# Task 2 : Python and Lambda
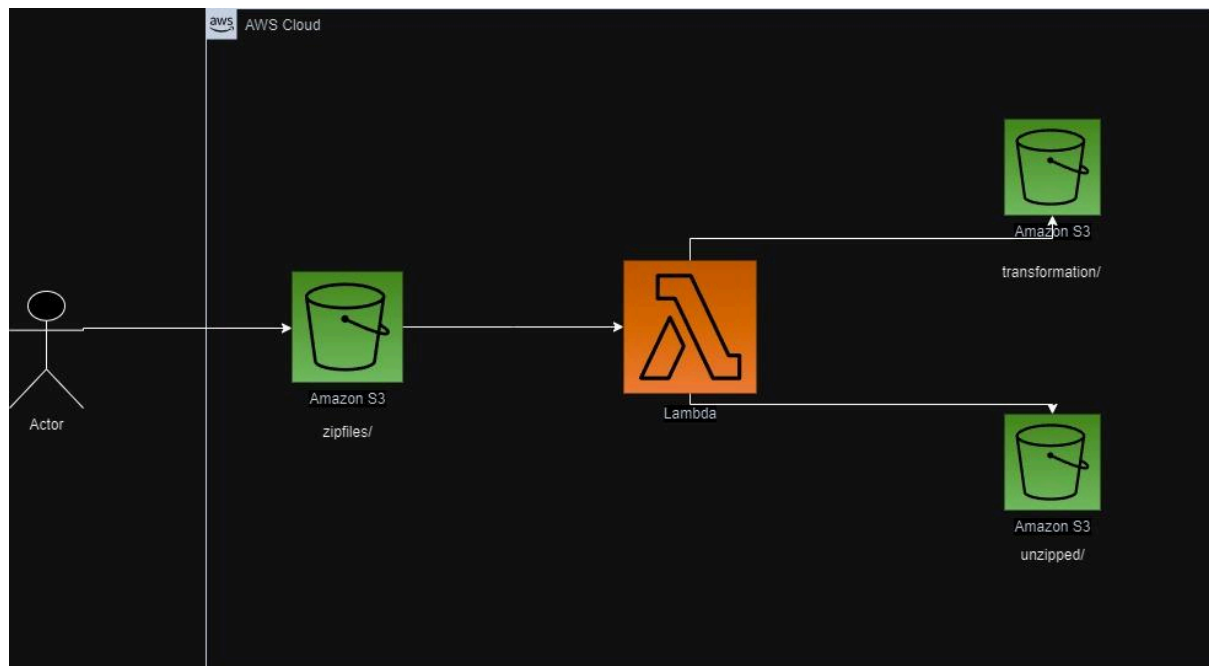
Work with RXNORM file,

1. Scrap the latest RXNORM file from NLM webpage
2. Download the latest RXNORM file with api_key
3. Create a log file for the downloaded file
4. Add header into each rff from RXNORM.xlsx
5. Add CODE_SET & VERSION_MONTH column with default values RxNorm and version month from downloaded filename
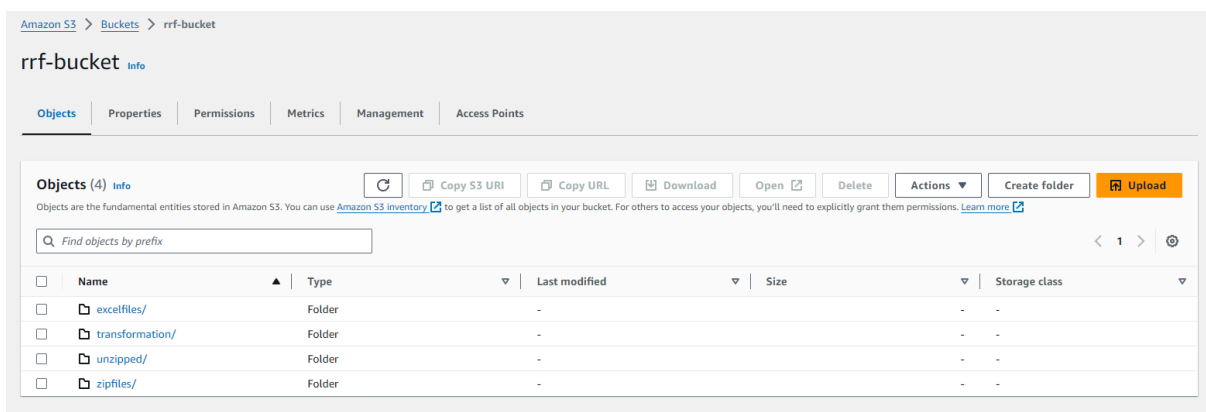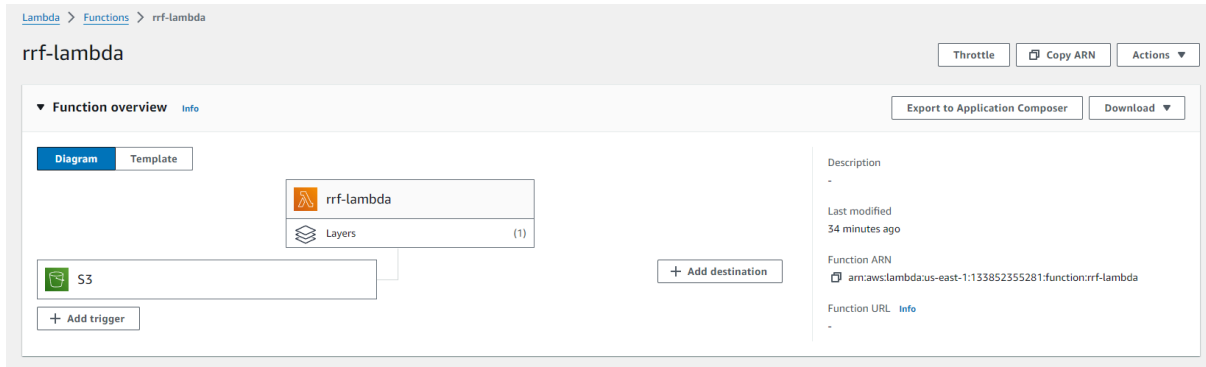6. Convert dates into YYYY-MM-DD

7. Save files as txt delimited by Pipe(|)
8. Validate row_count between original and converted files

RXNORM.xlsx

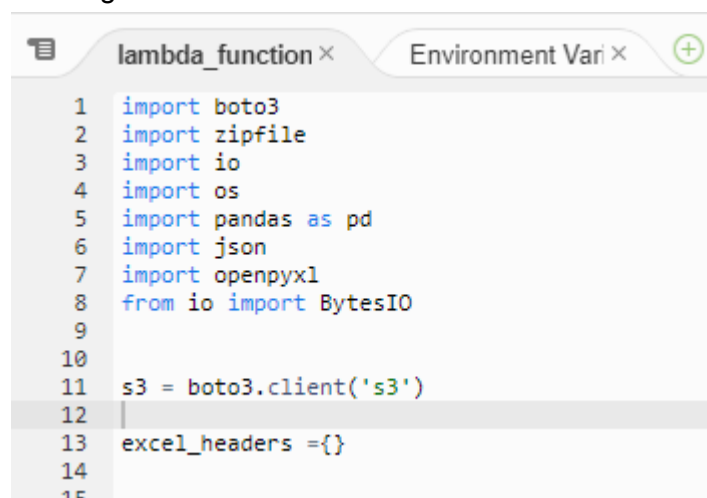**The Cloud architecture of the given project is designed as follow:**

**Step 1: Proper setup of aws was done:** lambda was created and s3 bucket was setup with four folders- zipfiles/ (to input the zipdata), excelfiles/(to input the excel file), unzipped/(to hold the unzipped files) and finally transformation/(to hold the final transformed data in .txt format)



Here, unzipped and transformation folder are created after lambda function is triggered

**Step 2: The code structure**

A. The imports and excel_headers dictionary was initiated to hold the headers to be assigned to the rrf files

```
1  import boto3
2  import zipfile
3  import io
4  import os
5  import pandas as pd
6  import json
7  import openpyxl
8  from io import BytesIO
9
10
11  s3 = boto3.client('s3')
12
13  excel_headers ={}
14
15
```

B. Function to read the excel_file and update the excel_header

```python
14
15
16  def read_excel_from_s3(bucket):
17      try:
18          folder_path = 'excelfiles/'
19          excel_file_name = 'RxNorm_Header.xlsx'
20          key = folder_path + excel_file_name
21
22          # Download the Excel file to the /tmp directory
23          local_excel_file = '/tmp/RxNorm_Header.xlsx'
24          s3.download_file(bucket, key, local_excel_file)
25
26          # Check if the Excel file exists
27          if os.path.exists(local_excel_file):
28              print(f"Excel file downloaded to: {local_excel_file}")
29
30
31
32
33          # Read the Excel file into an ExcelFile object
34          excel_file = pd.ExcelFile(local_excel_file)
35
36              # Get the sheet names
37          sheet_names = excel_file.sheet_names
38          print("Sheet names:", sheet_names)
39
40          for sheets in sheet_names:
41              # Read the data from the sheet into a DataFrame
42              sheets_data = excel_file.parse(sheets,header=None)
43              headers_data = sheets_data.iloc[:, 0].tolist()
44              excel_headers[sheets] = headers_data
45          print(f'excel_headers dictionary for sheet {sheet_names[0]}: {excel_headers[sheet_names[0]]}')
46
47
48      except Exception as e:
49          print(f"Error occurred: {e}")
```

C. Function to add Code Set and Version Month: Code set was set to RXNORM and version month was extracted from the zip file name which is RxNorm_full_02052024.zip
Hence Version Month is: 2024-05-02

```python
50
51  def code_set_and_version_month(zip_data, rrf_df):
52      try:
53          # Convert the zip data to a string to extract information
54          zip_filename = zip_data.decode('utf-8')
55          # Extract version month from the filename
56          version_month = os.path.splitext(zip_filename)[0].split('_')[-1]
57
58          # Convert version month to a more readable format
59          version_month = pd.to_datetime(version_month, format='%m%d%Y').strftime('%Y-%m-%d')
60          print(f"Version month: {version_month}")
61
62          # Add 'Code Set' and 'Version Month' columns to the DataFrame
63          rrf_df['Code Set'] = 'RxNorm'
64          rrf_df['Version Month'] = version_month
65      except Exception as e:
66          print(f"Error occurred while extracting version month: {e}")
67
68      return rrf_df
69
```

D. Function to apply header to rrf from excel_header

```python
70   def apply_header_to_rrf(file_name, rrf_df):
71       # Check if the corresponding Excel sheet exists
72       if file_name in excel_headers:
73           # Get the headers from the Excel sheet
74           excel_headers_list = excel_headers[file_name]
75           excel_headers_list = [header for header in excel_headers_list if header != 'SVER']
76           # Take names from the excel header list up to the length of the split DataFrame
77           excel_headers_list = excel_headers_list[:len(rrf_df.columns)]
78           # Set the correct header for the DataFrame
79           rrf_df.columns = excel_headers_list
80       return rrf_df
```

E. Function to convert date to indicated format
   a. Here, in the context of RXNSAB, the column SVER derived its value from extracting year part from the VSTART column as done below

```python
121  def process_date_columns(file_name,rrf_df):
122      date_columns = ['VSTART', 'VEND','CREATED_TIMESTAMP', 'UPDATED_TIMESTAMP', 'LAST_RELEASED']
123      for column in date_columns:
124          if column in rrf_df.columns:
125              if file_name == 'RXNSAB':
126                  # Apply the conversion function to each value in the column
127
128
129                  rrf_df[column] = rrf_df[column].apply(convert_date_format)
130                  # Extract year from the VSTART column after date conversion and save it directly as a string
131                  rrf_df['SVER'] = pd.to_datetime(rrf_df['VSTART'], format='%Y-%m-%d').dt.year.astype(str)
132                  # Reorder the columns to place 'SVER' before 'VSTART'
133                  # Reorder columns
134                  # Reorder columns
135                  sver_index = rrf_df.columns.get_loc('SVER')
136                  vstart_index = rrf_df.columns.get_loc('VSTART')
137                  sf_index = rrf_df.columns.get_loc('SF')
138
139                  # Remove 'SVER' from its original position
140
141                  column_sver = rrf_df.pop('SVER')
142
143                  # Insert 'SVER' after 'SF', before 'VSTART'
144                  if sver_index < vstart_index:
145                      rrf_df.insert(vstart_index - 1, 'SVER', column_sver)
146                  elif sver_index > vstart_index:
147                      rrf_df.insert(vstart_index, 'SVER', column_sver)
148              if file_name == 'RXNATOMARCHIVE':
149                  rrf_df[column] = rrf_df[column].apply(update_nato_date)
150      return rrf_df
```

b. Now, the anomalies in the VSTART column of RXNSAB is processed as below:

```python
81
82  def convert_date_format(value):
83
84      try:
85          # Try to parse the value into datetime format
86          parsed_date = pd.to_datetime(value, format='%Y_%m_%d').date()
87          # Extract only the date part
88          return parsed_date.strftime('%Y-%m-%d')
89      except ValueError:
90          if value == '2020':
91              return '2020-01-01'
92          elif value == '5.0_2024_01_04':
93              # Remove the float value and parse the remaining string
94
95              return convert_date_format('2024_01_04')
96          elif value == '2020AA':
97              return '2024-01-02'
98          elif value == '20AA_240205F':
99              return '2024-02-05'
100         else:
101             return value
102
```

c. Now, different date formats of RXNATOMARCHIVE is handled as below:

```python
102
103  def update_nato_date(value):
104      try:
105          # Attempt to parse the value using the first date format
106          parsed_date = pd.to_datetime(value, format='%m/%d/%Y %I:%M:%S %p').date()
107      except ValueError:
108          try:
109              # If the first format fails, attempt to parse using the second date format
110              parsed_date = pd.to_datetime(value, format='%d-%b-%y').date()
111          except ValueError:
112              # If both formats fail, return None or handle the error appropriately
113              return None  # Or handle the error appropriately
114          # Check if the parsed_date is NaT
115      if pd.isnull(parsed_date):
116          return '0000-00-00'  # Replace NaT with '0000-00-00'
117      else:
118          # Extract only the date part and return it in the desired format
119          return parsed_date.strftime('%Y-%m-%d')
120
```

**F.** Now the function to save the final transformed data as .txt delimited by pipe

```
151
152  def save_as_txt_file(rrf_df, file_name, bucket_name):
153      # Construct the filename for the output text file
154      transformation_folder = 'transformation/'
155
156      # Convert DataFrame to CSV format in memory
157      csv_buffer = io.StringIO()
158      rrf_df.to_csv(csv_buffer, sep='|', index=False)
159
160      # Upload the CSV buffer to S3
161      s3_key = transformation_folder + file_name + '.txt'
162      s3.put_object(Bucket=bucket_name, Key=s3_key, Body=csv_buffer.getvalue())
163
164      print(f"Transformed data saved to: s3://{bucket_name}/{s3_key}")
165
166
```

**G.** The main function where the previous functions are called:

      **a.** Here zip files are read from the proper location and the last pipe is ignored by splitting

```
167
168  def read_and_relocate_rrf_files(s3, bucket, key):
169      try:
170          zip_response = s3.get_object(Bucket=bucket, Key=key)
171          zip_data = zip_response['Body'].read()
172
173          # Wrap the zip data in a BytesIO object
174          zip_file = BytesIO(zip_data)
175
176          file_path = 'rrf'
177          unzipped_folder = 'unzipped/'
178
179          with zipfile.ZipFile(zip_file, 'r') as zip_ref:
180              for file_info in zip_ref.infolist():
181                  if file_info.filename.startswith(file_path) and not file_info.filename.endswith('/'):
182                      filename = os.path.basename(file_info.filename)
183                      print(f"The {filename} is read from zip file.")
184
185                      with zip_ref.open(file_info) as source_file:
186                          file_content = source_file.read().decode('utf-8')
187                          if file_content.endswith('|'):
188                              file_content = file_content[:-1]
189                          file_content_io = io.StringIO(file_content)
190                          rrf_df = pd.read_csv(file_content_io, delimiter='|', header=None)
191                          rrf_df = rrf_df.iloc[:, :-1]
192                          print(f"Row count before transformation: {rrf_df.shape[0]}")
193
```

      **b.** Previous functions are called to transform according to the requirements and unzipped files are placed in proper folder

```
192                          print(f"Row count before transformation: {rrf_df.shape[0]}")
193
194                          file_name = os.path.splitext(filename)[0]
195                          apply_header_to_rrf(file_name, rrf_df)
196                          rrf_df = process_date_columns(file_name, rrf_df)
197                          code_set_and_version_month(zip_data, rrf_df)
198
199                          print(f"Row count of {file_name} after transformation: {rrf_df.s
200
201                          pd.set_option('display.max_columns', None)
202                          print(rrf_df.head(5))
203
204                          # Save the transformed DataFrame to a text file
205                          save_as_txt_file(rrf_df, file_name, bucket)
206
207
208                          # Upload the unzipped file to the 'unzipped' folder
209                          unzipped_key = unzipped_folder + filename
210                          s3.put_object(Bucket=bucket, Key=unzipped_key, Body=file_content)
211                          print(f"Unzipped file saved to: s3://{bucket}/{unzipped_key}")
212
213      except Exception as e:
214          print(f"Error occurred: {e}")
```

H. The main lambda function

```
215
216  def lambda_handler(event, context):
217      bucket = event['Records'][0]['s3']['bucket']['name']
218      key = event['Records'][0]['s3']['object']['key']
219
220
221      # This is the function that relocate the rrf files from zip file
222      read_excel_from_s3(bucket)
223      read_and_relocate_rrf_files(s3,bucket,key)
224
```

## Step 3: Results and analysis

A. Zip file is uploaded to zipfiles/ in order to trigger the lambda and excel file is loaded to excelfiles/

B. unzipped/ and transformation/ folders were created after trigger

## unzipped/

Copy S3 URI

**Objects** | Properties

**Objects (9)** Info

Objects are the fundamental entities stored in Amazon S3. You can use Amazon S3 inventory ⎘ to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. Learn more ⎘

🔍 Find objects by prefix

| | Name ▲ | Type ▽ | Last modified ▽ | Size ▽ | Storage class ▽ |
|---|---|---|---|---|---|
| ☐ | 📄 RXNATOMARCHIVE.RRF | RRF | March 10, 2024, 16:04:16 (UTC+05:45) | 71.4 MB | Standard |
| ☐ | 📄 RXNCONSO.RRF | RRF | March 10, 2024, 16:04:26 (UTC+05:45) | 118.6 MB | Standard |
| ☐ | 📄 RXNCUI.RRF | RRF | March 10, 2024, 16:04:28 (UTC+05:45) | 1.7 MB | Standard |
| ☐ | 📄 RXNCUICHANGES.RRF | RRF | March 10, 2024, 16:04:27 (UTC+05:45) | 14.9 KB | Standard |
| ☐ | 📄 RXNDOC.RRF | RRF | March 10, 2024, 16:04:28 (UTC+05:45) | 214.2 KB | Standard |
| ☐ | 📄 RXNREL.RRF | RRF | March 10, 2024, 16:05:21 (UTC+05:45) | 484.4 MB | Standard |
| ☐ | 📄 RXNSAB.RRF | RRF | March 10, 2024, 16:05:27 (UTC+05:45) | 9.8 KB | Standard |
| ☐ | 📄 RXNSAT.RRF | RRF | March 10, 2024, 16:06:09 (UTC+05:45) | 498.7 MB | Standard |
| ☐ | 📄 RXNSTY.RRF | RRF | March 10, 2024, 16:06:16 (UTC+05:45) | 18.4 MB | Standard |

## transformation/

Copy S3 URI

**Objects** | Properties

**Objects (9)** Info

Objects are the fundamental entities stored in Amazon S3. You can use Amazon S3 inventory ⎘ to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. Learn more ⎘

🔍 Find objects by prefix

| | Name ▲ | Type ▽ | Last modified ▽ | Size ▽ | Storage class ▽ |
|---|---|---|---|---|---|
| ☐ | 📄 RXNATOMARCHIVE.txt | txt | March 10, 2024, 16:04:15 (UTC+05:45) | 63.2 MB | Standard |
| ☐ | 📄 RXNCONSO.txt | txt | March 10, 2024, 16:04:25 (UTC+05:45) | 118.7 MB | Standard |
| ☐ | 📄 RXNCUI.txt | txt | March 10, 2024, 16:04:28 (UTC+05:45) | 1.6 MB | Standard |
| ☐ | 📄 RXNCUICHANGES.txt | txt | March 10, 2024, 16:04:27 (UTC+05:45) | 15.0 KB | Standard |
| ☐ | 📄 RXNDOC.txt | txt | March 10, 2024, 16:04:28 (UTC+05:45) | 211.0 KB | Standard |
| ☐ | 📄 RXNREL.txt | txt | March 10, 2024, 16:05:15 (UTC+05:45) | 522.3 MB | Standard |
| ☐ | 📄 RXNSAB.txt | txt | March 10, 2024, 16:05:27 (UTC+05:45) | 10.0 KB | Standard |
| ☐ | 📄 RXNSAT.txt | txt | March 10, 2024, 16:06:03 (UTC+05:45) | 497.9 MB | Standard |
| ☐ | 📄 RXNSTY.txt | txt | March 10, 2024, 16:06:16 (UTC+05:45) | 18.1 MB | Standard |

C. Here is the RXNATOMARCHIVE.rrf files after transformation of dates in the txt format

Raw | **Formatted**

```
RXAUI|AUI|STR|ARCHIVE_TIMESTAMP|CREATED_TIMESTAMP|UPDATED_TIMESTAMP|CODE|IS_BRAND|LAT|LAST_RELEASED|SAUI|VSAB|RXCUI|SAB|TTY|MERGED_TO_RXCUI|Code Set|Version Month
947|A10335796|Mesna|2020-04-27|2005-03-10|2020-04-27|44||ENG|2020-04-06||RXNORM_19AB_200406F|44|RXNORM|IN|44|RxNorm|2024-02-05
1424|A10334758|beta-Alanine|2020-04-27|2005-03-10|2020-04-27|61||ENG|2020-04-06||RXNORM_19AB_200406F|61|RXNORM|IN|61|RxNorm|2024-02-05
1684|A10334529|4-Aminobenzoic Acid|2020-04-27|2005-03-10|2020-04-27|74||ENG|2020-04-06||RXNORM_19AB_200406F|74|RXNORM|IN|74|RxNorm|2024-02-05
2192|A16791816|Eicosapentaenoic Acid|2020-04-27|2005-03-10|2020-04-27|90||ENG|2020-04-06||RXNORM_19AB_200406F|90|RXNORM|PIN|90|RxNorm|2024-02-05
2265|A10334531|5-Hydroxytryptophan|2020-04-27|2005-03-10|2020-11-06|94||ENG|2020-04-06||RXNORM_19AB_200406F|94|RXNORM|IN|94|RxNorm|2024-02-05
2311|A16793037|Ticlopidine Hydrochloride|2020-04-27|2005-03-10|2020-04-27|97||ENG|2020-04-06||RXNORM_19AB_200406F|97|RXNORM|PIN|97|RxNorm|2024-02-05
2332|A10334533|6-Aminocaproic Acid|2020-04-27|2005-03-10|2020-04-27|99||ENG|2020-04-06||RXNORM_19AB_200406F|99|RXNORM|IN|99|RxNorm|2024-02-05
2453|A10334534|6-Mercaptopurine|2010-10-21|2005-03-10|2010-10-21|103||ENG|2010-10-04||RXNORM_10AA_101004F|103|RXNORM|IN|103|RxNorm|2024-02-05
2663|A10336065|Oxyquinoline|2020-04-27|2005-03-10|2020-04-27|110||ENG|2020-04-06||RXNORM_19AB_200406F|110|RXNORM|IN|110|RxNorm|2024-02-05
4330|A10334539|Acebutolol|2020-04-27|2005-03-10|2020-04-27|149||ENG|2020-04-06||RXNORM_19AB_200406F|149|RXNORM|IN|149|RxNorm|2024-02-05
4414|A10334540|Acenocoumarol|2020-04-27|2005-03-10|2020-04-27|154||ENG|2020-04-06||RXNORM_19AB_200406F|154|RXNORM|IN|154|RxNorm|2024-02-05
4458|A10334542|Acepromazine|2020-04-27|2005-03-10|2020-04-27|155||ENG|2020-04-06||RXNORM_19AB_200406F|155|RXNORM|IN|155|RxNorm|2024-02-05
4565|A10334544|Acetanilide|2006-11-15|2005-03-10|2006-11-15|162||ENG|0000-00-00||RXNORM_06AC_061012F|162|RXNORM|IN|162|RxNorm|2024-02-05
4655|A10334545|Acetazolamide|2020-04-27|2005-03-10|2020-04-27|167||ENG|2020-04-06||RXNORM_19AB_200406F|167|RXNORM|IN|167|RxNorm|2024-02-05
4714|A10338597|Acetic Acid|2020-04-27|2005-03-10|2020-04-27|168||ENG|2020-04-06||RXNORM_19AB_200406F|168|RXNORM|IN|168|RxNorm|2024-02-05
4855|A10334546|Acetohexamide|2020-04-27|2005-03-10|2020-04-27|173||ENG|2020-04-06||RXNORM_19AB_200406F|173|RXNORM|IN|173|RxNorm|2024-02-05
4998|A10334547|Acetone|2020-04-27|2005-03-10|2020-04-27|178||ENG|2020-04-06||RXNORM_19AB_200406F|178|RXNORM|IN|178|RxNorm|2024-02-05
```

D. Here is the RXNSAB.rrf files transformed after adding SVER column and changing the date format

Raw | Formatted

```
1  VCUI|RCUI|VSAB|RSAB|SON|SF|SVER|VSTART|VEND|IMETA|RMETA|SLC|SCC|SRL|TFR|CFR|CXTY|TTYL|ATNL|LAT|CENC|CURVER|SABIN|SSN|SCIT|Code set|Code Set|Version Month
2  C5233827|C1140218|MMSL_2024_01_01|MMSL|Multum MediSource Lexicon|MMSL|2024|2024-01-01|||2020AA|;;;Multum Information Services;3200 Cherry Creek South Drive, Suite 300;;Denver;CO;United States;80209;888-633
   -4772 x1420;;;http://www.multum.com/|;;Multum Information Services;3200 Cherry Creek South Drive, Suite 300;;Denver;CO;United States;80209;888-633-4772 x1420;;;http://www.multum.com/|1|||GN,MTH_RXN_BD,CD
   ,IN,MS,SC,BD,BN|DCSA,DDF,DDFA,DHIC,DPC,DRT,DRTA,DST,NDC,TYPE|ENG|UTF-8|Y|Y|Multum|;;;;Medisource Lexicon;;;January 01, 2024;Denver, CO;Multum Information Services, Inc.;;;;;;;;|RxNorm|2024-02-05
3  C5233830|C1140182|MMX_2024_01_02|MMX|Micromedex RED BOOK|MMX|2024|2024-01-02|||2020AA|;;;Micromedex;6200 South Syracuse Way, Suite 300;;Englewood;CO;United States;80111-4740;(800) 525-9083;;info@mdx.com;https
   ://www.ibm.com/us-en/marketplace/micromedex-with-watson/||3|||MTH_RXN_CD,BD,CD,MTH_RXN_BD|DCSA,DDF,DESI_DESC,DRT,DST,LABELER,MMX_RXO,NDC,UPC|ENG|UTF-8|Y|Y|Micromedex|;;;;Micromedex RED BOOK;;;January 02,
   2024;;;;;;;;;;|RxNorm|2024-02-05
4  C5233828|C1140261|NDDF_2024_01_03|NDDF|FDB MedKnowledge (Formerly NDDF Plus)|NDDF|2024|2024-01-03|||2020AA|;;First Databank Customer Support;701 Gateway Blvd, Suite 600;;South San Francisco;CA;United States
   ;94080;800-633-3453;;cs@fdbhealth.com;|;;First Databank Customer Support;701 Gateway Blvd, Suite 600;;South San Francisco;CA;United States;94080;800-633-3453;;cs@fdbhealth.com;|3|||CDC,CDA,CDD,IN
   ,MTH_RXN_CDC,DF|NDC|ENG|UTF-8|Y|Y|FDB MedKnowledge|;;;;FDB MedKnowledge (formerly NDDF Plus);;;January 03, 2024;South San Francisco, CA;First Databank;;;;;;;;|RxNorm|2024-02-05
5  C5233835|C1140284|RXNORM_20AA_240205F|RXNORM|RxNorm Vocabulary|RXNORM|2024|2024-02-05|||2020AA||RxNorm Customer Service;U.S. National Library of Medicine;8600 Rockville Pike;;Bethesda;MD;United States;20894
   ;(888) FIND-NLM;;rxnorminfo@nlm.nih.gov;https://www.nlm.nih.gov/research/umls/rxnorm/|RxNorm Customer Service;U.S. National Library of Medicine;8600 Rockville Pike;;Bethesda;MD;United States;20894;(888)
   FIND-NLM;;rxnorminfo@nlm.nih.gov;https://www.nlm.nih.gov/research/umls/rxnorm/|0|||SCDF,DFG,SBD,PIN,TMSY,BN,SBDC,MIN,SBDG,PSN,SCDC,GPCK,DF,ET,IN,SBDF,SCD,SY,BPCK,SCDG,SCDGP,SCDFP,SBDFP|AMBIGUITY_FLAG,NDC
   ,ORIG_CODE,ORIG_SOURCE,RXN_ACTIVATED,RXN_AI,RXN_AM,RXN_AVAILABLE_STRENGTH,RXN_BN_CARDINALITY,RXN_BOSS_FROM,RXN_BOSS_STRENGTH_DENOM_UNIT,RXN_BOSS_STRENGTH_DENOM_VALUE,RXN_BOSS_STRENGTH_NUM_UNIT
   ,RXN_BOSS_STRENGTH_NUM_VALUE,RXN_HUMAN_DRUG,RXN_IN_EXPRESSED_FLAG,RXN_OBSOLETED,RXN_QUALITATIVE_DISTINCTION,RXN_QUANTITY,RXN_STRENGTH,RXN_VET_DRUG,RXTERM_FORM|ENG|UTF-8|Y|Y|RxNorm work done by the
   National Library of Medicine|;;;;RxNorm;;;META2020AA Full Update 2024_02_05;Bethesda, MD;National Library of Medicine;;;;;;;|RxNorm|2024-02-05
6  C3531723|C2720507|SNOMEDCT_US_2023_06_30|SNOMEDCT_US|US Edition of SNOMED CT|SNOMEDCT|2023|2023-06-30|||2020AA||National Library Of Medicine;NLM is a Charter Member of SNOMED International on behalf of the U
   .S.;National Library of Medicine;8600 Rockville Pike;;Bethesda;MD;United States;20894;1-888-FIND-NLM (1-888-346-3656);;https://support.nlm.nih.gov/support/create-case/;|National Library Of Medicine;NLM is
   a Charter Member of SNOMED International on behalf of the U.S.;National Library Of Medicine;8600 Rockville Pike;;Bethesda;MD;United States;20894;1-888-FIND-NLM (1-888-346-3656);;https://support.nlm.nih
   .gov/support/create-case/;|9|||FN,SY,PT,PTGB,SYGB|ENG|UTF-8|Y|Y|US Edition of SNOMED CT|;;International Health Terminology Standards Development Organisation;;US Edition of SNOMED CT;;;June 30, 2023
   ;Copenhagen, Denmark;International Health Terminology Standards Development Organisation (IHTSDO);;;;;;;;;|RxNorm|2024-02-05
7  C5233837|C1140288|VANDF_2023_12_29|VANDF|Veterans Health Administration National Drug File|VANDF|2023|2023-12-29|||2020AA||Michael Lincoln, M.D.;U.S. Department of Veterans Affairs, Veterans Health
   Administration;;Washington;DC;United States;;;;michael.lincoln@med.va.gov;http://www.pbm.va.gov/default.aspx|Michael Lincoln, M.D.;;;;;;;;;;;michael.lincoln@med.va.gov;http://www.pbm.va.gov/default
   .aspx|0|||CD,MTH_RXN_CD,AB,IN,PT|DCSA,DDF,DRUG_CLASS_TYPE,EXCLUDE_DI_CHECK,NDC,NDF_TRANSMIT_TO_CMOP,NFI,NF_INACTIVATE,NF_NAME,PARENT_CLASS,SNGL_OR_MULT_SRC_PRD,VAC,VA_CLASS_NAME,VA_DISPENSE_UNIT
   ,VA_GENERIC_NAME,VMO|ENG|UTF-8|Y|Y|National Drug File|;;;;Veterans Health Administration National Drug File;;;December 29, 2023;Washington, DC;U.S. Department of Veterans Affairs;;;;;;;|RxNorm|2024-02-05
8  C5233836|C1812643|MTHSPL_2024_01_27|MTHSPL|Metathesaurus FDA Structured Product Labels|MTHSPL|2024|2024-01-27|2006-12-21||2020AA||RxNorm Customer Service;U.S. National Library of Medicine;8600 Rockville Pike
   ;;Bethesda;MD;United States;20894;(888) FIND-NLM;;rxnorminfo@nlm.nih.gov;https://www.nlm.nih.gov/research/umls/rxnorm/|RxNorm Customer Service;U.S. National Library of Medicine;8600 Rockville Pike
   ;;Bethesda;MD;United States;20894;(888) FIND-NLM;;rxnorminfo@nlm.nih.gov;https://www.nlm.nih.gov/research/umls/rxnorm/|0|||DP,SU,MTH_RXN_DP|ANADA,ANDA,BLA,COATING,COLOR,COLORTEXT,CONDITIONAL_NADA,DCSA
   ,DIETARY_SUPPLEMENT,DM,SPL_ID_EXEMPT_DEVICE,IMPRINT_CODE,LABELER_LABEL_TYPE,LEGALLY_MARKETED,UNAPPROVED_NEW_ANIMAL_DRUGS_FOR_MINOR_SPECIES,MARKETING_CATEGORY,MARKETING_EFFECTIVE_TIME_HIGH
```

E. Finally validating the row counts as before and after:

| | | |
|---|---|---|
| ▶ | 2024-03-10T17:30:09.626+05:45 | Row count before transformation: 1133065 |
| ▶ | 2024-03-10T17:30:09.627+05:45 | Version month: 2024-02-05 |
| ▶ | 2024-03-10T17:30:09.640+05:45 | Row count of RXNCONSO after transformation: 1133065 |
| ▶ | 2024-03-10T17:30:22.414+05:45 | Row count before transformation: 153 |
| ▶ | 2024-03-10T17:30:22.414+05:45 | Version month: 2024-02-05 |
| ▶ | 2024-03-10T17:30:22.415+05:45 | Row count of RXNCUICHANGES after transformation: 153 |
| ▶ | 2024-03-10T17:30:22.559+05:45 | Row count before transformation: 30046 |
| ▶ | 2024-03-10T17:30:22.559+05:45 | Version month: 2024-02-05 |
| ▶ | 2024-03-10T17:30:22.560+05:45 | Row count of RXNCUI after transformation: 30046 |
| ▶ | 2024-03-10T17:30:22.946+05:45 | The RXNDOC.RRF is read from zip file. |
| ▶ | 2024-03-10T17:30:22.953+05:45 | Row count before transformation: 3445 |
| ▶ | 2024-03-10T17:30:22.954+05:45 | Version month: 2024-02-05 |
| ▶ | 2024-03-10T17:30:22.954+05:45 | Row count of RXNDOC after transformation: 3445 |
| ▶ | 2024-03-10T17:30:39.524+05:45 | Row count before transformation: 7154306 |
| ▶ | 2024-03-10T17:30:39.525+05:45 | Version month: 2024-02-05 |
| ▶ | 2024-03-10T17:30:39.624+05:45 | Row count of RXNREL after transformation: 7154306 |
| ▶ | 2024-03-10T17:31:53.314+05:45 | The RXNSAB.RRF is read from zip file. |
| ▶ | 2024-03-10T17:31:53.567+05:45 | Row count before transformation: 13 |
| ▶ | 2024-03-10T17:31:53.573+05:45 | Version month: 2024-02-05 |
| ▶ | 2024-03-10T17:31:53.573+05:45 | Row count of RXNSAB after transformation: 13 |

| | | |
|---|---|---|
| ▶ | 2024-03-10T17:32:09.994+05:45 | Row count before transformation: 7222404 |
| ▶ | 2024-03-10T17:32:09.995+05:45 | Version month: 2024-02-05 |
| ▶ | 2024-03-10T17:32:10.096+05:45 | Row count of RXNSAT after transformation: 7222404 |
| ▶ | 2024-03-10T17:33:08.069+05:45 | The RXNSTY.RRF is read from zip file. |
| ▶ | 2024-03-10T17:33:08.811+05:45 | Row count before transformation: 461874 |
| ▶ | 2024-03-10T17:33:08.812+05:45 | Version month: 2024-02-05 |
| ▶ | 2024-03-10T17:33:08.818+05:45 | Row count of RXNSTY after transformation: 461874 |