

1.0 PySpark - Instalando a biblioteca PySpark no Google Colab

```
1 ! pip install pyspark
```

```
Requirement already satisfied: pyspark in /usr/local/lib/python3.11/dist-packages (3.5.1)
Requirement already satisfied: py4j==0.10.9.7 in /usr/local/lib/python3.11/dist-packages (from pyspark) (0.10.9.7)
```

```
1 ! pip install findspark
```

```
Collecting findspark
  Downloading findspark-2.0.1-py2.py3-none-any.whl.metadata (352 bytes)
  Downloading findspark-2.0.1-py2.py3-none-any.whl (4.4 kB)
Installing collected packages: findspark
Successfully installed findspark-2.0.1
```

1.1 Importante as bibliotecas:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import plotly.express as px
6 import findspark
7 findspark.init()
8
9 from pyspark.sql import SparkSession
10 spark = SparkSession.builder.master("local[*]").getOrCreate()
11 from pyspark.sql.functions import concat, col, to_date, lpad, when, sum as _sum, round, avg

1 df = spark.read.csv('/content/PNAD_COVID_consolidado.csv', sep = ',', inferSchema = True, header = True)
2 # header para falar que tem cabeçalho
3 # inferSchema para forçar a rodar os dados importados
4
5 print('df.count : ', df.count()) # para mostrar o total de linhas
6 print('df. col ct : ', len(df.columns)) # para mostrar o total de colunas
7 print('df.columns : ', df.columns) # para listar todas as colunas
```

```
df.count : 1154279
df. col ct : 146
df.columns : ['Ano', 'UF', 'CAPITAL', 'RM_RIDE', 'V1008', 'V1012', 'V1013', 'V1016', 'Estrato', 'UPA', 'V1022', 'V1023', 'V1030',
```

2.0 Usando SQL no PySpark:

```
1 df.describe().show() # para mostrar as estatísticas do DataFrame
```

```
...
```

```
1 df.printSchema()
```

```
...
```

✓ 2.1 Remover as duplicatas:

```
1 df = df.dropDuplicates()
2 print('Total de linhas:', df.count())
3 print('Total de colunas: ', df.columns)
```

...

✓ 2.1 Renomeando as colunas:

```
1 df1 = df \
2     .withColumnRenamed('UF', 'unidade_federacao') \
3     .withColumnRenamed('CAPITAL', 'capital_estado') \
4     .withColumnRenamed('RM_RIDE', 'regiao_metropolitana') \
5     .withColumnRenamed('V1008', 'numero_domicilio') \
6     .withColumnRenamed('V1012', 'semana_no_mes') \
7     .withColumnRenamed('V1013', 'mes_pesquisa') \
8     .withColumnRenamed('V1016', 'numero_entrevista_domicilio') \
9     .withColumnRenamed('Estrato', 'estrato_amostrado') \
10    .withColumnRenamed('UPA', 'unidade_primaria_amostragem') \
11    .withColumnRenamed('V1023', 'tipo_area') \
12    .withColumnRenamed('V1030', 'projecao_populacao') \
13    .withColumnRenamed('V1031', 'peso_sem_pos_estratificacao') \
14    .withColumnRenamed('V1032', 'peso_com_pos_estratificacao') \
15    .withColumnRenamed('posest', 'dominio_projecao') \
16    .withColumnRenamed('A001', 'numero_ordem_morador') \
17    .withColumnRenamed('A001A', 'condicao_no_domicilio') \
18    .withColumnRenamed('A001B1', 'dia_nascimento') \
19    .withColumnRenamed('A001B2', 'mes_nascimento') \
20    .withColumnRenamed('A001B3', 'ano_nascimento') \
21    .withColumnRenamed('A002', 'idade_morador') \
22    .withColumnRenamed('A003', 'sexo') \
23    .withColumnRenamed('A004', 'cor_ou_raca') \
24    .withColumnRenamed('A005', 'escolaridade') \
25    .withColumnRenamed('B0011', 'semana_passada_febre') \
```

```
26 .withColumnRenamed('B0012', 'semana_passada_tosse') \
27 .withColumnRenamed('B0013', 'semana_passada_dor_garganta') \
28 .withColumnRenamed('B0014', 'semana_passada_dificuldade_respirar') \
29 .withColumnRenamed('B0015', 'semana_passada_dor_cabeca') \
30 .withColumnRenamed('B0016', 'semana_passada_dor_peito') \
31 .withColumnRenamed('B0017', 'semana_passada_nausea') \
32 .withColumnRenamed('B0018', 'semana_passada_nariz_entupido_escorrendo') \
33 .withColumnRenamed('B0019', 'semana_passada_fadiga') \
34 .withColumnRenamed('B00110', 'semana_passada_dor_olhos') \
35 .withColumnRenamed('B00111', 'semana_passada_perda_cheiro_sabor') \
36 .withColumnRenamed('B00112', 'semana_passada_dor_muscular') \
37 .withColumnRenamed('B002', 'procurou_estab_saude') \
38 .withColumnRenamed('B0031', 'recup_sintomas_ficou_casa') \
39 .withColumnRenamed('B0032', 'recup_sintomas_ligou_prof_saude') \
40 .withColumnRenamed('B0033', 'recup_sintomas_tomou_remedio_conta') \
41 .withColumnRenamed('B0034', 'recup_sintomas_tomou_remedio_medico') \
42 .withColumnRenamed('B0035', 'recup_sintomas_visita_prof_sus') \
43 .withColumnRenamed('B0036', 'recup_sintomas_visita_prof_particular') \
44 .withColumnRenamed('B0041', 'atendimento_posto_ubs') \
45 .withColumnRenamed('B0042', 'atendimento_pronto_socorro_sus') \
46 .withColumnRenamed('B0043', 'atendimento_hospital_sus') \
47 .withColumnRenamed('B0044', 'atendimento_ambulatorio_privado') \
48 .withColumnRenamed('B0045', 'atendimento_pronto_socorro_privado') \
49 .withColumnRenamed('B0046', 'atendimento_hospital_privado') \
50 .withColumnRenamed('B005', 'internado_um_dia_ou_mais') \
51 .withColumnRenamed('B006', 'sedado_entubado_ventilado') \
52 .withColumnRenamed('B007', 'tem_plano_saude') \
53 .withColumnRenamed('C001', 'trabalhou_ou_bico_semana') \
54 .withColumnRenamed('C002', 'afastado_trabalho') \
55 .withColumnRenamed('C003', 'motivo_afastamento') \
56 .withColumnRenamed('C004', 'continuou_remunerado') \
57 .withColumnRenamed('C005', 'tempo_afastado') \
58 .withColumnRenamed('C0051', 'meses_afastado_menos_um_ano') \
59 .withColumnRenamed('C0052', 'meses_afastado_um_a_dois_anos') \
60 .withColumnRenamed('C0053', 'anos_afastado_dois_ou_mais') \
61 .withColumnRenamed('C006', 'mais_de_um_trabalho') \
62 .withColumnRenamed('C007', 'tipo_ocupacao') \
63 .withColumnRenamed('C007A', 'area_trabalho') \
64 .withColumnRenamed('C007B', 'carteira_assinada_ou_estatutario') \
65 .withColumnRenamed('C007C', 'cargo_funcao_trabalho') \
66 .withColumnRenamed('C007D', 'atividade_empresa') \
67 .withColumnRenamed('C007E', 'numero_empregados_negocio') \
68 .withColumnRenamed('C008', 'horas_trabalhadas_habitual') \
69 .withColumnRenamed('C009', 'horas_trabalhadas_semana') \
70 .withColumnRenamed('C010', 'recebia_habitualmente') \
71 .withColumnRenamed('C0101', 'recebia_em_dinheiro') \
72 .withColumnRenamed('C01011', 'faixa_rendimento_dinheiro') \
73 .withColumnRenamed('C01012', 'valor_rendimento_dinheiro') \
74 .withColumnRenamed('C0103', 'recebia_beneficios') \
75 .withColumnRenamed('C0104', 'nao_remunerado') \
76 .withColumnRenamed('C011A', 'recebia_efetivamente') \
77 .withColumnRenamed('C011A1', 'recebia_efetivamente_dinheiro') \
78 .withColumnRenamed('C011A11', 'faixa_rendimento_efetivo_dinheiro') \
79 .withColumnRenamed('C011A12', 'valor_rendimento_efetivo_dinheiro') \
80 .withColumnRenamed('C012', 'trabalho_mesmo_local') \
81 .withColumnRenamed('C013', 'trabalho_remoto') \
82 .withColumnRenamed('C014', 'contribui_inss') \
83 .withColumnRenamed('C015', 'procurou_trabalho') \
84 .withColumnRenamed('C016', 'motivo_nao_procurou_trabalho') \
85 .withColumnRenamed('C017A', 'gostaria_trabalhar') \
86 .withColumnRenamed('D0011', 'recebeu_aposentadoria_pensao') \
87 .withColumnRenamed('D0013', 'valor_aposentadoria_pensao') \
88 .withColumnRenamed('D0021', 'recebeu_pensao_doacao_mesada') \
89 .withColumnRenamed('D0023', 'valor_pensao_doacao_mesada') \
90 .withColumnRenamed('D0031', 'recebeu_bolsa_familia') \
91 .withColumnRenamed('D0033', 'valor_bolsa_familia') \
92 .withColumnRenamed('D0041', 'recebeu_bpc_loas') \
93 .withColumnRenamed('D0043', 'valor_bpc_loas') \
94 .withColumnRenamed('D0051', 'recebeu_auxilio_emergencial') \
95 .withColumnRenamed('D0053', 'valor_auxilio_emergencial') \
96 .withColumnRenamed('D0061', 'recebeu_seguro_desemprego') \
97 .withColumnRenamed('D0063', 'valor_seguro_desemprego') \
98 .withColumnRenamed('D0071', 'recebeu_outros_rendimentos') \
99 .withColumnRenamed('D0073', 'valor_outros_rendimentos') \
100 .withColumnRenamed('F001', 'condicao_domicilio') \
101 .withColumnRenamed('F0021', 'valor_aluguel_pago') \
102 .withColumnRenamed('F0022', 'faixa_valor_aluguel') \
103 .withColumnRenamed('F0061', 'quem_respondeu_questionario') \
```

```
104     .withColumnRenamed('F006', 'numero_ordem_informante')
105 df1.show(10)
```

...

✓ 2.3 Agrupando coluna dia_nascimento, mes_nascimento e ano_nascimento para colocar apenas a data de nascimento:

```
1 # vamos formatar o número para que ele sempre tenha dois dígitos:
2 df1 = df1.withColumn('dia_formatado', lpad(col('dia_nascimento').cast('string'), 2, '0'))
3 df1 = df1.withColumn('mes_formatado', lpad(col('mes_nascimento').cast('string'), 2, '0'))

1 # agora, agrupando em uma única coluna:
2 df2 = df1.withColumn('data_nascimento', concat(col('dia_formatado'), col('mes_formatado'), col('ano_nascimento')))
```

```
1 df2.show(10)
```

...

```
1 df3 = df2.withColumn('data_nasc', to_date(col('data_nascimento'), 'ddMMyyyy'))
```

```
1 df3.show(10)
```

...

✓ 2.4 Apagando as colunas:

- dia_nascimento
- mes_nascimento
- ano_nascimento

- dia_formatado
- mes_formatado

```
1 # vamos verificar o número de colunas antes:
2 print('df3. col ct : ', len(df3.columns))
```

...

```
1 # apagando as colunas:
2 df4 = df3.drop(*['dia_nascimento', 'mes_nascimento', 'ano_nascimento', 'dia_formatado', 'mes_formatado'])
3 df4.show(5)
```

...

```
1 # apagando a coluna data_nascimento:
2 df5 = df4.drop(*['data_nascimento', 'numero_formatado'])
3 df5.show(5)
```

...

```
1 # verificando o número de colunas depois:
2 print('df5. col ct : ', len(df5.columns))
```

...

✓ 2.5 Seleccionando algumas colunas:

✓ 2.5.1 Contando o número de pesquisa realizada em cada mês:

```
1 # Agrupar pelo mês da pesquisa e contar o número de entrevistas
2 numero_pesquisa_meses = df5.groupBy('mes_pesquisa').count().orderBy('mes_pesquisa')
3 numero_pesquisa_meses.show()
```

...

```
1 # contando o número total de pessoas:
2 total_pessoas = df5.count()
3 print(f"Total de pessoas: {total_pessoas}")
```

...

✓ 2.5.2 Agrupando por cor/raça:

```
1 # criando uma coluna com o nome da cor/raça df_raca = df.withColumn(
2 df_raca = df5.withColumn(
3     'cor_ou_raca_nome',
4     when(col('cor_ou_raca') == '1', 'Branca')
```

```

5     .when(col('cor_ou_raca') == '2', 'Preta')
6     .when(col('cor_ou_raca') == '3', 'Amarela')
7     .when(col('cor_ou_raca') == '4', 'Parda')
8     .when(col('cor_ou_raca') == '5', 'Indígena')
9     .otherwise('Ignorado')
10 )
11
12 # Contar o número de pessoas por cor/raça
13 tabela_raca = df_raca.groupBy('cor_ou_raca_nome').count().orderBy('count', ascending=False)
14 tabela_raca.show()
...

```

✓ 2.5.3 Agrupando por escolaridade:

```

1 # criando uma coluna descritiva de escolaridade
2 df_escolaridade = df5.withColumn(
3     'escolaridade_nome',
4     when(col('escolaridade') == '1', 'Sem instrução')
5     .when(col('escolaridade') == '2', 'Fundamental incompleto')
6     .when(col('escolaridade') == '3', 'Fundamental completo')
7     .when(col('escolaridade') == '4', 'Médio incompleto')
8     .when(col('escolaridade') == '5', 'Médio completo')
9     .when(col('escolaridade') == '6', 'Superior incompleto')
10    .when(col('escolaridade') == '7', 'Superior completo')
11    .when(col('escolaridade') == '8', 'Pós-graduação')
12    .otherwise('Ignorado')
13 )
14
15 # Conta o número de pessoas por escolaridade
16 tabela_escolaridade = df_escolaridade.groupBy('escolaridade_nome').count().orderBy('count', ascending=False)
17 tabela_escolaridade.show()
...

```

✓ 2.5.4 Agrupando por estado:

```

1 # criando uma coluna com o nome do estado
2 df_estado = df5.withColumn(
3     'estado_nome',
4     when(col('unidade_federacao') == '11', 'Rondônia')
5     .when(col('unidade_federacao') == '12', 'Acre')
6     .when(col('unidade_federacao') == '13', 'Amazonas')
7     .when(col('unidade_federacao') == '14', 'Roraima')
8     .when(col('unidade_federacao') == '15', 'Pará')
9     .when(col('unidade_federacao') == '16', 'Amapá')
10    .when(col('unidade_federacao') == '17', 'Tocantins')
11    .when(col('unidade_federacao') == '21', 'Maranhão')
12    .when(col('unidade_federacao') == '22', 'Piauí')
13    .when(col('unidade_federacao') == '23', 'Ceará')
14    .when(col('unidade_federacao') == '24', 'Rio Grande do Norte')
15    .when(col('unidade_federacao') == '25', 'Paraíba')
16    .when(col('unidade_federacao') == '26', 'Pernambuco')
17    .when(col('unidade_federacao') == '27', 'Alagoas')
18    .when(col('unidade_federacao') == '28', 'Sergipe')
19    .when(col('unidade_federacao') == '29', 'Bahia')

```

```

20     .when(col('unidade_federacao') == '31', 'Minas Gerais')
21     .when(col('unidade_federacao') == '32', 'Espírito Santo')
22     .when(col('unidade_federacao') == '33', 'Rio de Janeiro')
23     .when(col('unidade_federacao') == '35', 'São Paulo')
24     .when(col('unidade_federacao') == '41', 'Paraná')
25     .when(col('unidade_federacao') == '42', 'Santa Catarina')
26     .when(col('unidade_federacao') == '43', 'Rio Grande do Sul')
27     .when(col('unidade_federacao') == '50', 'Mato Grosso do Sul')
28     .when(col('unidade_federacao') == '51', 'Mato Grosso')
29     .when(col('unidade_federacao') == '52', 'Goiás')
30     .when(col('unidade_federacao') == '53', 'Distrito Federal')
31 )
32
33 # Contar o número de pessoas por estado
34 tabela_estado = df_estado.groupBy('estado_nome').count().orderBy('count', ascending=False)
35 tabela_estado.show(27)

```

...

✓ 2.5.5 Agora, vamos agrupar por sintomas:

```

1 # Lista de sintomas e seus nomes descritivos
2 sintomas = [
3     ('semana_passada_febre', 'Febre'),
4     ('semana_passada_tosse', 'Tosse'),
5     ('semana_passada_dor_garganta', 'Dor de garganta'),
6     ('semana_passada_dificuldade_respirar', 'Dificuldade para respirar'),
7     ('semana_passada_dor_cabeca', 'Dor de cabeça'),
8     ('semana_passada_dor_peito', 'Dor no peito'),
9     ('semana_passada_nausea', 'Náusea'),
10    ('semana_passada_nariz_entupido_escorrendo', 'Nariz entupido ou escorrendo'),
11    ('semana_passada_fadiga', 'Fadiga'),
12    ('semana_passada_dor_olhos', 'Dor nos olhos'),
13    ('semana_passada_perda_cheiro_sabor', 'Perda de cheiro ou sabor'),
14    ('semana_passada_dor_muscular', 'Dor muscular')
15 ]
16
17 # Número total de pessoas (para porcentagem)
18 total_pessoas
19
20 # Lista para armazenar resultados
21 resultados = []
22
23 for codigo, nome in sintomas:
24     # Conta quantas pessoas responderam '1' (Sim) para o sintoma
25     count_sim = df5.filter(col(codigo) == '1').count()
26     perc_sim = __builtins__.round((count_sim / total_pessoas) * 100, 2)
27     resultados.append((nome, count_sim, float(perc_sim)))
28
29 # Cria um DataFrame com os resultados

```

```
30 result_df = spark.createDataFrame(resultados, ['Sintoma', 'Número absoluto', 'Porcentagem (%)'])
31 result_df.show(truncate=False)
```

...

```
1 from pyspark.sql.functions import col
2
3 # Lista de sintomas e seus nomes descritivos
4 sintomas = [
5     ('semana_passada_febre', 'Febre'),
6     ('semana_passada_tosse', 'Tosse'),
7     ('semana_passada_dor_garganta', 'Dor de garganta'),
8     ('semana_passada_dificuldade_respirar', 'Dificuldade para respirar'),
9     ('semana_passada_dor_cabeca', 'Dor de cabeça'),
10    ('semana_passada_dor_peito', 'Dor no peito'),
11    ('semana_passada_nausea', 'Náusea'),
12    ('semana_passada_nariz_entupido_escorrendo', 'Nariz entupido ou escorrendo'),
13    ('semana_passada_fadiga', 'Fadiga'),
14    ('semana_passada_dor_olhos', 'Dor nos olhos'),
15    ('semana_passada_perda_cheiro_sabor', 'Perda de cheiro ou sabor'),
16    ('semana_passada_dor_muscular', 'Dor muscular')
17 ]
18
19 # Número total de pessoas (para porcentagem)
20 total_pessoas
21
22 # Lista para armazenar resultados
23 resultados = []
24
25 for codigo, nome in sintomas:
26     # Conta quantas pessoas responderam '1' (Sim) para o sintoma
27     count_sim = df5.filter(col(codigo) == '1').count()
28     perc_sim = __builtins__.round((count_sim / total_pessoas) * 100, 2)
29
30     # Conta quantas pessoas responderam '2' (Não) para o sintoma
31     count_nao = df5.filter(col(codigo) == '2').count()
32     perc_nao = __builtins__.round((count_nao / total_pessoas) * 100, 2)
33
34     resultados.append((nome, count_sim, perc_sim, count_nao, perc_nao))
35
36 # Cria um DataFrame com os resultados
37 result_df = spark.createDataFrame(resultados, [
38     'Sintoma',
39     'Número (Sim)', 'Porcentagem (Sim) %',
40     'Número (Não)', 'Porcentagem (Não) %'
41 ])
42 result_df.show(truncate=False)
```

...

✓ 3.0 Construindo alguns gráficos:

✓ 3.1 Gráfico do número de pesquisa realizada em cada mês:

```
1 # Converta o DataFrame Spark para pandas
2 numero_pesquisa_meses_pd = numero_pesquisa_meses.toPandas()
3
4 # Gera o gráfico de barras
5 plt.figure(figsize=(10, 5))
6 plt.bar(numero_pesquisa_meses_pd['mes_pesquisa'], numero_pesquisa_meses_pd['count'], color='blue', width=0.3)
7
8 plt.xlabel('Mês da Pesquisa')
9 plt.ylabel('Número de Entrevistas')
10 plt.title('Número de Entrevistas por Mês')
11 plt.grid(True, linestyle='--', color='black', linewidth=0.8)
12 plt.gca().set_facecolor('lightgray')
13 plt.tight_layout()
14 plt.show()
```

...

Clique duas vezes (ou pressione "Enter") para editar

✓ 3.2 Gráfico por cor/raça:

```
1 # Convertendo para pandas para plotar o gráfico
2 tabela_raca_pd = tabela_raca.toPandas()
3
4 # Calcular porcentagem
5 total = tabela_raca_pd['count'].sum()
6 tabela_raca_pd['porcentagem'] = (tabela_raca_pd['count'] / total) * 100
7
8 # Plotar o gráfico
9 plt.figure(figsize=(12, 6))
10 plt.bar(tabela_raca_pd['cor_ou_raca_nome'], tabela_raca_pd['porcentagem'], color='blue',width=0.3)
11 plt.xlabel('Cor ou Raça')
12 plt.ylabel('Porcentagem (%)')
13 plt.title('Distribuição por Cor ou Raça')
14 plt.xticks(rotation=45, ha='right')
15 plt.grid(True, linestyle='--', color='black', linewidth=0.8)
16 plt.gca().set_facecolor('lightgray') # Fundo cinza
17 plt.tight_layout()
18 plt.show()
```

...

✓ 3.3 Gráfico por escolaridade:

```
1 # Converter o DataFrame Spark para pandas
2 tabela_escolaridade_pd = tabela_escolaridade.toPandas()
3
4 # Gráfico de barras com barras finas
5 plt.figure(figsize=(12, 6))
6 plt.bar(tabela_escolaridade_pd['escolaridade_nome'], tabela_escolaridade_pd['count'], color='blue', width=0.3)
7 plt.xlabel('Escolaridade')
8 plt.ylabel('Número de Pessoas')
9 plt.title('Distribuição por Escolaridade')
10 plt.xticks(rotation=45, ha='right')
11 plt.grid(True, linestyle='--', color='black', linewidth=0.8)
12 plt.gca().set_facecolor('lightgray')
13 plt.tight_layout()
14 plt.show()
```

...

✓ 3.4 Gráfico por estado:

```
1 import matplotlib.pyplot as plt
2
3 # Converter o DataFrame Spark para pandas
4 tabela_estado_pd = tabela_estado.toPandas()
```

```
5
6 # Gráfico de barras com barras finas
7 plt.figure(figsize=(14, 6))
8 plt.bar(tabela_estado_pd['estado_nome'], tabela_estado_pd['count'], color='blue')
9 plt.xlabel('Estado')
10 plt.ylabel('Número de Pessoas')
11 plt.title('Distribuição de Pessoas por Estado')
12 plt.xticks(rotation=45, ha='right')
13 plt.grid(True, linestyle='--', color='black', linewidth=0.8)
14 plt.gca().set_facecolor('lightgray') # Fundo cinza
15 plt.tight_layout()
16 plt.show()

...
```

✓ Gráfico por sintomas:

```
1 # Converta o DataFrame Spark para pandas
2 result_pd = result_df.toPandas()
3
4 # Crie o gráfico de barras para a porcentagem de pessoas que tiveram o sintoma
5 plt.figure(figsize=(12, 6))
6 plt.bar(result_pd['Sintoma'], result_pd['Porcentagem (Sim) %'], color='blue', width=0.3)
7 plt.xticks(rotation=45, ha='right')
8 plt.xlabel('Sintoma')
9 plt.ylabel('Porcentagem (%)')
10 plt.title('Porcentagem de Pessoas com Sintomas na Semana Passada')
11 plt.grid(True, linestyle='--', color='gray', linewidth=0.8)
12 plt.gca().set_facecolor('lightgray')
13 plt.tight_layout()
14 plt.show()
15

...
```

```

1 import numpy as np
2
3 x = np.arange(len(result_pd['Sintoma']))
4 width = 0.35
5
6 fig, ax = plt.subplots(figsize=(14, 6))
7 bars1 = ax.bar(x - width/2, result_pd['Porcentagem (Sim) %'], width, label='Sim', color='blue')
8 bars2 = ax.bar(x + width/2, result_pd['Porcentagem (Não) %'], width, label='Não', color='green')
9
10 ax.set_xlabel('Sintoma')
11 ax.set_ylabel('Porcentagem (%)')
12 ax.set_title('Porcentagem de Pessoas com e sem Sintomas na Semana Passada')
13 ax.set_xticks(x)
14 ax.set_xticklabels(result_pd['Sintoma'], rotation=45, ha='right')
15 ax.legend()
16 plt.tight_layout()
17 plt.grid(True, linestyle='--', color='gray', linewidth=0.8)
18 plt.gca().set_facecolor('lightgray')
19 plt.show()
20
...

```

✓ 4.0 Analisando os dados:

✓ 4.1 Analisando se os sintomas tem relação com escolaridade e/ou renda:

```

1 # Lista das colunas de sintomas (conforme o dicionário PNAD COVID-19)
2 colunas_sintomas = [
3     'semana_passada_febre', 'semana_passada_tosse', 'semana_passada_dor_garganta',
4     'semana_passada_dificuldade_respirar', 'semana_passada_dor_cabeca', 'semana_passada_dor_peito',
5     'semana_passada_nausea', 'semana_passada_nariz_entupido_escorrendo', 'semana_passada_fadiga',
6     'semana_passada_dor_olhos', 'semana_passada_perda_cheiro_sabor', 'semana_passada_dor_muscular'
7 ]
8
9 # Expressão para somar sintomas com valor 1 (Sim)
10 expr_sintomas = sum([when(col(c) == 1, 1).otherwise(0) for c in colunas_sintomas])
11
12 # Adiciona a coluna descritiva de escolaridade e o número de sintomas
13 df_novo = df5.withColumn(
14     'escolaridade_nome',
15     when(col('escolaridade') == 1, 'Sem instrução')
16     .when(col('escolaridade') == 2, 'Fundamental incompleto')
17     .when(col('escolaridade') == 3, 'Fundamental completo')

```

```
18     .when(col('escolaridade') == 4, 'Médio incompleto')
19     .when(col('escolaridade') == 5, 'Médio completo')
20     .when(col('escolaridade') == 6, 'Superior incompleto')
21     .when(col('escolaridade') == 7, 'Superior completo')
22     .when(col('escolaridade') == 8, 'Pós-graduação')
23     .otherwise('Ignorado')
24 ).withColumn(
25     'numero_sintomas', expr_sintomas
26 )
27
28 # Seleciona as colunas desejadas
29 df_final = df_novo.select('idade_morador', 'escolaridade_nome', 'numero_sintomas')
30
31 df_final.show(10)

...
```

1 Comece a programar ou gere código com IA.