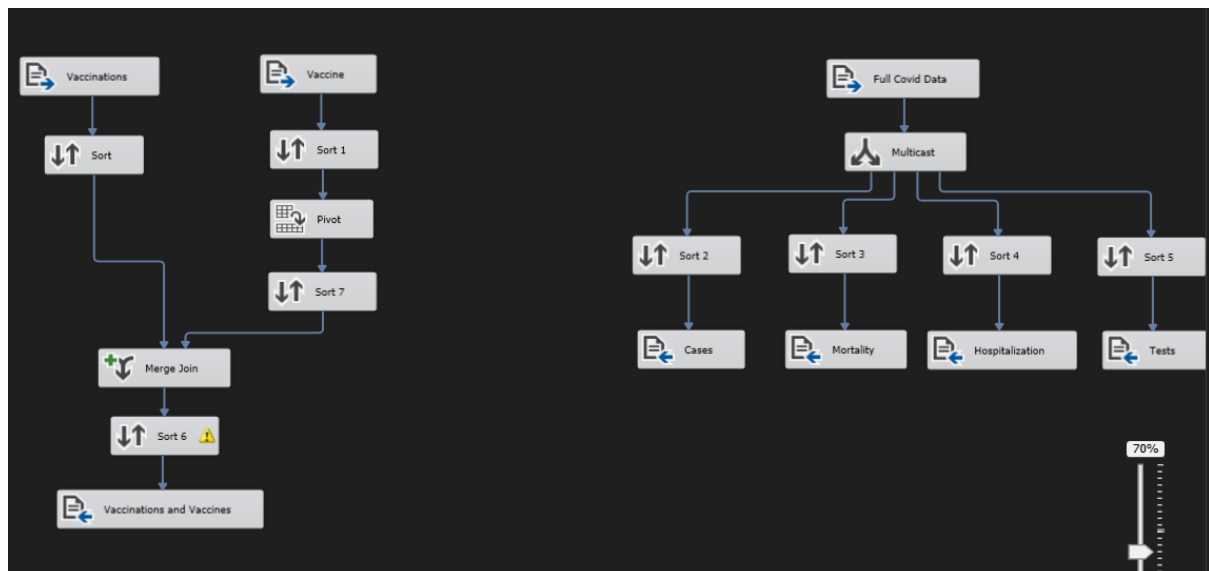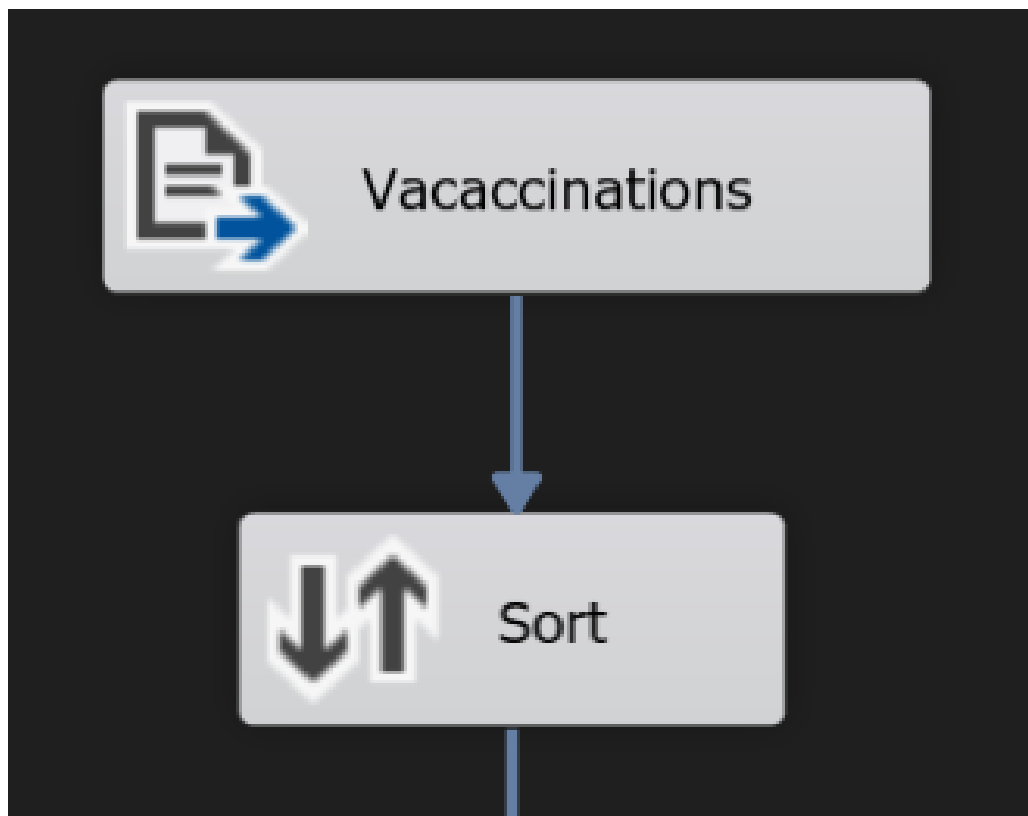## Preprocessing Stage

### Aim of the process

The full preprocessing stage is done to obtain files ready for data cleansing done with the scripts written in python. The whole data flow process at this stage is shown below:
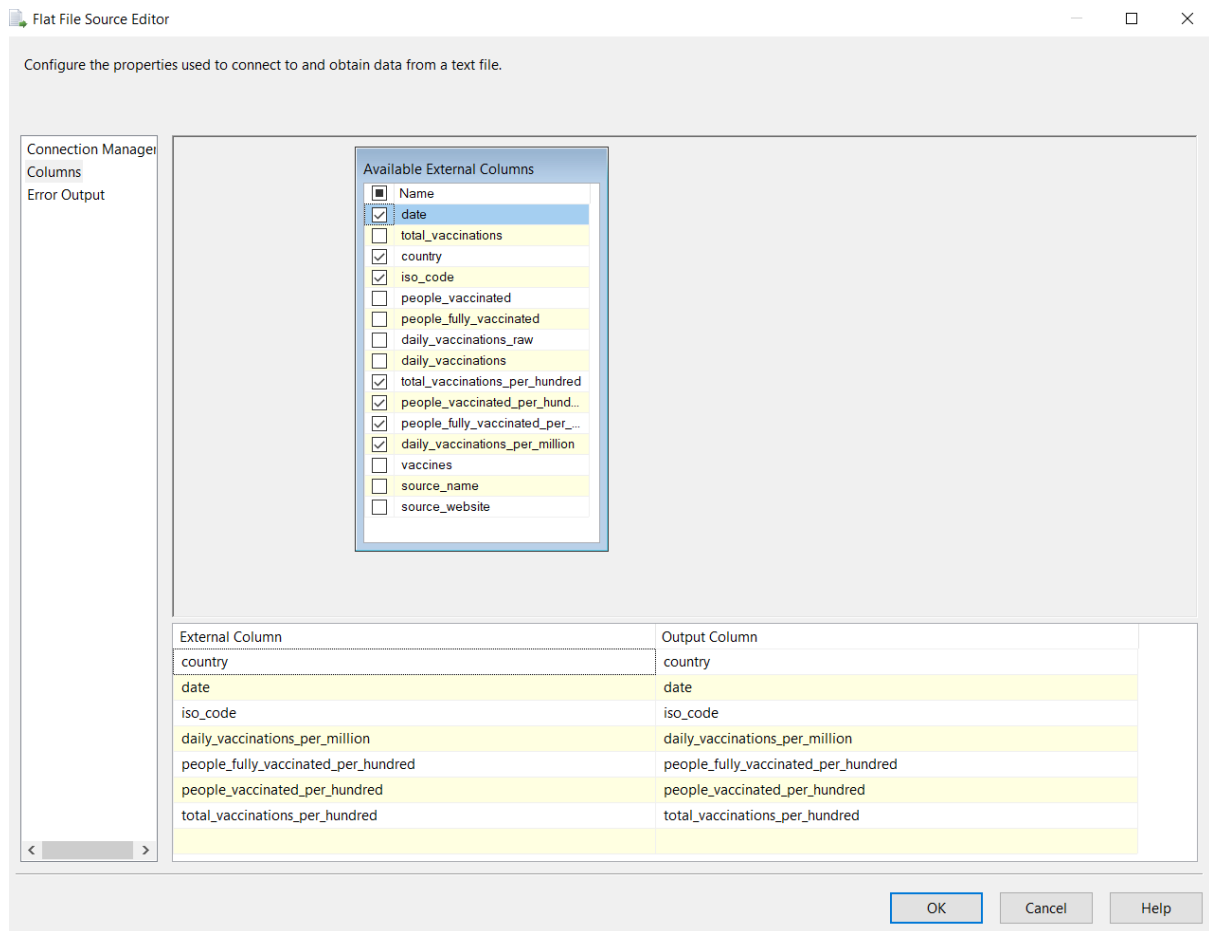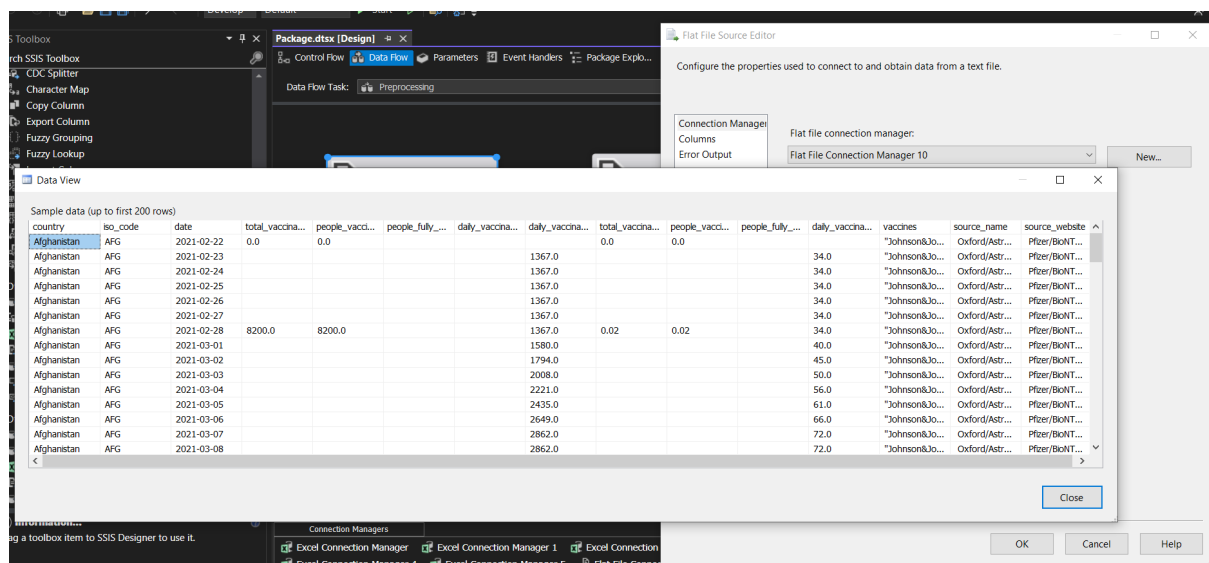


### Defining Data Sources

Defining 'Vaccinations' Flat File Source (country_vaccinations.csv) and then Sorting the data:

Selecting the appropriate columns:



Data preview (as we can see, there is a lot of missing and not continuous data, which will be cleaned also by using python scripts):

The data are sorted based on the 'country' and 'date' columns:



The same process is repeated for the 'Vaccine' (country_vaccinations_by_manufacturer.csv) and 'Full Covid Data' (cowid-covid-data-full.csv) Flat File Sources:

**Flat File Source Editor**

Configure the properties used to connect to and obtain data from a text file.

Connection Manager
**Columns**
Error Output

Available External Colu...

- ☑ Name
- ☑ location
- ☑ date
- ☑ vaccine
- ☑ total_vaccinations

| External Column | Output Column |
|---|---|
| location | location |
| date | date |
| vaccine | vaccine |
| total_vaccinations | total_vaccinations |
| | |

OK    Cancel    Help

---

**Flat File Source Editor**

Configure the properties used to connect to and obtain data from a text file.

Conne
Colum
Error

New...

**Data View**

Sample data (up to first 200 rows)

| location | date | vaccine | total_vaccinat... |
|---|---|---|---|
| Argentina | 2020-12-29 | Moderna | 2 |
| Argentina | 2020-12-29 | Oxford/Astra... | 3 |
| Argentina | 2020-12-29 | Sinopharm/B... | 1 |
| Argentina | 2020-12-29 | Sputnik V | 20481 |
| Argentina | 2020-12-30 | Moderna | 2 |
| Argentina | 2020-12-30 | Oxford/Astra... | 3 |
| Argentina | 2020-12-30 | Sinopharm/B... | 1 |
| Argentina | 2020-12-30 | Sputnik V | 40583 |
| Argentina | 2020-12-31 | Moderna | 2 |
| Argentina | 2020-12-31 | Oxford/Astra... | 3 |
| Argentina | 2020-12-31 | Sinopharm/B... | 1 |
| Argentina | 2020-12-31 | Sputnik V | 43388 |
| Argentina | 2021-01-01 | Moderna | 2 |
| Argentina | 2021-01-01 | Oxford/Astra... | 5 |
| Argentina | 2021-01-01 | Sinopharm/B... | 1 |
| Argentina | 2021-01-01 | Sputnik V | 43513 |

Close

OK    Cancel    Help

## Pivot

In the 'Vaccine' data values from the 'vaccine' column are transferred to separate columns and then sorted:



**Pivot Key:**
Values in the input data from this column will become new column names in the output

vaccine

**Set Key:**
Identifies a group of input rows that will get pivoted into one output row. The input data must be sorted on this column

date

**Pivot Value**
Values from this column will be mapped into the new pivot output columns

total_vaccinations

☐ Ignore un-matched Pivot Key values and report them after DataFlow execution

**Generate pivot output columns from values:**

Hint: choose to 'Ignore' un-matched Pivot Key values, execute this DataFlow in the debugger and copy the value list reported in the debugger's Output Window
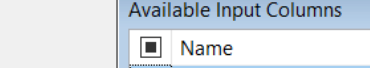
[value1],[value2],[value3]

**Existing pivoted output columns:**

cansino_total_vaccinations
covaxin_total_vaccinations
johnson&johnson_total_vaccinations
moderna_total_vaccinations
novavax_total_vaccinations
astrazeneca_total_vaccinations
pfizer_total_vaccinations
sinopharm_total_vaccinations
sinovac_total_vaccinations
sputnikv_total_vaccinations

Generate Columns Now

OK          Cancel          Help
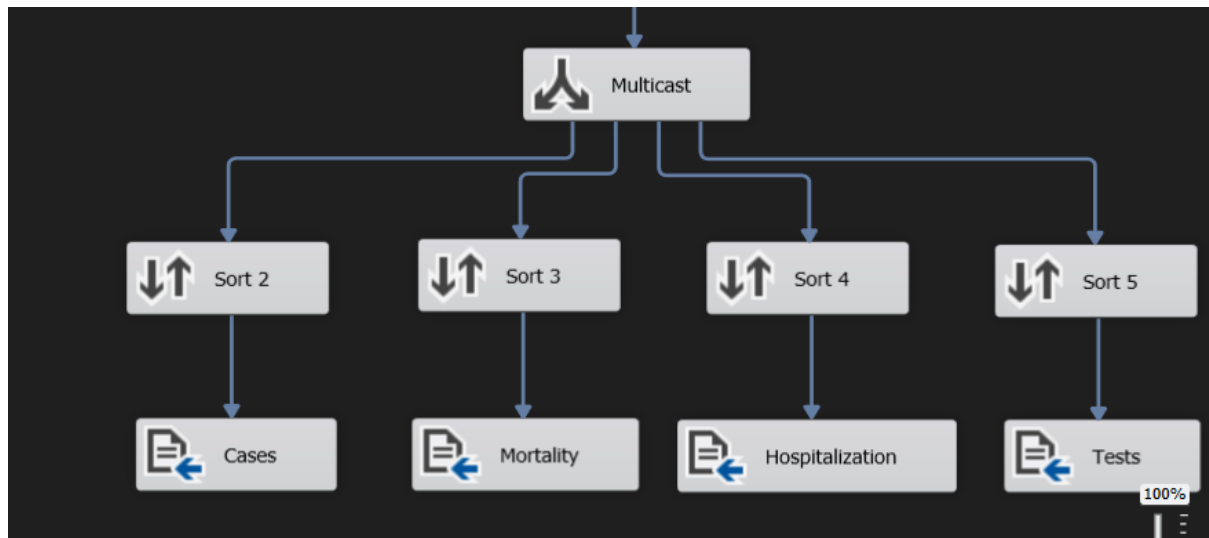
## Sort Transformation Editor

Specify the columns to sort, and set their sort type and their sort order. All nonselected columns are copied unchanged.

**Available Input Columns**

| ■ | Name | Pass Thr... |
|---|------|-------------|
| ☑ | date | ■ |
| ☐ | cansino_total_vaccinations | ☑ |
| ☐ | covaxin_total_vaccinations | ☑ |
| ☐ | johnson&johnson_total_vaccinations | ☑ |
| ☐ | moderna_total_vaccinations | ☑ |
| ☐ | novavax_total_vaccinations | ☑ |

| Input Column | Output Alias | Sort Type | Sort Order | Con |
|--------------|--------------|-----------|------------|-----|
| date | date | ascending | 2 | |
| location | location | ascending | 1 | |

☐ Remove rows with duplicate sort values

[ OK ]  [ Cancel ]  [ Help ]

### Multicast

Data from the 'Full Covid Data' are separated, sorted (location and date) and then sent to separate files, which will be processed by python scripts:

**Merge Join**

Data from the 'Vaccinations' and 'Vaccine' are merged using an inner join, then sorted and transferred to separate file: