# SARB Features

January 24, 2021

```python
[1]: import pandas as pd
     import numpy as np
     import matplotlib.pyplot as plt

     feature_set_sarb = pd.read_csv('data/sarb_features_data.csv').
      ↪drop(['unemployment rate'], axis=1).set_index('Date')
     target = pd.read_csv('data/sarb_target_data.csv').set_index('Date')
```

```python
[2]: import warnings

     warnings.filterwarnings('ignore')
```

## 0.1 The full feature set

*These feature were accessed from the South African Reserve Bank.*

*There are **147 features in total**, these cover a significant portfion of the South African economy*

**The data from 1922-01-01 to 2020-01-01** if it used for unemployment forecasting, deleting redudant observations is helpful

```python
[3]: feature_set_sarb
```

```
[3]:            Final consumption expenditure by general government   \
     Date
     1922-01-01                                              NaN
     1922-02-01                                              NaN
     1922-03-01                                              NaN
     1922-04-01                                              NaN
     1922-05-01                                              NaN
     …                                                       …
     2019-09-01                                              NaN
     2019-10-01                                              NaN
     2019-11-01                                              NaN
     2019-12-01                                              NaN
     2020-01-01                                              NaN

                Consolidated general government: Revenue   \
     Date
```

```
1922-01-01                                      NaN
1922-02-01                                      NaN
1922-03-01                                      NaN
1922-04-01                                      NaN
1922-05-01                                      NaN
...                                             ...
2019-09-01                                      NaN
2019-10-01                                      NaN
2019-11-01                                      NaN
2019-12-01                                      NaN
2020-01-01                                      NaN

            Foreign liabilities: Total portfolio investment  \
Date
1922-01-01                                                NaN
1922-02-01                                                NaN
1922-03-01                                                NaN
1922-04-01                                                NaN
1922-05-01                                                NaN
...                                                       ...
2019-09-01                                                NaN
2019-10-01                                                NaN
2019-11-01                                                NaN
2019-12-01                                                NaN
2020-01-01                                                NaN

            Foreign liabilities: Portfolio investment: Equity securities  \
Date
1922-01-01                                                NaN
1922-02-01                                                NaN
1922-03-01                                                NaN
1922-04-01                                                NaN
1922-05-01                                                NaN
...                                                       ...
2019-09-01                                                NaN
2019-10-01                                                NaN
2019-11-01                                                NaN
2019-12-01                                                NaN
2020-01-01                                                NaN

            Domestic output: All groups   \
Date
1922-01-01                             NaN
1922-02-01                             NaN
1922-03-01                             NaN
1922-04-01                             NaN
1922-05-01                             NaN
```

```
...                              ...
2019-09-01                       114.3
2019-10-01                       114.6
2019-11-01                       114.3
2019-12-01                       114.5
2020-01-01                        NaN


            Final consumption expenditure by households: Total    \
Date
1922-01-01                                               NaN
1922-02-01                                               NaN
1922-03-01                                               NaN
1922-04-01                                               NaN
1922-05-01                                               NaN
...                                                       ...
2019-09-01                                               NaN
2019-10-01                                               NaN
2019-11-01                                               NaN
2019-12-01                                               NaN
2020-01-01                                               NaN


            Gross fixed capital formation   \
Date
1922-01-01                            NaN
1922-02-01                            NaN
1922-03-01                            NaN
1922-04-01                            NaN
1922-05-01                            NaN
...                                   ...
2019-09-01                            NaN
2019-10-01                            NaN
2019-11-01                            NaN
2019-12-01                            NaN
2020-01-01                            NaN


            SDDS - Financial derivative liabilities  \
Date
1922-01-01                                      NaN
1922-02-01                                      NaN
1922-03-01                                      NaN
1922-04-01                                      NaN
1922-05-01                                      NaN
...                                             ...
2019-09-01                                      NaN
2019-10-01                                      NaN
2019-11-01                                      NaN
2019-12-01                                      NaN
```

```
2020-01-01                                                    NaN


            Foreign liabilities: Portfolio investment: Debt securities  \
Date
1922-01-01                                                    NaN
1922-02-01                                                    NaN
1922-03-01                                                    NaN
1922-04-01                                                    NaN
1922-05-01                                                    NaN
…                                                             …
2019-09-01                                                    NaN
2019-10-01                                                    NaN
2019-11-01                                                    NaN
2019-12-01                                                    NaN
2020-01-01                                                    NaN


            Change in inventories   …  \
Date                                …
1922-01-01                    NaN   …
1922-02-01                    NaN   …
1922-03-01                    NaN   …
1922-04-01                    NaN   …
1922-05-01                    NaN   …
…                             …    …
2019-09-01                    NaN   …
2019-10-01                    NaN   …
2019-11-01                    NaN   …
2019-12-01                    NaN   …
2020-01-01                    NaN   …


            Physical volume of manufacturing production: Total   \
Date
1922-01-01                                                  NaN
1922-02-01                                                  NaN
1922-03-01                                                  NaN
1922-04-01                                                  NaN
1922-05-01                                                  NaN
…                                                           …
2019-09-01                                                 99.0
2019-10-01                                                101.5
2019-11-01                                                100.0
2019-12-01                                                  NaN
2020-01-01                                                  NaN


            Remuneration per worker in non-agricultural: Total   \
Date
1922-01-01                                                  NaN
```

```
1922-02-01                                        NaN
1922-03-01                                        NaN
1922-04-01                                        NaN
1922-05-01                                        NaN
…                                                 …
2019-09-01                                        NaN
2019-10-01                                        NaN
2019-11-01                                        NaN
2019-12-01                                        NaN
2020-01-01                                        NaN

            Consolidated general government: Non-financial assets - Net  \
Date
1922-01-01                                        NaN
1922-02-01                                        NaN
1922-03-01                                        NaN
1922-04-01                                        NaN
1922-05-01                                        NaN
…                                                 …
2019-09-01                                        NaN
2019-10-01                                        NaN
2019-11-01                                        NaN
2019-12-01                                        NaN
2020-01-01                                        NaN

            Consolidated general government: Cash surplus / deficit  \
Date
1922-01-01                                        NaN
1922-02-01                                        NaN
1922-03-01                                        NaN
1922-04-01                                        NaN
1922-05-01                                        NaN
…                                                 …
2019-09-01                                        NaN
2019-10-01                                        NaN
2019-11-01                                        NaN
2019-12-01                                        NaN
2020-01-01                                        NaN

            CPI Headline   Gross domestic expenditure   \
Date
1922-01-01          0.6                          NaN
1922-02-01          0.6                          NaN
1922-03-01          0.6                          NaN
1922-04-01          0.6                          NaN
1922-05-01          0.6                          NaN
…                    …                            …
```

```
2019-09-01          113.4                              NaN
2019-10-01          113.4                              NaN
2019-11-01          113.5                              NaN
2019-12-01          113.8                              NaN
2020-01-01            NaN                              NaN


            Net cash-flow from operating activities   \
Date
1922-01-01                                      NaN
1922-02-01                                      NaN
1922-03-01                                      NaN
1922-04-01                                      NaN
1922-05-01                                      NaN
…                                                …
2019-09-01                                      NaN
2019-10-01                                      NaN
2019-11-01                                      NaN
2019-12-01                                      NaN
2020-01-01                                      NaN


            Non-agricultural employment: Total   \
Date
1922-01-01                                  NaN
1922-02-01                                  NaN
1922-03-01                                  NaN
1922-04-01                                  NaN
1922-05-01                                  NaN
…                                            …
2019-09-01                                  NaN
2019-10-01                                  NaN
2019-11-01                                  NaN
2019-12-01                                  NaN
2020-01-01                                  NaN


            Consolidated general government: Expense  Residual item
Date
1922-01-01                                      NaN            NaN
1922-02-01                                      NaN            NaN
1922-03-01                                      NaN            NaN
1922-04-01                                      NaN            NaN
1922-05-01                                      NaN            NaN
…                                                …              …
2019-09-01                                      NaN            NaN
2019-10-01                                      NaN            NaN
2019-11-01                                      NaN            NaN
2019-12-01                                      NaN            NaN
2020-01-01                                      NaN            NaN
```

```
[1432 rows x 147 columns]
```

## 0.2 Additional Feature or Target

Depending on the purpose the **unemployment rate** be considered as a *target*, if you want to forecast it. Or is another feature if you want to forecast something else

```
[4]: target
```

```
[4]:            unemployment rate
     Date
     1922-01-01              NaN
     1922-02-01              NaN
     1922-03-01              NaN
     1922-04-01              NaN
     1922-05-01              NaN
     ...                     ...
     2019-09-01             29.1
     2019-10-01              NaN
     2019-11-01              NaN
     2019-12-01             29.1
     2020-01-01              NaN

     [1432 rows x 1 columns]
```

# 1 Important Note: Missing Data from mixing frequencies

### 1.0.1 Ensure to review SARB Feature List.pdf to see the frequency of each feature that was accessed fro the SARB (https://github.com/rudzanimulaudzi/sarb_feature_set/blob/main/SARB%20Feature%20List

*The missing data occurs because we are merging data that monthly and data that is quartely, hence all quartely data should be expected to have missing data.* This is normal when dealing with multiple time series.
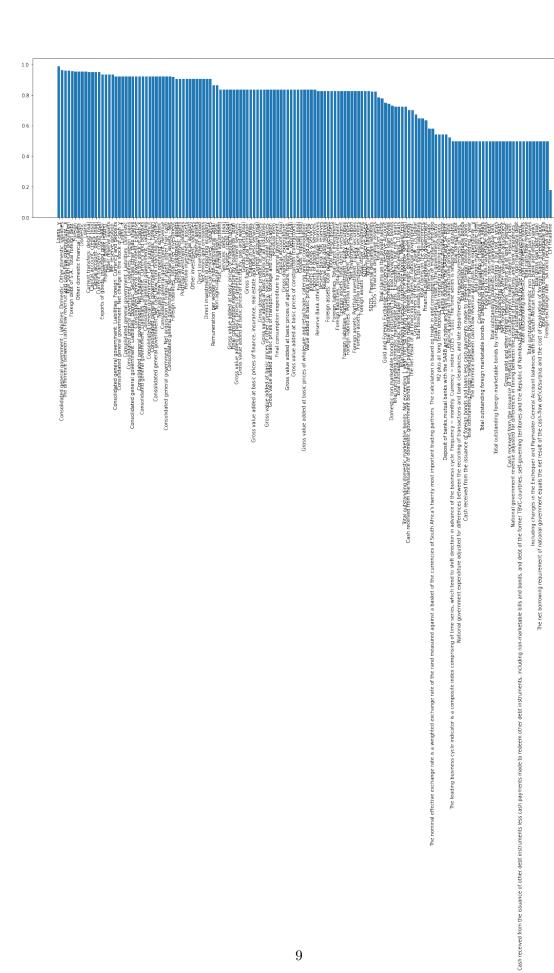
**In this case the data is from 1922. Some data observations should be deleted according to match the target variable. This is will improve the percentage of missing data.**

```
[7]: #Here we visualize the frequency of missingness of each feature
     feature = feature_set_sarb.isna().sum()/len(feature_set_sarb)

     feature = feature.sort_values(ascending=False)
     feature_df = pd.DataFrame(feature.index, columns=['Feature Name'])
     feature_df['Missing Frequency'] = np.array(feature.values)
     feature_df['Rank'] = feature_df['Missing Frequency'].rank(ascending=False)
     feature_df
```

```
[7]:                                    Feature Name  Missing Frequency  \
     0                                       Loans _y           0.986732
     1    Consolidated general government: Liabilities: …           0.959497
     2    The difference between cash-flow revenue and c…           0.958799
     3                     Total South African population           0.958101
     4            Foreign debt of S.A.: Total foreign debt           0.953212
     ..                                            …                 …
     142                     Total gross loan debt (nsa)           0.497207
     143  The net borrowing requirement of national gove…           0.497207
     144          Total outstanding domestic marketable bills           0.497207
     145     Foreign exchange rate : SA rand per USA dollar           0.496508
     146                                     CPI Headline           0.178771

          Rank
     0      1.0
     1      2.0
     2      3.0
     3      4.0
     4      5.0
     ..      …
     142  132.0
     143  132.0
     144  132.0
     145  146.0
     146  147.0

     [147 rows x 3 columns]
```

```python
[8]: plt.figure(figsize=(20,6))
     plt.xticks(rotation=90)
     plt.bar(feature_df['Feature Name'], feature_df['Missing Frequency'])
```

```
[8]: <BarContainer object of 147 artists>
```

### 1.0.2 Here I give an example of how to deal with missing data using Forward Fill approach. See https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.fillna.html

```
[14]: #Data imputation strategy is foward fill but there many other options,
      #choose based on your research needs and what gives higher accuracy
      x_values_ffill = feature_set_sarb.fillna(method='ffill')
      y_values_ffill = target.fillna(method='ffill')
```

```
[15]: # Remove all data points before unemployment rate data is available.␣
      ↪Unemployment rate is my target variable.
      valid_start = y_values_ffill.first_valid_index()
      y_values_ffill = y_values_ffill[valid_start : ]
      x_values_ffill = x_values_ffill[valid_start : ]
```

```
[16]: #I fill with NA here to avoid any features that might be NA i.e. insurance
      x_values_ffill = x_values_ffill.fillna(0)
      x_values_ffill.isna().sum()
```

```
[16]: Final consumption expenditure by general government           0
      Consolidated general government: Revenue                     0
      Foreign liabilities: Total portfolio investment              0
      Foreign liabilities: Portfolio investment: Equity securities 0
      Domestic output: All groups                                  0
                                                                   ..
      Gross domestic expenditure                                   0
      Net cash-flow from operating activities                      0
      Non-agricultural employment: Total                           0
      Consolidated general government: Expense                     0
      Residual item                                                0
      Length: 147, dtype: int64
```

```
[17]: x_values_ffill
```

```
[17]:            Final consumption expenditure by general government   \
      Date
      1970-03-01                                          142014.0
      1970-04-01                                          142014.0
      1970-05-01                                          142014.0
      1970-06-01                                          142014.0
      1970-07-01                                          142014.0
      …                                                        …
      2019-09-01                                          653236.0
      2019-10-01                                          653236.0
```

```
2019-11-01                                                653236.0
2019-12-01                                                653236.0
2020-01-01                                                653236.0

            Consolidated general government: Revenue  \
Date
1970-03-01                                      0.0
1970-04-01                                      0.0
1970-05-01                                      0.0
1970-06-01                                      0.0
1970-07-01                                      0.0
…                                               …
2019-09-01                                 462964.0
2019-10-01                                 462964.0
2019-11-01                                 462964.0
2019-12-01                                 462964.0
2020-01-01                                 462964.0

            Foreign liabilities: Total portfolio investment   \
Date
1970-03-01                                            2.0
1970-04-01                                            2.0
1970-05-01                                            2.0
1970-06-01                                            2.0
1970-07-01                                            2.0
…                                                     …
2019-09-01                                         3313.0
2019-10-01                                         3313.0
2019-11-01                                         3313.0
2019-12-01                                         3313.0
2020-01-01                                         3313.0

            Foreign liabilities: Portfolio investment: Equity securities  \
Date
1970-03-01                                                  2.0
1970-04-01                                                  2.0
1970-05-01                                                  2.0
1970-06-01                                                  2.0
1970-07-01                                                  2.0
…                                                           …
2019-09-01                                               2036.0
2019-10-01                                               2036.0
2019-11-01                                               2036.0
2019-12-01                                               2036.0
2020-01-01                                               2036.0

            Domestic output: All groups   \
```

```
Date
1970-03-01                          4.1
1970-04-01                          4.2
1970-05-01                          4.2
1970-06-01                          4.2
1970-07-01                          4.2
…                                    …
2019-09-01                        114.3
2019-10-01                        114.6
2019-11-01                        114.3
2019-12-01                        114.5
2020-01-01                        114.5

            Final consumption expenditure by households: Total    \
Date
1970-03-01                                          459049.0
1970-04-01                                          459049.0
1970-05-01                                          459049.0
1970-06-01                                          459049.0
1970-07-01                                          459049.0
…                                                        …
2019-09-01                                         1961051.0
2019-10-01                                         1961051.0
2019-11-01                                         1961051.0
2019-12-01                                         1961051.0
2020-01-01                                         1961051.0

            Gross fixed capital formation    \
Date
1970-03-01                      176103.0
1970-04-01                      176103.0
1970-05-01                      176103.0
1970-06-01                      176103.0
1970-07-01                      176103.0
…                                    …
2019-09-01                      613640.0
2019-10-01                      613640.0
2019-11-01                      613640.0
2019-12-01                      613640.0
2020-01-01                      613640.0

            SDDS - Financial derivative liabilities  \
Date
1970-03-01                                     0.0
1970-04-01                                     0.0
1970-05-01                                     0.0
1970-06-01                                     0.0
```

```
1970-07-01                                           0.0
…                                                    …
2019-09-01                                         118.0
2019-10-01                                         118.0
2019-11-01                                         118.0
2019-12-01                                         118.0
2020-01-01                                         118.0


            Foreign liabilities: Portfolio investment: Debt securities  \
Date
1970-03-01                                                        0.0
1970-04-01                                                        0.0
1970-05-01                                                        0.0
1970-06-01                                                        0.0
1970-07-01                                                        0.0
…                                                                 …
2019-09-01                                                     1277.0
2019-10-01                                                     1277.0
2019-11-01                                                     1277.0
2019-12-01                                                     1277.0
2020-01-01                                                     1277.0


            Change in inventories    …  \
Date                                 …
1970-03-01             18617.0  …
1970-04-01             18617.0  …
1970-05-01             18617.0  …
1970-06-01             18617.0  …
1970-07-01             18617.0  …
…                          …   …
2019-09-01             -9526.0  …
2019-10-01             -9526.0  …
2019-11-01             -9526.0  …
2019-12-01             -9526.0  …
2020-01-01             -9526.0  …


            Physical volume of manufacturing production: Total   \
Date
1970-03-01                                                38.3
1970-04-01                                                41.5
1970-05-01                                                39.8
1970-06-01                                                40.6
1970-07-01                                                41.7
…                                                        …
2019-09-01                                                99.0
2019-10-01                                               101.5
2019-11-01                                               100.0
```

```
2019-12-01                                                     100.0
2020-01-01                                                     100.0


            Remuneration per worker in non-agricultural: Total   \
Date
1970-03-01                                                       0.9
1970-04-01                                                       0.9
1970-05-01                                                       0.9
1970-06-01                                                       0.9
1970-07-01                                                       0.9
…                                                                …
2019-09-01                                                     173.4
2019-10-01                                                     173.4
2019-11-01                                                     173.4
2019-12-01                                                     173.4
2020-01-01                                                     173.4


            Consolidated general government: Non-financial assets - Net   \
Date
1970-03-01                                                       0.0
1970-04-01                                                       0.0
1970-05-01                                                       0.0
1970-06-01                                                       0.0
1970-07-01                                                       0.0
…                                                                …
2019-09-01                                                  -26886.0
2019-10-01                                                  -26886.0
2019-11-01                                                  -26886.0
2019-12-01                                                  -26886.0
2020-01-01                                                  -26886.0


            Consolidated general government: Cash surplus / deficit   \
Date
1970-03-01                                                       0.0
1970-04-01                                                       0.0
1970-05-01                                                       0.0
1970-06-01                                                       0.0
1970-07-01                                                       0.0
…                                                                …
2019-09-01                                                  -82087.0
2019-10-01                                                  -82087.0
2019-11-01                                                  -82087.0
2019-12-01                                                  -82087.0
2020-01-01                                                  -82087.0


            CPI Headline   Gross domestic expenditure   \
Date
```

```
1970-03-01          1.6                965734.0
1970-04-01          1.6                965734.0
1970-05-01          1.6                965734.0
1970-06-01          1.6                965734.0
1970-07-01          1.6                965734.0
…                   …                  …
2019-09-01         113.4              3222339.0
2019-10-01         113.4              3222339.0
2019-11-01         113.5              3222339.0
2019-12-01         113.8              3222339.0
2020-01-01         113.8              3222339.0


           Net cash-flow from operating activities   \
Date
1970-03-01                                     0.0
1970-04-01                                     0.0
1970-05-01                                     0.0
1970-06-01                                     0.0
1970-07-01                                     0.0
…                                              …
2019-09-01                                 -55201.0
2019-10-01                                 -55201.0
2019-11-01                                 -55201.0
2019-12-01                                 -55201.0
2020-01-01                                 -55201.0


           Non-agricultural employment: Total   \
Date
1970-03-01                              49.2
1970-04-01                              49.2
1970-05-01                              49.2
1970-06-01                              49.2
1970-07-01                              49.2
…                                       …
2019-09-01                             106.8
2019-10-01                             106.8
2019-11-01                             106.8
2019-12-01                             106.8
2020-01-01                             106.8


           Consolidated general government: Expense  Residual item
Date
1970-03-01                                     0.0        169951.0
1970-04-01                                     0.0        169951.0
1970-05-01                                     0.0        169951.0
1970-06-01                                     0.0        169951.0
1970-07-01                                     0.0        169951.0
```

```
...                                                       ...              ...
2019-09-01                                           518165.0           3938.0
2019-10-01                                           518165.0           3938.0
2019-11-01                                           518165.0           3938.0
2019-12-01                                           518165.0           3938.0
2020-01-01                                           518165.0           3938.0

[794 rows x 147 columns]
```

# 2 Options for dealing with mixed frequencies and high dimensional data from SARB

1. Only use data with the same data frequency i.e. delete monthy data and only use quartely because unemployment is quartely OR delete all quarterly and use monthly but using sampling techniques / imputation to fill missing data
2. Use sampling techniques, upsample all monthly data to be quartely so that all your data is now in the same frequency
3. Use data imputation techniques, I gave an example above. You must use the same one for your features and target
4. Remove low variance and duplicate feature using statistical scores

These four should you leave you with the right data to model either the South African unemployment rate (which I suggest you use) or any other macroeconomic variable in the data.

**Also remember your features and target are in different dataframes. Use pd.concat to merge these two i.e pd.concat([x_values_ffill, y_values_ffill], axis=1)**

[ ]: