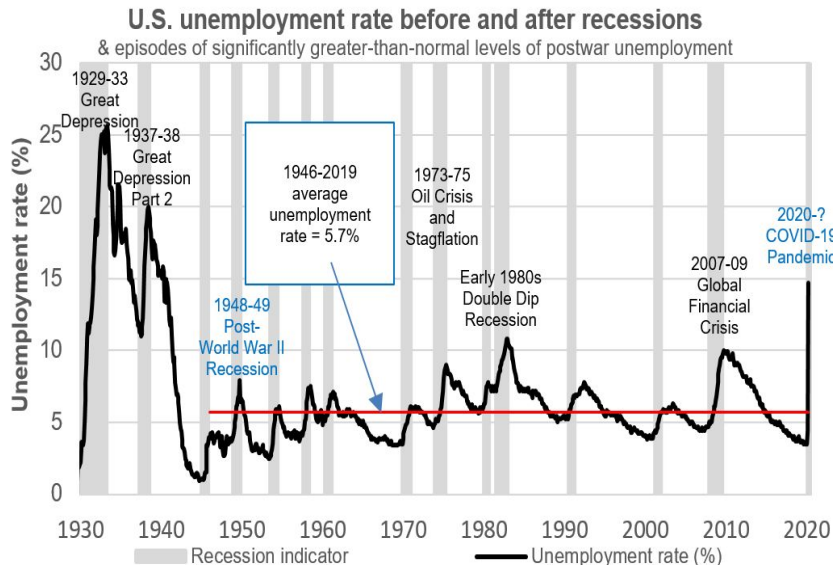
Source: Bloomberg, 24 January 2021, 00:01 GMT

BBC

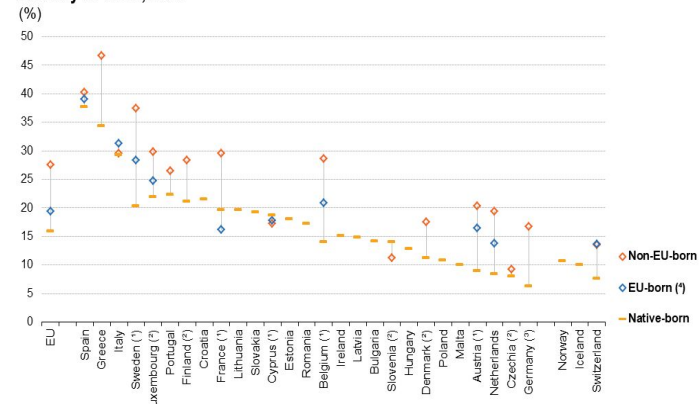
# Specifications (4)데이터 시각화

## 3. Line + Stacked Bar Chart

- Interactive Feature
  - 관심국가 highlight, 외에 나머지 국가 hover로 확인 가능
  - 분산, 가장 높은 수치의 데이터, ATH 확인
  - 국가 선택 및 hover 시 detailed index number 확인 가능
- Prediction
  - 2021 Q4 ~ 2022 Q4 실업률 예측



**Youth unemployment rates for the population aged 15-24 years, by country of birth, 2020**



# Development Environment

---

## Data Gathering

- Crawling (python)
- API

## Data Processing and Analysis

- AWS services via CLI and console
- boto 3 library (python)
- Pyspark
- AWS Glue Crawling / Data Catalog
- AWS EMR (Hadoop, Spark 등)
- AWS Athena + SQL
- AWS Sagemaker
- vanilla python

## Data Visualization

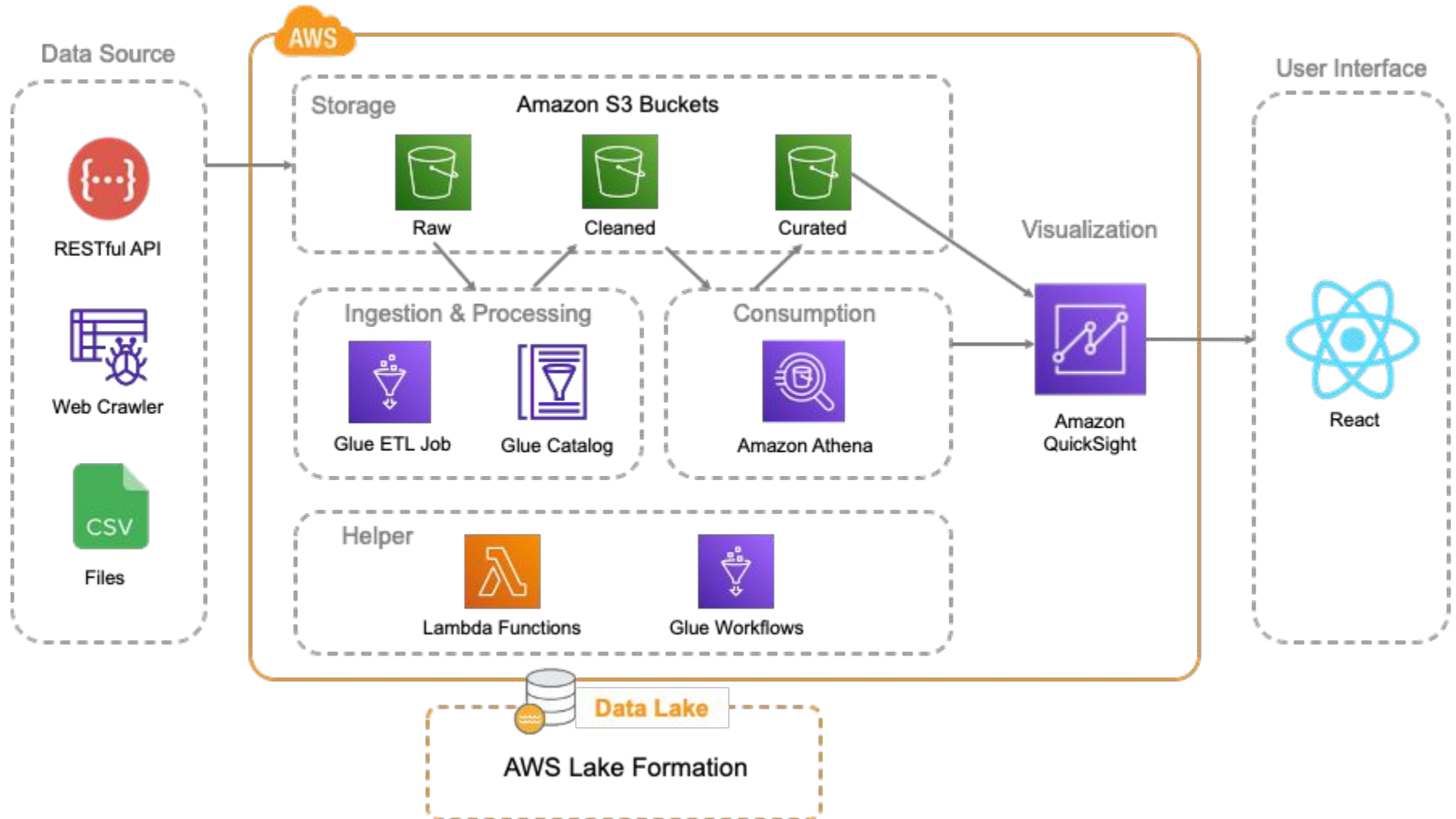
- AWS Quicksight
- python libraries (folium, bokeh)

## Web Framework

- React (javascript)

and more...

# Architecture



# Current Status

---

## AWS 사용자 계정 관리

- AWS 계정 루트 사용자 생성하고 역할을 만들어 IAM 사용자에게 권한 부여

## AWS 데이터 레이크

- COVID-19 데이터 레이크용 Amazon S3 버킷 생성:
  - Raw 데이터, cleaned 데이터, curated 데이터
- AWS Lake Formation 구조 잡기:
  - COVID-19 데이터 레이크
  - 데이터 레이크 루트 위치로 S3경로를 등록
  - 데이터베이스를 만들어 메타데이터 저장

## 간단한 데이터 수집 및 프로세싱

- python, AWS console 활용해 데이터 S3에 업로드
- Glue 활용한 간단한 데이터 프로세싱

# Division and Assignment of Work

---

## 주제별 데이터 전처리와 분석

1. 대기오염: 김현수
2. 실업률: 레이첼 간
3. 주가: 나혜주



# Further Plans

---

## 향후 계획

- 웹크롤링을 통해서 대오염도, 주가 및 실업률에 관한 데이터를 수집하여 데이터베이스에 저장
- Amazon Athena로 저장된 데이터 분석
- AWS QuickSight로 데이터를 시각화하여 COVID-19 확산과 각 주제별 데이터 간의 상관관계 나타내기
- React 웹 애플리케이션 개발
- AWS QuickSight 임베딩 SDK를 사용하여 React 앱에서 대시보드 구현

# Schedule

내용	9월		10월				11월				12월		
	3주	4주	1주	2주	3주	4주	1주	2주	3주	4주	1주	2주	3주
프로젝트 소개 및 자료조사													
주제확정 및 Infrastructure 구축													
프로젝트 스펙 발표													
주제별 데이터 수집 및 분석													
프로젝트 중간 발표													
데이터 시각화													
React 앱 개발 및 대시보드 구현													
Optimization													
프로젝트 최종 발표													

감사합니다