

Continuando a coleta de dados, inserimos o arquivo <https://dados.gov.br/dataset/cadastro-nacional-de-reclamacoes-fundamentadas-procons-sindec1> (<https://dados.gov.br/dataset/cadastro-nacional-de-reclamacoes-fundamentadas-procons-sindec1>), acessado em 25/07/2021. Os arquivos do endereço [https://dadosabertos.pgfn.gov.br/2021\\_trimestre\\_02/Dados\\_abertos\\_Nao\\_Previdenciario.zip](https://dadosabertos.pgfn.gov.br/2021_trimestre_02/Dados_abertos_Nao_Previdenciario.zip) ([https://dadosabertos.pgfn.gov.br/2021\\_trimestre\\_02/Dados\\_abertos\\_Nao\\_Previdenciario.zip](https://dadosabertos.pgfn.gov.br/2021_trimestre_02/Dados_abertos_Nao_Previdenciario.zip)), acessado em 25/07/2021, foram concatenados no notebook 01, resultando no arquivo "concatenadoordenado.csv", de 4GB. Arquivo igualmente "upado" para este notebook.

In [1]:

```
1 import pandas as pd
2 from collections import Counter
```

## Inserindo o dataset do Sindec e o preparando para o merge com o "concatenadoordenado" oriundo do notebook 01

In [2]:

```
1 #Carregando o dataset do Sindec (Procon)
2 df_procon = pd.read_csv(r'C:\Users\73594253368\Desktop\Curso\Datasets\Procon\CRF2019Dac
```

```
b'Skipping line 117: expected 23 fields, saw 25\nSkipping line 1205: expected 23 fields, saw 25\nSkipping line 1285: expected 23 fields, saw 24\nSkipping line 3596: expected 23 fields, saw 25\nSkipping line 4134: expected 23 fields, saw 24\nSkipping line 5483: expected 23 fields, saw 25\nSkipping line 8028: expected 23 fields, saw 24\nSkipping line 8307: expected 23 fields, saw 24\nSkipping line 8824: expected 23 fields, saw 25\nSkipping line 9516: expected 23 fields, saw 24\nSkipping line 9696: expected 23 fields, saw 24\nSkipping line 9698: expected 23 fields, saw 24\nSkipping line 10545: expected 23 fields, saw 25\nSkipping line 11232: expected 23 fields, saw 25\nSkipping line 11503: expected 23 fields, saw 25\nSkipping line 11548: expected 23 fields, saw 24\nSkipping line 11554: expected 23 fields, saw 25\nSkipping line 11703: expected 23 fields, saw 24\nSkipping line 12473: expected 23 fields, saw 24\nSkipping line 13879: expected 23 fields, saw 24\nSkipping line 16696: expected 23 fields, saw 24\nSkipping line 17302: expected 23 fields, saw 24\n'
```

In [3]:

```
1 df_procon.shape
```

Out[3]:

```
(17555, 23)
```

In [4]:

```
1 df_procon.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 17555 entries, 0 to 17554
Data columns (total 23 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   AnoCalendario                        17555 non-null  int64
1   DataArquivamento                  17552 non-null  object
2   DataAbertura                       17552 non-null  object
3   CodigoRegiao                       17555 non-null  int64
4   Regiao                             17555 non-null  object
5   UF                                  17555 non-null  object
6   strRazaoSocial                     17555 non-null  object
7   strNomeFantasia                    14034 non-null  object
8   Tipo                               17555 non-null  int64
9   NumeroCNPJ                         16028 non-null  float64
10  RadicalCNPJ                       15974 non-null  float64
11  RazaoSocialRFB                     14572 non-null  object
12  NomeFantasiaRFB                    7651 non-null   object
13  CNAEPrincipal                      14572 non-null  float64
14  DescCNAEPrincipal                  14502 non-null  object
15  Atendida                           17555 non-null  object
16  CodigoAssunto                      17541 non-null  float64
17  DescricaoAssunto                   17541 non-null  object
18  CodigoProblema                     45 non-null     float64
19  DescricaoProblema                  45 non-null     object
20  SexoConsumidor                     17546 non-null  object
21  FaixaEtariaConsumidor              17555 non-null  object
22  CEPConsumidor                      13921 non-null  float64
dtypes: float64(6), int64(3), object(14)
memory usage: 3.1+ MB
```

In [5]:

```
1 df_procon.head()
```

Out[5]:

	AnoCalendario	DataArquivamento	DataAbertura	CodigoRegiao	Regiao	UF	strf
0	2019	2019-10-04 11:12:54.000	2019-09-02 09:27:15.000	1	Norte	RO	
1	2019	2019-01-08 10:56:05.000	2018-12-04 15:19:18.000	1	Norte	RO	BANCO DC
2	2019	2019-08-15 15:14:14.000	2019-07-16 17:00:46.000	1	Norte	RO	CENTRAIS DE ROI
3	2019	2019-01-04 11:31:47.000	2018-04-19 10:09:02.000	1	Norte	RO	NOVA PRO PROFISSION
4	2019	2019-01-04 10:26:36.000	2018-08-30 09:46:37.000	1	Norte	RO	OI MO' EI

5 rows × 23 columns

In [6]:

```
1 #Preparando o campo CNPJ para o merge: retirar nullos, sinais e transformar para inteiro
2 df_procon = df_procon.dropna(subset=['NumeroCNPJ']).dropna(axis=1, how = 'all')
```

In [8]:

```
1 df_procon.shape
```

Out[8]:

(16028, 23)

In [ ]:

1

In [ ]:

```
1 df_procon.info()
```

In [ ]:

```
1 df_procon['NumeroCNPJ'] = df_procon['NumeroCNPJ'].replace('.', '')
2 df_procon['NumeroCNPJ'] = df_procon['NumeroCNPJ'].replace('/', '')
3 df_procon['NumeroCNPJ'] = df_procon['NumeroCNPJ'].replace('-', '')
```

In [ ]:

```
1 df_procon['NumeroCNPJ']
```

In [ ]:

```
1 df_procon['NumeroCNPJ'] = pd.to_numeric(df_procon['NumeroCNPJ'],downcast='integer')
```

In [ ]:

```
1 df_procon
```

In [ ]:

```
1 # Colocando index para permitir visualização, na hora do merge com o dataset da PFN, de  
2 # empresas em DAU  
3 # df_procon_2 é um ProconIndexado  
4 df_procon_2 = df_procon
```

In [ ]:

```
1 df_procon_2 = df_procon_2.rename_axis('index1').reset_index()
```

In [ ]:

```
1 df_procon_2
```

In [ ]:

```
1 df_procon_2.shape
```

In [ ]:

```
1 df_procon_2.head()
```

**As 27 tabelas de devedores da PFN foram concatenadas no dataset concatenadoordenado.csv. A seguir o tratamento para o merge com a base do Sindec**

In [ ]:

```
1 #Carregando o dataset da PFN  
2 df_pfn = pd.read_csv(r'C:\Users\73594253368\Desktop\Curso\Datasets\Procon\concatenadoordenado.csv')
```

In [ ]:

```
1 df_pfn.info()
```

In [ ]:

```
1 df_pfn.rename(columns = {'TIPO_PESSOA': 'Tipo_Pessoa'}, inplace=True)
```

Restringimos o data frame a apenas duas colunas, pois nosso objetivo com os dados da PFN é, a partir dos dados das demandas no Sindec, sabermos se a empresa constante nas demandas é listada, também, como devedora da Fazenda Nacional.

In [ ]:

```
1 df_pfn_colunas = df_pfn[['CPF_CNPJ', 'Tipo_Pessoa']]
```

In [ ]:

```
1 df_pfn_colunas.head()
```

In [ ]:

```
1 #Removendo os nulos
2 df_pfn_colunas = df_pfn_colunas.dropna(subset=['CPF_CNPJ']).dropna(axis=1, how = 'all')
3 df_pfn_colunas = df_pfn_colunas.dropna(subset=['Tipo_Pessoa']).dropna(axis=1, how = 'all')
```

In [ ]:

```
1 #Retirando as pessoas físicas
2 df_pfn_colunas = df_pfn_colunas[df_pfn_colunas.Tipo_Pessoa.str.contains('Pessoa jurídica')]
3 df_pfn_colunas.rename(columns = {'CPF_CNPJ': 'CNPJ'}, inplace=True)
```

In [ ]:

```
1 #Limpando os CNPJs
2 df_pfn_colunas['CNPJ'] = df_pfn_colunas['CNPJ'].str.replace('.', '')
3 df_pfn_colunas['CNPJ'] = df_pfn_colunas['CNPJ'].str.replace('/', '')
4 df_pfn_colunas['CNPJ'] = df_pfn_colunas['CNPJ'].str.replace('-', '')
5 df_pfn_colunas['CNPJ'] = pd.to_numeric(df_pfn_colunas['CNPJ'], downcast='integer')
```

In [ ]:

```
1 df_pfn_colunas.info()
```

In [ ]:

```
1 df_pfn_colunas['Tipo_Pessoa'].value_counts()
```

In [ ]:

```
1 df_pfn_colunas.info()
```

In [ ]:

```
1 df_pfn_colunas.shape
```

In [ ]:

```
1
```

In [ ]:

```
1
```

## Merge Procon Indexado com PFN duas colunas

In [ ]:

```
1
2 df_merged = pd.merge(df_procon_2, df_pfn_colunas, how='left',
3                       left_on=['NumeroCNPJ'],
4                       right_on=['CNPJ'])
```

In [ ]:

```
1 df_merged.shape
```

In [ ]:

```
1
```

In [ ]:

```
1 df_merged.info()
```

In [ ]:

```
1
```

In [ ]:

```
1 #Retirando os registros duplicados
2 df_merged = df_merged.drop_duplicates()
```

In [ ]:

```
1
```

In [ ]:

```
1 df_merged.info()
```

In [ ]:

```
1
```

In [ ]:

```
1 # Exportando o dataframe criado após os merges pra a etapa seguinte, a de tratamento de
2 # df_merged.to_csv(r'C:\Users\73594253368\Desktop\Curso\Datasets\Procon\df_merged.csv')
```

In [ ]:

```
1
```