



PUC Minas

PONTIFÍCIA UNIVERSIDADE CATÓLICA DE MINAS GERAIS

NÚCLEO DE EDUCAÇÃO A DISTÂNCIA

Pós-graduação *Lato Sensu* em Analytics e Business Intelligence

RELATÓRIO TÉCNICO

ORGANIZAÇÃO DOS DADOS DA PREVIDÊNCIA SOCIAL

Bruna Fernandes de Medeiros

Belo Horizonte

2021

SUMÁRIO

- 1. Introdução3**
 - 1.1. Contexto3
 - 1.2. Objetivos3
 - 1.3. Público alvo3
- 2. Modelo de Dados4**
 - 2.1. Modelo Dimensional4
 - 2.2. Fatos e Dimensões4
- 3. Integração, Tratamento e Carga de Dados5**
- 4. Camada de Apresentação15**
 - 4.1 Dashboard15
 - 4.2 Análises avançadas15
- 5. Registros de Homologação15**
- 5. Conclusões15**
- 6. Links16**

1. Introdução

1.1. Contexto

O Brasil é uma país que distribui muitos tipos de benefícios da previdência, são elas: por aposentadoria acidentária, por idade, invalidez, entre outras. Porém não é fornecida ferramenta que facilite o acesso a essas informações por categorias. E com isso não é possível diferenciar os valores destinados a cada tipo.

Entretanto as informações estão disponibilizadas em formato CSV de forma completa no site do governo, onde consta ano, sexo, quantidade de benefícios concedidos.

Sendo assim possível que as empresas de forma geral possam acessar essas informações mais detalhadas com gráficos e painéis para trabalhá-las com o objetivo de obter insights de forma rápida.

1.2. Objetivos

Tornar as informações sobre os benefícios concedidos pela previdência social mais fácil de serem analisadas e visualizadas fornecendo um dashboard interativo.

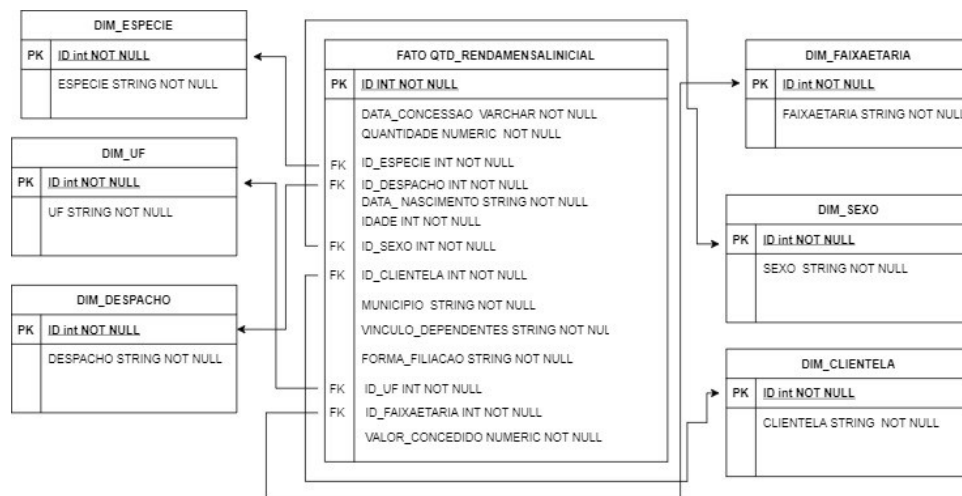
1.3. Público alvo

Pessoas interessadas em assuntos relacionados a previdência social, que tenham acesso à internet, redes sociais, e que estejam por dentro das notícias e se sintam dispostas a análises mais profundas do que o governo lhes oferece.

2. Modelo de Dados

2.1. Modelo Dimensional

O diagrama do modelo dimensional oferece uma visão sobre como os dados serão organizados para atender às necessidades identificadas por meio dos objetivos do projeto.



2.2. Fatos e Dimensões

Tabela Fato “QTD_RENDAMENSALINICIAL”: contém os registros da quantidade renda recebida de acordo com as colunas: ID; DATA_CONCESSAO; QUANTIDADE; ID_ESPECIE; ID_FAIXAETARIA; ID_DESPACHO; DATA_NASCIMENTO; ID_SEXO; ID_CLIENTELA; MUNICIPIO; VINCULO_DEPENDENTES; FORMA_FILIACAO, VALOR_CONCEDIDO; IDADE E ID_UF.

As dimensões são 5, divididas em:

“DIM_ESPECIE”, onde contém os registros do tipo do benefício. Exemplo: “Pensão por morte previdenciária”, “Auxílio salário maternidade”, “Auxílio doença previdenciário”, entre outros. Representada pelas colunas ID e ESPECIE.

“DIM_UF”, nessa tabela consta os estados brasileiros. Exemplo: Bahia, São Paulo, etc. Representada pelas colunas ID e UF.

“DIM_DESPACHO”, nessa tabela consta a forma de concessão. Exemplo: Concessão normal, decorrente de ação judicial, entre outras. Representada pelas colunas ID e DESPACHO.

“DIM_SEXO”, representa o sexo do beneficiário. Exemplo: feminino ou masculino. Representada pelas colunas ID e SEXO.

“DIM_CLIENTELA”, consta a área geográfica em que o beneficiário reside. Exemplo: urbano ou rural. Representada pelas colunas ID e CLIENTELA.

“DIM_FAIXAETARIA”, onde consta a divisão das idades. Exemplo: 15-25, 26-39.

3. Integração, Tratamento e Carga de Dados

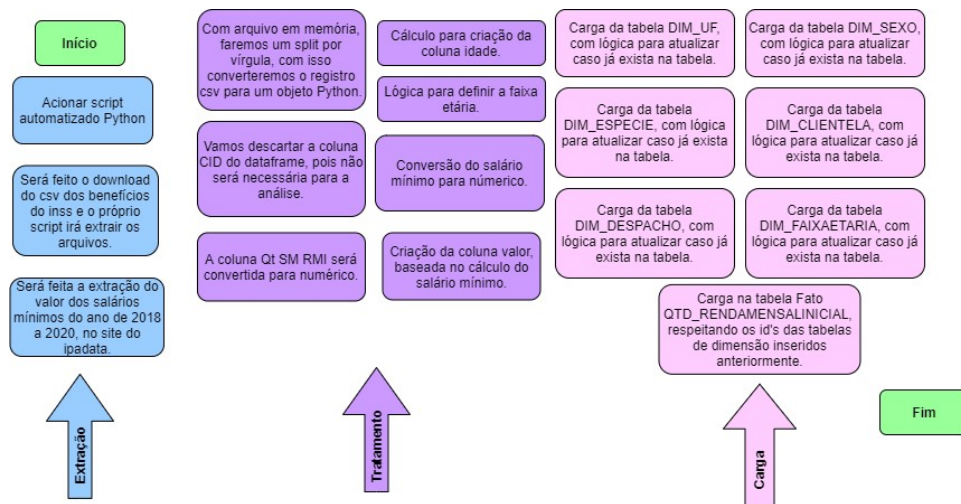
O processo de extração, tratamento e carga de dados será feito, nos arquivos do ano de 2018, 2019 e 2020. Com script Python, utilizando biblioteca Pandas, o arquivo será acessado no site: <https://dados.gov.br/dataset/inss-beneficios-concedidos> onde fará o download dos arquivos, também será feito download da base do salário mínimo no site: <http://www.ipeadata.gov.br/exibeserie.aspx?stub=1&serid1739471028=1739471028> finalizando a etapa de extração.

Para iniciar a etapa de tratamento com os arquivos em memória será feito um split por vírgula, convertendo o registro CSV para um objeto Python, descartando a coluna CID, pois não será interessante para a nossa análise. A coluna QT SM RMI deve ser convertida para numérico, pois usaremos a mesma para calcular o valor da concessão com base no salário mínimo criando a coluna Valor. Criação das colunas ID_FAIXAETARIA e IDADE com base na coluna DT NASCIMENTO do CSV, para melhor análise e enriquecimento dos dados no processo de ETL que será orquestrado pelo Python3.

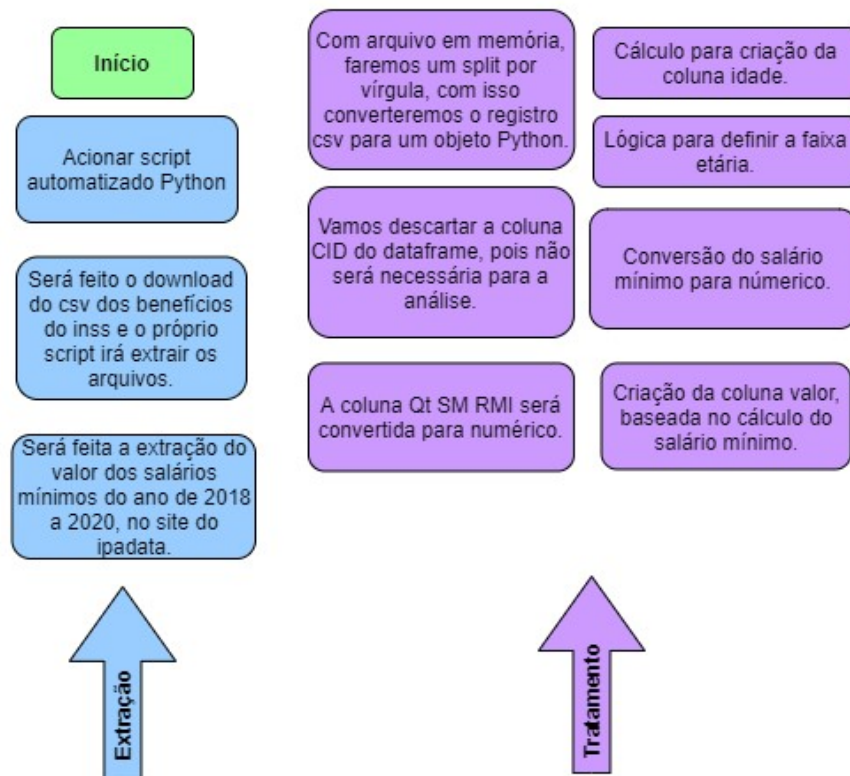
Já no processo de carga de dados, iremos inserir nas tabelas de dimensão, com a lógica para atualizar caso já exista na tabela e ao final do processo iremos

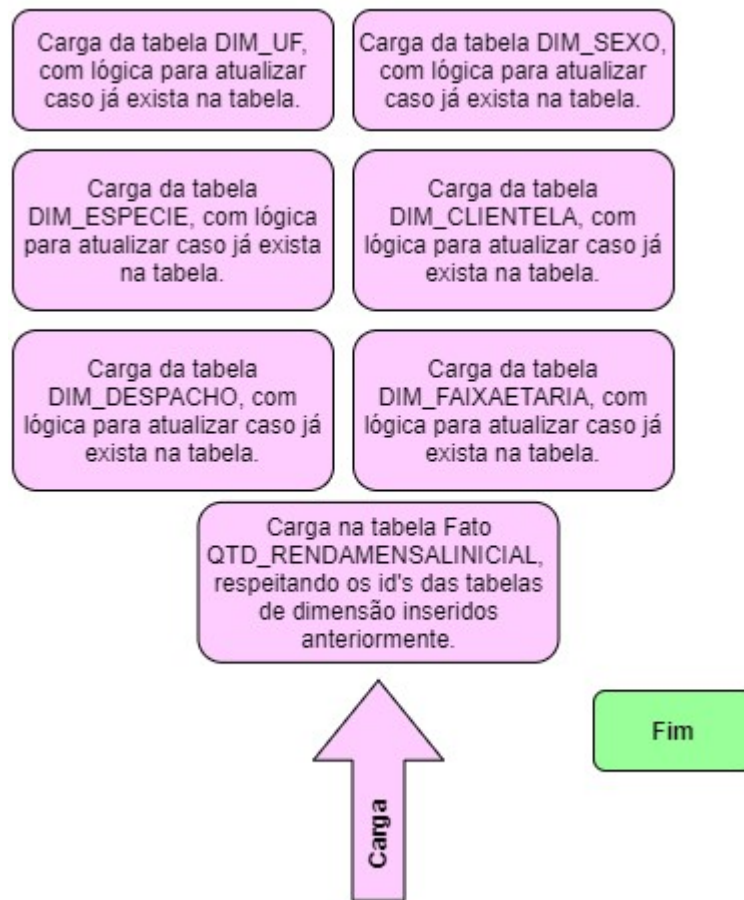
popular a tabela de fato respeitando os ids das tabelas de dimensão inseridos anteriormente. Todo o processo de carga será feito no banco de dados relacional MySQL.

Abaixo fluxograma das etapas:

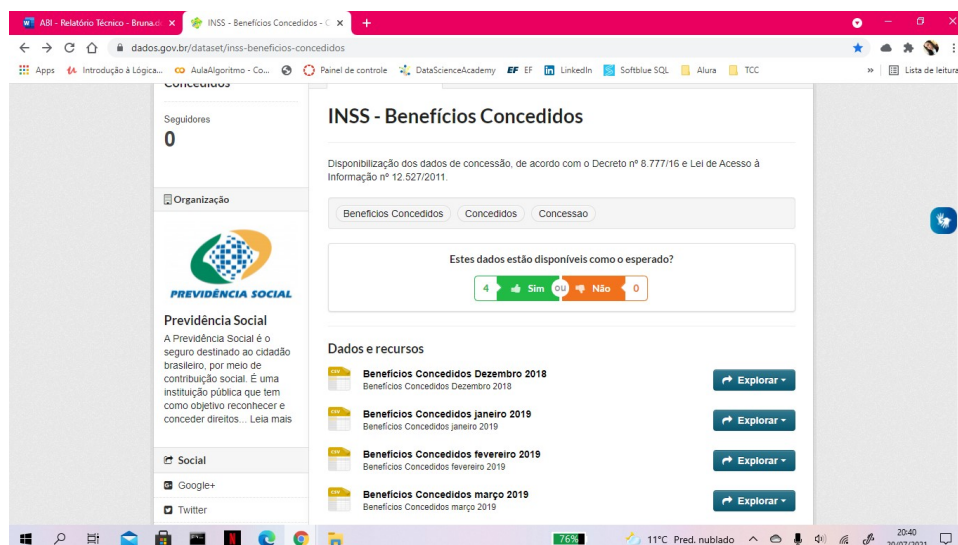


Fluxograma com zoom das etapas:





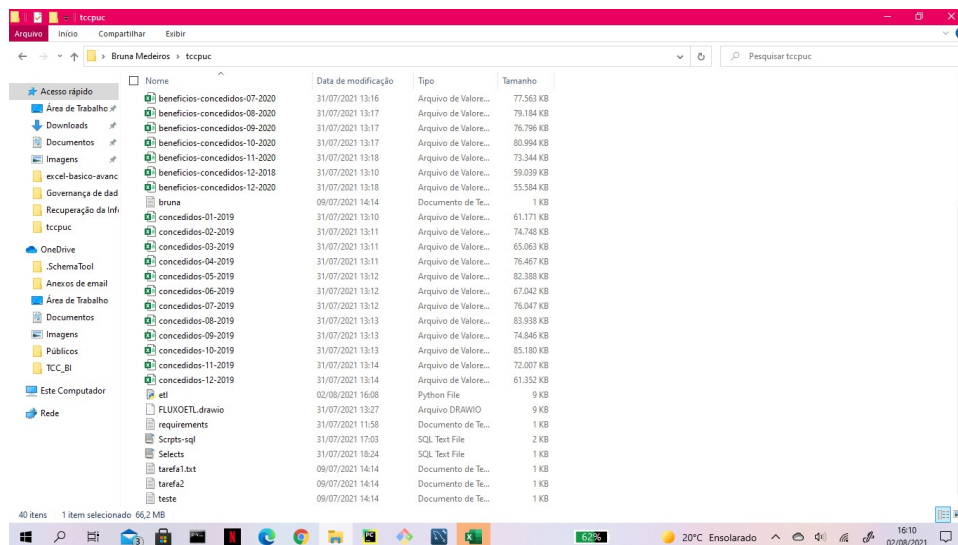
Abaixo print do site onde os arquivos foram extraídos:



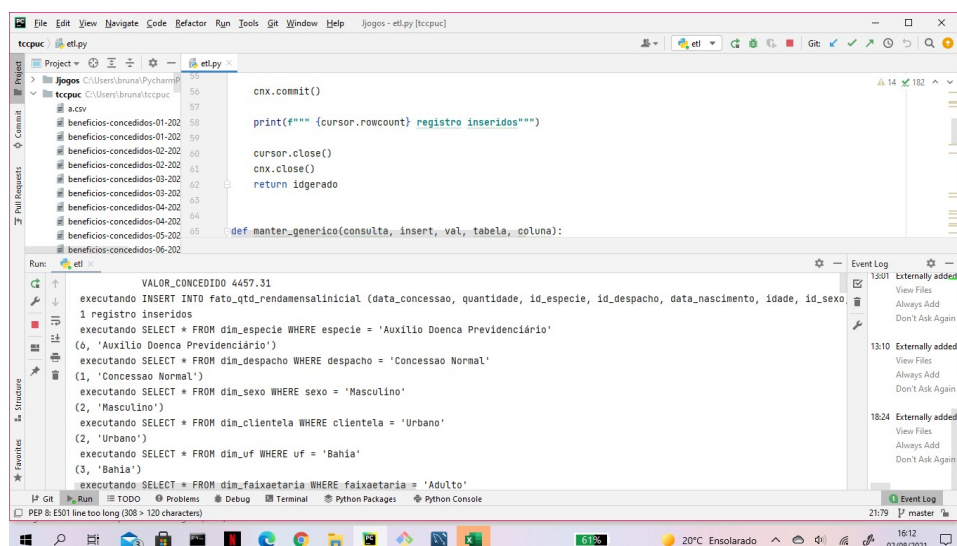
Scripts da criação das tabelas no MySql que faz parte da carga de dados:

```
4  create table dim_especie(  
5  id int not null auto_increment,  
6  especie varchar(1000) not null,  
7  primary key(id));  
8  
9  create table dim_sexo(  
10 id int not null auto_increment,  
11 sexo varchar(1000) not null,  
12 primary key(id));  
13  
14 create table dim_uf(  
15 id int not null auto_increment,  
16 uf varchar(1000) not null,  
17 primary key(id));  
18  
19 create table dim_despacho(  
20 id int not null auto_increment,  
21 despacho varchar(1000) not null,  
22 primary key(id));  
23  
24 create table dim_clientela(  
25 id int not null auto_increment,  
26 clientela varchar(1000) not null,  
27 primary key(id));  
28  
29 create table fato_qtd_rendamensalinicial(  
30 id int not null auto_increment,  
31 data_concessao varchar(1000) not null,  
32 quantidade numeric not null,  
33 id_especie int not null,  
34 id_despacho int not null,  
35 data_nascimento varchar(1000) not null,  
36 id_sexo int not null,  
37 id_clientela int not null,  
38 municipio varchar(1000) not null,  
39 vinculo_dependentes varchar(1000) not null,  
40 forma_filiacao varchar(1000) not null,  
41 id_uf int not null,  
42 primary key(id));  
43  
44 alter table fato_qtd_rendamensalinicial add foreign key(id_especie) references dim_especie(id);  
45 alter table fato_qtd_rendamensalinicial add foreign key(id_despacho) references dim_despacho(id);  
46 alter table fato_qtd_rendamensalinicial add foreign key(id_sexo) references dim_sexo(id);  
47 alter table fato_qtd_rendamensalinicial add foreign key(id_clientela) references dim_clientela(id);  
48 alter table fato_qtd_rendamensalinicial add foreign key(id_uf) references dim_uf(id);
```

Diretório do script etl.py após o download automatizado dos CSVs diretamente do site do governo:



Logs do script de etl.py em execução:



Link do github contendo os códigos fonte <https://github.com/brunamedeiros/tccpuc> .

Trecho do código criado para baixar os arquivos do site do INSS para a máquina local. Primeiro é realizado um crawl e com o HTML em memória utilizando a funcionalidade de ler CSS do BeautifulSoup, os links de download são extraídos facilmente e assim os CSVs são baixados pelo método baixar_arquivo.

```

def baixar_arquivo(url):
    local_filename = url.split('/')[-1]
    r = requests.get(url, verify=False)
    f = open(local_filename, 'wb')
    for chunk in r.iter_content(chunk_size=512 * 1024):
        if chunk: # filter out keep-alive new chunks
            f.write(chunk)
    f.close()
    return

def baixarcsvinss():
    page = requests.get(URL_INSS)
    soup = BeautifulSoup(page.text, 'html.parser')
    links = soup.find_all('a', attrs={"class": "resource-url-analytics"})

    for l in links:
        csv = l.attrs['href']
        csv = csv.split('url=')[1]
        print(f"" Baixando {csv} "")
        baixar_arquivo(csv)
        print(f"" Sucesso ao baixar {csv} "")

```

Trecho do código que acessa o site do governo do salário mínimo e realiza a raspagem do HTML, e converte essa informação para uma estrutura amigável em Python utilizando a biblioteca BeautifulSoup.

```

def baixarsalariominimo():
    print('baixando salario p1')
    page = requests.get(URL_SALARIO_MINIMO)
    print('baixando salario p2')
    soup = BeautifulSoup(page.text, 'html.parser')
    trs = soup.find_all('tr', attrs={"class": "dxgvDataRow"})
    print('baixando salario p3')
    salarios = []
    for tr in trs:
        tds = tr.find_all('td', attrs={"class": "dxgv"})
        if len(tds) == 2:
            td_ano = tds[0].text
            td_valor = tds[1].text
            if td_ano >= ANO_CORTE_SALARIO:
                td_valor_float = float(td_valor.replace(".", "").replace(",", "."))
                print(f"" ano {td_ano} valor {td_valor_float} "")
                salarios.append({"ano": td_ano, "valor": td_valor_float})
    return salarios

```


Parte do código contendo o coração da ETL, o método `etl_beneficios` recebe a lista de salários e o nome do CSV com isso utiliza a biblioteca Pandas para auxiliar no parse do mesmo. Para cada linha do arquivo carrega as tabelas no banco de dimensões e fato `fato_qtd_rendamensal` inicial no banco de dados MySQL.

```
def etl_beneficios(salarios):
    # csv = 'concedidos-10-2019.csv'
    csv = 'beneficios-concedidos-01-2020.csv'
    print(f"carregando para o banco de dados o csv de beneficios {csv}")
    try:
        data = pd.read_csv(csv, sep=';', low_memory=False, encoding='utf-8')
    except UnicodeDecodeError as e:
        data = pd.read_csv(csv, sep=';', low_memory=False, encoding='ISO-8859-1')
    print(data.columns)
    c_compe = 'Competência concessão'
    c_especie = 'Espécie'
    c_despacho = 'Despacho'
    c_dtnasc = 'Dt Nascimento'
    c_sex = 'Sexo.'
    c_cli = 'Clientela'
    c_mun = 'Mun Resid'
    c_vin = 'Vínculo dependentes'
    c_forma = 'Forma Filiação'
    c_uf = 'UF'
    c_qt = 'Qt SM RMI'
    for i, row in data.iterrows():
        data_concessao = ler_linha_csv(row, c_compe, 0)
        especie = ler_linha_csv(row, c_especie, 1)
        id_especie = manter_generico(SELECT_DIM_GENERIC, INSERT_DIM_GENERIC, (especie,), 'dim_especie', 'especie')

        despacho = ler_linha_csv(row, c_despacho, 4)
        id_despacho = manter_generico(SELECT_DIM_GENERIC, INSERT_DIM_GENERIC, (despacho,), 'dim_despacho', 'despacho')
        data_nascimento = ler_linha_csv(row, c_dtnasc, 5)
        sexo = ler_linha_csv(row, c_sex, 6)
        id_sexo = manter_generico(SELECT_DIM_GENERIC, INSERT_DIM_GENERIC, (sexo,), 'dim_sexo', 'sexo')
        clientela = ler_linha_csv(row, c_cli, 7)
        id_clientela = manter_generico(SELECT_DIM_GENERIC, INSERT_DIM_GENERIC, (clientela,), 'dim_clientela', 'clientela')
        municipio = ler_linha_csv(row, c_mun, 8)
        vinculo_dependentes = ler_linha_csv(row, c_vin, 9)
        forma_filiacao = ler_linha_csv(row, c_forma, 10)
        uf = ler_linha_csv(row, c_uf, 11)
        id_uf = manter_generico(SELECT_DIM_GENERIC, INSERT_DIM_GENERIC, (uf,), 'dim_uf', 'uf')
        qtd = ler_linha_csv(row, c_qt, 12)
        quantidade = float(qtd.replace(".", "").replace(",", "."))
        vlsalariominimodomes = get_no_array(salarios, "ano", converter_mesano(data_concessao))["valor"]
        valor_concedido = quantidade * vlsalariominimodomes
        idade = calcular_idade(data_nascimento)
        faixaetaria = regra_faixa_etaria(idade)
        id_faixaetaria = manter_generico(SELECT_DIM_GENERIC, INSERT_DIM_GENERIC, (faixaetaria,), 'dim_faixaetaria', 'faixaetaria')

    print(
        f""" DATA_CONCESSAO {data_concessao} ESPECIE {especie} DESPACHO {despacho} DATA_NASCIMENTO
        {data_nascimento} SEXO {sexo} CLIENTELA {clientela} MUNICIPIO {municipio}
        VINCULO_DEPE {vinculo_dependentes} FORMA_FI {forma_filiacao} UF {uf} QTD {quantidade}
        IDADE {idade}

        """
    )
    print(
        f""" DATA_CONCESSAO {data_concessao} ESPECIE {especie} DESPACHO {despacho} DATA_NASCIMENTO
        {data_nascimento} SEXO {sexo} CLIENTELA {clientela} MUNICIPIO {municipio}
        VINCULO_DEPE {vinculo_dependentes} FORMA_FI {forma_filiacao} UF {uf} QTD {quantidade}
        IDADE {idade}
        FAIXA_ETARIA {faixaetaria}
        VALOR_CONCEDIDO {valor_concedido}"""
    )
    manter_generico(None, INSERT_FATO, (data_concessao, quantidade, id_especie, id_despacho, data_nascimento, idade, id_sexo, id_clientela, id_municipio, id_vinculo_dependentes, id_forma_filiacao, id_uf, id_faixaetaria, valor_concedido))
```

Tabelas após a execução do script etl.py :

```
10 • select * from dim_clientela;
```

Result Grid		Filter Rows:	Edit:	Export/Impc
id	clientela			
1	Rural			
2	Urbano			
NULL	NULL			

```
13  
14 • select * from dim_despacho;  
15
```

Result Grid		Filter Rows:	Edit:	Export/Import:	Wrap Cell Content:
id	despacho				
1	Concessao Normal				
2	Conc. Base Artigo 27 Inciso II do Rbps				
3	Conc. Decorrente Revisao Administrativa				
4	Concessao Decorrente de Acao Judicial				
5	Concessao com Conversao Tempo de Servico				
6	Conc. com Base Artigo 35 da Lei 8213/91				
7	Concessao em Fase Recursal				
8	Conc. com Base no Artigo 183 do Rbps				
9	Conc. com Base no Artigo 180 do Rbps				
10	Conc. s/Verificacao da Perda Qualidade				
NULL	NULL				

```
20 • select * from dim_especie;
```

Result Grid		Filter Rows:	Edit:	Export/Import:	Wrap Cell
id	especie				
4	Auxilio Acidente Previdenciário				
5	Aposentadoria por Idade				
6	Auxilio Doenca Previdenciário				
7	Amp. Social Pessoa Portadora Deficiencia				
8	Aposentadoria por Tempo de Contribuição				
9	Aposentadoria Invalidez Previdenciária				
10	Amparo Social ao Idoso				
11	Aposentadoria Especial				
12	Aposent. Tempo de Serviço de Professor				
13	Aposent. Invalidez Acidente Trabalho				
14	Auxilio Reclusão				

```
21 • select * from dim_sexo;
```

Result Grid		Filter Rows:	Edit:	Export
id	sexo			
1	Feminino			
2	Masculino			
NULL	NULL			

22 • `select * from dim_uf;`

23

24

Result Grid | Filter Rows: | Edit: | Export

	id	uf
1	1	Alagoas
2	2	Amazonas
3	3	Bahia
4	4	Ceará
5	5	Mato Grosso do Sul
6	6	Espírito Santo
7	7	Goiás
8	8	Maranhão
9	9	Mato Grosso
10	10	Minas Gerais
11	11	Pará

25 • `select * from dim_faixaetaria;`

Result Grid | Filter Rows: | Edit: | Export

	id	faixaetaria
1	1	Adulto
2	2	Idoso
3	3	Criança
4	4	Adolescente
5	NULL	NULL

Tabela fato_qtd_rendamensalinicial :

26 • `select * from fato_qtd_rendamensalinicial;`

Result Grid | Filter Rows: | Edit: | Export/Import: | Wrap Cell Contents: | Fetch rows: |

	id	data_concessao	quantidade	id_especie	id_despacho	data_nascimento	idade	id_sexo	id_clientela	municipio	vinculo_dependentes	forma_filiacao
1	202001	1	1	1	1	25/01/1988	33	1	1	09149-MA-Bom Jesus das Selvas	Filho	Segurado Especial
2	202001	1	2	1	1	07/05/1981	40	2	2	02075-AL-Porto Real do Colégio	Não Informado	Empregado
3	202001	1	3	1	1	06/06/1959	62	1	1	22002-SE-Aquidabã	Cônjuge	Segurado Especial
4	202001	1	3	1	1	25/11/1974	47	1	1	09170-MA-Igarapé do Meio	Cônjuge	Segurado Especial
5	202001	1	3	1	1	14/09/1999	22	1	1	05002-CE-Acaraú	Filho	Segurado Especial
6	202001	1	1	1	1	12/09/1976	45	1	1	09054-MA-Itapecuru Mirim	Filho	Segurado Especial
7	202001	1	1	1	1	08/05/1992	29	1	1	05066-CE-Itapiúna	Filho	Segurado Especial
8	202001	1	3	1	1	29/12/1931	90	2	1	09121-MA-Turialga	Cônjuge	Segurado Especial
9	202001	1	1	1	1	11/11/1992	29	1	1	05095-CE-Novo Oriente	Filho	Segurado Especial
10	202001	1	1	1	1	15/08/1990	31	1	2	09097-MA-Sambaíba	Filho	Desempregado

4. Camada de Apresentação

4.1 Dashboard

Descreva os painéis de controle (dashboard) criados detalhando as visões estratégica, tática e operacional e como são interligadas. Descreva os indicadores e suas fórmulas de cálculo, as métricas, as dimensões de análise, as visualizações empregadas (mapas, gráficos, gauges, entre outros), bem como os filtros aplicáveis à solução, com base nas dimensões utilizadas.

4.2 Análises avançadas

Padrões e/ou Tendências: análises de dados criadas a partir da aplicação de técnicas de aprendizado de máquina. As análises podem ser independentes ou incorporadas no dashboard criado anteriormente.

5. Registros de Homologação

Testes da solução desenvolvida mostrando que o dado apresentado no dashboard é o mesmo dos sistemas fonte. Estes testes são apresentados, normalmente, por meio de consultas SQL feitas contra as fontes de dados e pelo confronto dos resultados com o que é exibido nas visualizações de dados apresentadas nos dashboards

5. Conclusões

Este item deve apresentar os seguintes itens: (1) análise crítica apresentando os achados mais relevantes nos dados, feitos a partir do uso do dashboard e das experiências adquiridas no processo de desenvolvimento; (2) proposta de intervenção, possíveis ações a serem tomadas por gestores do contexto analisado no intuito de melhorar o desempenho da organização; e, por fim, (3) as lições aprendidas no processo que possam ajudar nos projetos futuros.

Aponte, ainda, as limitações do trabalho e possíveis pontos de extensão para que outros possam utilizar como ponto de partida em novos projetos.

6. Links

Aqui devem ser disponibilizar os links para os artefatos criados (processos de carga, códigos-fonte, painéis, entre outros) e para o vídeo com a apresentação de 5 minutos e para o repositório contendo os códigos fontes ou os artefatos construídos no projeto.

REFERÊNCIAS

<http://www.ipeadata.gov.br/exibeserie.aspx?stub=1&serid1739471028=1739471028>

<https://dados.gov.br/dataset/inss-beneficios-concedidos> 22.07.21 20:27

Como um projeto aplicado, as referências, levantadas no processo de revisão bibliográfica, são opcionais e estimulamos que o responsável pelo projeto apresente aquilo que for relevante para complementar o entendimento do trabalho. Dessa forma, relacione-as de acordo com o modelo a seguir.

SOBRENOME DO AUTOR, Nome do autor. **Título do livro ou artigo.**
Cidade: Editora, ano.

SOBRENOME DO AUTOR, Nome do autor. **Título do livro ou artigo.**
Cidade: Editora, ano.

SOBRENOME DO AUTOR, Nome do autor. **Título do livro ou artigo.**
Cidade: Editora, ano.

SOBRENOME DO AUTOR, Nome do autor. **Título do livro ou artigo.**
Cidade: Editora, ano.

SOBRENOME DO AUTOR, Nome do autor. **Título do livro ou artigo.**
Cidade: Editora, ano.

SOBRENOME DO AUTOR, Nome do autor. **Título do livro ou artigo.**
Cidade: Editora, ano.

Fornecer um dashboard que facilite a análise sobre os benefícios concedidos da previdência social a partir do ano de 2009, com informações divididas por ano, valor, clientela (urbano, rural), sexo, quantidade benefícios concedidos, grupo/principais espécies (tipo do benefício: aposentadoria acidentária, por idade, invalidez, etc.).

Baseado em dados reais extraídos da base de dados do governo federal.