## I.     Setting up a Data Engineering Infrastructure: A step-by-step guide

**Installations and configurations on windows 10**

   a. **Apache Nifi**
   b. **Apache Airflow**
   c. **Elasticsearch**
   d. **Kibana**
   e. **PostgreSQL**
   f. **pgAdmin 4**

   1. **Installing and configuring Apache Nifi**

Curl https://mirrors.estointernet.in/apache/nifi/1.12.1/nifi-1.12.1-bin.tar.gz #cmd

https://www.clearpeaks.com/installing-apache-nifi-on-windows/ #for window

**Install a valid Java Development Kit**

Download JDK from https://www.oracle.com/ae/java/technologies/javase/javase8u211-later-archive-downloads.html

Installation

Open ~ **edit the system variable environment** ~ from the search bar ~ click on **environment Variable ~** select Path **c://windows \system32 .. .. wwbem** then click **edit ~ new ~ then add java path then okay – open cmd and type javac to confirm installation**

**Download**: JDK Development Kit ~
https://www.oracle.com/ke/java/technologies/downloads/#java22

**Starting Nifi**

Download Nifi from https://nifi.apache.org/download/

Extract the compressed Nifi file, then go to bin, then click or run nifi.bat file it should start nifi environment ... then access the **Nifi UI** from *localhost:8443/nifi/*

## Successful after many trials

Generated Username [2dc3de96-af49-46b8-ab1d-b2c52483e711]

Generated Password [t9IzC9H3BESmG7Pl6IyvvqBEsKHg+n8K]

**Problem at a glance ~** latest version did not process an ID



Figure 1: Nifi Issues

## Processed ID

A Kibe Creation

## Successful: https://localhost:8443/nifi/

## 2. Installing and configuring Apache Airflow

**Tools needed**

- Docker
- Visual studio Code

**-Download and install docker – follow the serious SQL procedure (very hectic and stubborn this one)**

**- Set up Apache Airflow**

**Save file; https://airflow.apache.org/docs/apache-airflow/2.5.1/docker-compose.yaml as a .yaml file otherwise it will not work. You can use VS code to edit. Save file in a folder then open said folder from VS code environment.**

**Create a .env file then paste below code and save.**

AIRFLOW_IMAGE_NAME= apache/airflow:2.4.2

AIRFLOW_UID=50000

Open a terminal and run ***docker compose up -d***, these should pull all relevant resources and create an airflow container on docker, executable through ***Localhost:8080/Login/***

**Problem**



Figure 2: Airflow Issues

### 3. Installing and configuring Elasticsearch

curl https://artifacts.elastic.co/downloads/ Elasticsearch/elasticsearch-7.6.0-darwin-x86_64.tar.gz --output elasticsearch.tar.gz **#cmd**

**Download** Elasticsearch zip file from: https://www.elastic.co/downloads/elasticsearch then extract the file or https://www.elastic.co/downloads/past-releases#elasticsearch

Then go bin in the folder look for Elasticsearch type window batch file, then run a command line from there with (**elasticsearch.bat**), it will now create our application, accessible on a URL on port 9200 at https://localhost:9200 or 9300

Play around with ***Elasticsearch-service.bat remove*** to try different versions if a process fails

sc delete elasticsearch-service-x86_64 **or** Control Panel\System and Security\Administrative Tools **or** services.msc ~ on win + r **or** on #cmd

**Successful**



**Setting up a password**

Run CMD as an ADMIN, then move to the bin with Elasticsearch batch file, copy the path then paste on cmd e.g. cd C:\Users\ ... ... \Elasticsearch\elasticsearch-8.11.4\bin

Then, below command to reset a new password with username: elastic

```
-Elasticsearch-reset-password -u elastic --interactive
```

- Elasticsearch-reset-password -u elastic --interactive –verbose with debug

**Do not close** the previous cmd running Elasticsearch application.

**ERROR: Failed to determine the health of the cluster. , with exit code 69**

4. **Installing and configuring Kibana**

**Download Kibana:** https://www.elastic.co/downloads/kibana

Basically, our Kibana run on port 5601 at default.

Then follow below steps to launch

**Run** a CMD from the Kibana folder, basically, select the whole path on the search environment then type CMD, this should launch a working shell automatically, with path C:\Users\PC\ 'your directory' ... \kibana-8.14.0> .. **On CMD** environment Run **.\bin\kibana.ba,** this should install Kibana. On Firefox (Recommended) open **http://localhost:5601/?code=043252**

**Successful but with conflicting address issues as below**

**Figure 4a: Kibana Issues**



Figure 4b: Kibana Issues

5.  **Installing and configuring PostgreSQL**

**Download PostgreSQL:** https://www.enterprisedb.com/downloads/postgres-postgresql-downloads

Then install normally, this will take a few minutes Port = 5432, remember password = ***k#$e***

**Launches successfully as below**

6.  **Installing pgAdmin 4**

**pgAdmin 4** has been installed as a PostgreSQL component above.

**Problem**

**Asks configure a static port in one case, in another case it asks to return to default, then now launches successfully.**

**Figure 6: Postgres issue**

## Successful

## II. Reading and Writing Files

Explore the process of reading and writing CSV files using Python's built-in CSV library and pandas Dataframes. Compare the two approaches in terms of performance, ease of use, and handling of large datasets. Provide a scenario where one approach might be preferred over the other.

### a. How to read and write CSV files using Python's built-in CSV library

Sample code to **write** CSV file s in python's built in CSV library

```
[43] from faker import Faker
     import csv

     fake = Faker()
     header = ['name', 'age', 'street', 'city', 'state', 'zip', 'lng', 'lat']

     with open('data.csv', 'w', newline='') as output:
         mywriter = csv.writer(output)
         mywriter.writerow(header)
         for _ in range(1000):
             mywriter.writerow([fake.name(),fake.random_int(min=18, max=80, step=1),
                 fake.street_address(),fake.city(),fake.state(), fake.zipcode(),
                 fake.longitude(),fake.latitude()])

     print("Data generation complete! File: data.csv")

     Data generation complete! File: data.csv
```

Sample codes to **read** CSV file s in python's built in CSV library

```
Reading CSVs

import csv

with open('data.csv', newline='') as f:
    myreader = csv.DictReader(f)
    headers = myreader.fieldnames   # `next(myreader)` isn't needed with `DictReader`
    for row in myreader:
        print(row['name'])

Kyle Bennett
Jonathan Walters
Teresa Small
Breanna Evans
```

## b. How to Read and write CSVs using pandas Dataframes

Sample code to **read** CSVs using pandas Dataframes

```python
import pandas as pd
df=pd.read_csv('data.csv') #(df=pd.read_csv()('data.CSV'))wrongly written oinn the book
df.head(10)
```

| | name | age | street | city | state | zip | lng | lat |
|---|---|---|---|---|---|---|---|---|
| 0 | Katelyn Kim | 41 | 08099 Amanda Lane | Port Lisamouth | Michigan | 30630 | -59.356943 | 69.868939 |
| 1 | Steven Ferguson | 33 | 254 Gomez Isle | New Darin | Texas | 47992 | 93.826506 | 53.347920 |
| 2 | John White | 50 | 33956 Smith Mountain Suite 550 | Sharpshire | Iowa | 27073 | 41.482072 | -15.715067 |
| 3 | Maria Davis | 43 | 44520 Michael Walks | Andreashire | Tennessee | 75634 | -30.146793 | 16.119954 |
| 4 | Ronald Rodriguez | 23 | 0653 Savannah Ports | Rachelshire | Pennsylvania | 38182 | -89.041924 | -52.676035 |
| 5 | Mariah Butler | 59 | 5977 Shannon Summit | Marcusland | Nevada | 32363 | -35.253967 | 20.387420 |
| 6 | Nicholas Cole | 34 | 782 Johnson Ranch Suite 846 | East Angelaborough | Mississippi | 21857 | 117.343581 | -46.041201 |
| 7 | Ricardo Hernandez | 53 | 75997 Adams Terrace | Lake Susan | Alaska | 37949 | -174.563356 | -13.828719 |
| 8 | Heidi Simpson | 64 | 6768 Amanda Mills | Timothyville | North Carolina | 83956 | -14.817070 | -31.078382 |

Sample code to **write** CSVs using pandas Dataframes

### Create a DataFrame in Python

```python
[48] data={'Name':['Paul','Bob','Susan','Yolanda'],
     'Age':[23,45,18,21]}
     df=pd.DataFrame(data)
     df.to_csv('fromdf.CSV',index=False)
```

```python
import pandas as pd
df=pd.read_csv('fromdf.CSV') #(df=pd.read_csv()('data.CSV'))wrongly written oinn the book
df.head(10)
```

| | Name | Age |
|---|---|---|
| 0 | Paul | 23 |
| 1 | Bob | 45 |

### c. Writing JSON with Python

Sample code to write JSON file in Python

```
[3]  from faker import Faker
     import json
     output=open('data.JSON','w')
     fake=Faker()

[4]  alldata={}
     alldata['records']=[]

     for x in range(1000):
       data={"name":fake.name(),"age":fake.random_int
       (min=18, max=80, step=1),
       "street":fake.street_address(),
       "city":fake.city(),"state":fake.state(),
       "zip":fake.zipcode(),
       "lng":float(fake.longitude()),
       "lat":float(fake.latitude())}
       alldata['records'].append(data)

[7]  json.dump(alldata,output)
```

Sample code to read JSON files in Python

```
# Open the file in read mode ('r')
with open("data.JSON", "r") as f:
    data = json.load(f)
    print(data['records'][0])
    print(data['records'][0]['name'])
```
```
{'name': 'Holly Potter', 'age': 66, 'street': '6111 Conrad Light Apt. 806', 'city': 'West Luisberg', 'state': 'Alabama', 'zip':
Holly Potter
```

### d. Reading and writing JSON with Dataframes

Sample code read Json files

```python
import pandas as pd
df=pd.read_json('data.JSON')
df.head(10)
```

| | records |
|---|---|
| 0 | {'name': 'Holly Potter', 'age': 66, 'street': ... |
| 1 | {'name': 'Anna Navarro', 'age': 52, 'street': ... |
| 2 | {'name': 'Leah Rodriguez', 'age': 67, 'street'... |
| 3 | {'name': 'Melissa Hill', 'age': 64, 'street': ... |
| 4 | {'name': 'Joshua Thompson', 'age': 34, 'street... |
| 5 | {'name': 'Andrew Dawson', 'age': 35, 'street':... |
| 6 | {'name': 'Joshua Smith', 'age': 71, 'street': ... |
| 7 | {name': 'Kimberly Smith' 'age': 57 'street' |

Pass the orient parameter, which determines the format of the JSON that is returned

Sample results as below

```python
df.head(2).to_json(orient='records')
```

```
'[{"records":[{"name":"Holly Potter","age":66,"street":"6111 Conrad Light Apt. 806","city":"West Luisberg
t":-53.21111},{"name":"Anna Navarro","age":52,"street":"25638 Dawson Wall Suite 288","city":"Monroecheste
at":28.489776},{"name":"Leah Rodriguez","age":67,"street":"0198 Parker Island Suite 892","city":"Russells
09288,"lat":54.556264},{"name":"Melissa Hill","age":64,"street":"4751 Bowers Oval Suite 136","city":"Lisa
8,"lat":32.8358755},{"name":"Joshua Thompson","age":34,"street":"96612 Eddie Hill","city":"Ronaldton","st
1.05052},{"name":"Andrew Dawson","age":35,"street":"858 James Hills Suite 443","city":"Mcdonaldside","sta
-83.186551},{"name":"Joshua Smith","age...'
```
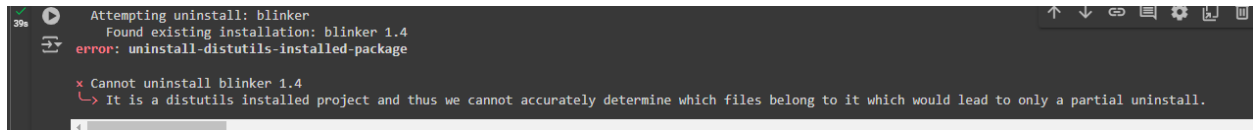
Default Json file looks like below

```python
df.head(2).to_json()
```

```
'{"records":{"0":[{"name":"Holly Potter","age":66,"street":"6111 Conrad Light Apt. 806","city":"West Luisberg","s
7,"lat":-53.21111},{"name":"Anna Navarro","age":52,"street":"25638 Dawson Wall Suite 288","city":"Monroechester",
35,"lat":28.489776},{"name":"Leah Rodriguez","age":67,"street":"0198 Parker Island Suite 892","city":"Russellshir
106.209288,"lat":54.556264},{"name":"Melissa Hill","age":64,"street":"4751 Bowers Oval Suite 136","city":"Lisamou
9098,"lat":32.8358755},{"name":"Joshua Thompson","age":34,"street":"96612 Eddie Hill","city":"Ronaldton","state":
t":31.05052},{"name":"Andrew Dawson","age":35,"street":"858 James Hills Suite 443","city":"Mcdonaldside","state":
at":-83.186551},{"name":"Joshua Smith",...'
```

## Building a CSV to a JSON data pipeline

### Wanted to install apache-airflow on my IDE but ran into below error



The system was configuring the environment in preparation for apache-airflow installation. I tried to uninstall it through conda remove blinker and still ran into an error as below.



**Conclusion**

Basically, this project requires a very robust system, for example core i7 and above, unfortunately I was using core i5. So, as I was approaching towards using the whole Data engineering infrastructure at once to perform ETL, my system was overwhelmed. However, the project is perfectly doable (Chapter 4: Working with Databases, Chapter 5: Cleaning and Transforming Data, Chapter 6: Building a 311 Data Pipeline) if and when the data engineering infrastructure is configured correctly. Continuing to struggle to make progress with my weak system is an unsatisfactory endeavor.