

FIAP + alura

Turma - 2TDAT

FASE 3 - BIG DATA
Grupo 82

Jonathan Abner Jerônimo de Freitas

jonathanabner2015@gmail.com

Osvaldo Caio Oliveira dos Santos

osvaldocaio@hotmail.com

Ranielli Santos

raniellidos@hotmail.com

Thalita Mendes Maina Begliomini

thalitamaina@gmail.com

_ROTEIRO

- 01 - Definições
- 02 - Sobre o Projeto
- 03 - Dashboard
- 04 - Considerações Finais
- 05 - Refências





DEFINIÇÕES

PROBLEMA

CORONAVÍRUS

01

_ DEFINIÇÃO DO PROBLEMA

Imagine agora que você foi contratado(a) como Expert em Data Analytics por um grande hospital para entender como foi o comportamento da população na época da pandemia da COVID-19 e quais indicadores seriam importantes para o planejamento, caso haja um novo surto da doença.

Apesar de ser contratado(a) agora, a sua área observou que a utilização do estudo do PNAD-COVID 19 do IBGE seria uma ótima base para termos boas respostas ao problema proposto, pois são dados confiáveis.

Porém, não será necessário utilizar todas as perguntas realizadas na pesquisa para enxergar todas as oportunidades ali postas. É sempre bom ressaltar que há dados triviais que precisam estar no projeto, pois auxiliam muito na análise dos dados:

/ Características clínicas dos sintomas;

/ Características da população;

/ Características econômicas da sociedade.

_ DEFINIÇÃO DO PROBLEMA

PNAD-COVID-19 do IBGE O Head de Dados pediu para que você entrasse na base de dados do PNAD-COVID-19 do IBGE e organizasse esta base para análise, utilizando Banco de Dados em Nuvem e trazendo as seguintes características:

- a. Utilização de no máximo 20 questionamentos realizados na pesquisa;
- b. Utilizar 3 meses para construção da solução;
- c. Caracterização dos sintomas clínicos da população;
- d. Comportamento da população na época da COVID-19;
- e. Características econômicas da Sociedade;

Seu objetivo será trazer uma breve análise dessas informações, como foi a organização do banco, as perguntas selecionadas para a resposta do problema e quais seriam as principais ações que o hospital deverá tomar em caso de um novo surto de COVID-19.

Dica: Leiam com atenção a base de dados

_ SOBRE A PNAD COVID19

O que é?

Objetiva estimar o número de pessoas com sintomas referidos associados à síndrome gripal e monitorar os impactos da pandemia da COVID-19 no mercado de trabalho brasileiro.

A coleta da Pesquisa Nacional por Amostra de Domicílios - PNAD COVID19 teve início em 4 de maio de 2020, com entrevistas realizadas por telefone em, aproximadamente, 48 mil domicílios por semana, totalizando cerca de 193 mil domicílios por mês, em todo o Território Nacional. A amostra é fixa, ou seja, os domicílios entrevistados no primeiro mês de coleta de dados permanecerão na amostra nos meses subsequentes, até o fim da pesquisa.

— CORONAVÍRUS

A Covid-19 é uma infecção respiratória aguda causada pelo coronavírus SARS-CoV-2, potencialmente grave, de elevada transmissibilidade e de distribuição global. O SARS-CoV-2 é um betacoronavírus descoberto em amostras de lavado broncoalveolar obtidas de pacientes com pneumonia de causa desconhecida na cidade de Wuhan, província de Hubei, China, em dezembro de 2019. Pertence ao subgênero Sarbecovírus da família Coronaviridae e é o sétimo coronavírus conhecido a infectar seres humanos.

De acordo com as evidências mais atuais, o SARS-CoV-2, da mesma forma que outros vírus respiratórios, é transmitido principalmente por três modos: contato, gotículas ou por aerossol.

A principal medida de prevenção contra formas graves da covid-19 é a vacina.

A campanha de vacinação contra a covid-19 foi iniciada em janeiro de 2021!



SOBRE O PROJETO

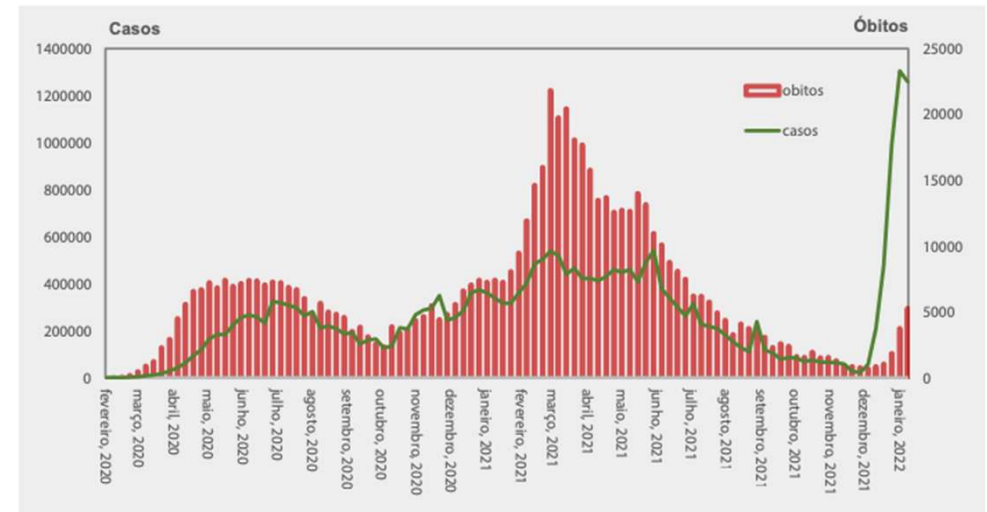
02

_ SOBRE O PROJETO

Banco de dados - O banco de dados foi disponibilizado no site para pesquisa no PNAD-COVID-19 do IBGE.

Foi inserido 3 tabelas de dados referentes aos meses de julho, agosto e setembro de 2020.

Utilizado estes meses, devido ao maior período de alta e posterior queda no ano de 2020. Todos estes dados foram incluídos e estruturados no Google Cloud.



Estrutura - O banco de dados SQL 'pos_tech3_caio' criado foi composto por 3 tabelas PNAD 6, PNAD 7 e PNAD 8 com uma tabela que armazena todos os dados da pesquisa.

A pesquisa é extensa contendo 1.151.956 linhas e 20 colunas, já nosso campo de pesquisa se limitou a trabalhar com as variáveis mostradas no próximo slide: Dicionário.

DICIONÁRIO

Código das variáveis escolhidas para análise e sua descrição

Código da variável	nº	Descrição
Ano		Ano de referência
UF		Unidade da Federação
V1013		Mês da pesquisa
V1022		Situação do domicílio
A002	A2	Idade do morador
A003	A3	Sexo
A004	A4	Cor ou raça
A005	A5	Escolaridade
B002	B2	Por causa disso, foi a algum estabelecimento de saúde?
B0031	B3	Providência tomada para recuperar dos sintomas foi ficar em casa
B0032	B3	Providência tomada para recuperar dos sintomas foi ligar para algum profissional de saúde
B0042	B4	Local que buscou atendimento foi pronto socorro do SUS/UPA
B0043	B4	Local que buscou atendimento foi hospital do SUS
B005	B5	Ao procurar o hospital, teve que ficar internado por um dia ou mais
B007	B7	Tem algum plano de saúde médico, seja particular, de empresa ou de órgão público
C003	C3	Qual o principal motivo deste afastamento temporário?
C013	C13	Na semana passada, o(a) Sr(a) estava em trabalho remoto (home office ou teletrabalho)?
C016	C16	Qual o principal motivo de não ter procurado trabalho na semana passada?
D0031	D1	Rendimentos de Programa Bolsa Família
D0061	D1	Seguro desemprego
F001	F1	Este domicílio é:

Respostas: As opções de respostas e regras são mostradas no link: [PNAD COVID19 - Questionário](#)

No trabalho as opções apontadas como 'Não Aplicável' se dá ao fato de serem ignoradas no momento da pesquisa, seja devido as regras das pesquisas ou usuário.

_GOOGLE CLOUD

Carregado dados para o Google Cloud conforme podemos verificar no destaque amarelo.

Alterado por código as variáveis da tabela para descrição conforme imagem abaixo.

The screenshot displays the Google Cloud console interface. On the left, the 'Explorer' sidebar shows a project named 'PostecProjetoTres'. Under 'Consultas salvas (8)', there are several saved queries, including 'pos_tech3_caio' through 'pos_tech3_caio7'. A yellow box highlights the 'Conexões externas' (External connections) section, which lists several databases: 'PNAD', 'PNAD060708_clusters', 'PNAD06_clusters', 'PNAD07_clusters', 'PNAD08_clusters', 'PNAD06', 'PNAD07', and 'PNAD08'. The main area shows a SQL query named 'pos_tech3_caio' being executed. The query is a complex CASE statement that filters data from the 'PNAD.PNAD08' table based on various conditions. Below the query, the 'Resultados da consulta' (Query results) section displays a table with 7 rows of data. The table has columns for 'Linha' (Line), 'Estado' (State), 'Mes' (Month), 'Zona' (Zone), 'Idade' (Age), 'Sexo' (Sex), 'Raça' (Race), 'Escolaridade' (Education), and 'Estabelecimento saúde' (Health establishment). The results show data for the state of Rondônia in June, categorized by urban area, age, sex, race, and education level.

Google Cloud

PosTecProjetoTres

Explorer

+ ADICIONAR

Google Cloud

PosTecProjetoTres

Pesquise (/) recursos, documentos, produtos e muito mais

Pesquisa

pos_tech3_caio

EXECUTAR

SALVAR CONSULTA

COMPARTILHAR

PROGRAMAÇÃO

MAIS

Consulta concluída.

```
466 CASE D0061
467 WHEN 1 THEN 'sim'
468 WHEN 2 THEN 'não'
469 ELSE CAST(D0061 AS STRING)
470 END as 'Seguro desemprego',
471 CASE
472 WHEN F001 = 1 THEN 'Próprio - já pago'
473 WHEN F001 = 2 THEN 'Próprio - ainda pagando'
474 WHEN F001 = 3 THEN 'Alugado'
475 WHEN F001 = 4 THEN 'Cedido por empregador'
476 WHEN F001 = 5 THEN 'Cedido por familiar'
477 WHEN F001 = 6 THEN 'Cedido de outra forma'
478 WHEN F001 = 7 THEN 'Outra condição'
479 ELSE 'Não aplicável'
480 END as 'Domicílio'
481 FROM
482 'postecprojeto08.PNAD.PNAD08'
```

Resultados da consulta

SALVAR RESULTADOS

EXPLORAR DADOS

INFORMAÇÕES DO JOB	RESULTADOS	GRÁFICO	JSON	DETALHES DA EXECUÇÃO	GRÁFICO DE EXECUÇÃO			
Linha	Estado	Mes	Zona	Idade	Sexo	Raça	Escolaridade	Estabelecimento saúde
1	Rondônia	junho	urbana	13	homem	Parda	Fundamental incompleto	Não aplicável
2	Rondônia	junho	urbana	10	homem	Parda	Fundamental incompleto	Não aplicável
3	Rondônia	junho	urbana	4	mulher	Parda	Sem instrução	Não aplicável
4	Rondônia	junho	urbana	0	mulher	Branca	Sem instrução	Não aplicável
5	Rondônia	junho	urbana	6	mulher	Branca	Fundamental incompleto	Não aplicável
6	Rondônia	junho	urbana	9	mulher	Parda	Fundamental incompleto	Não aplicável
7	Rondônia	junho	urbana	1	homem	Parda	Sem instrução	Não aplicável

Resultados por página: 50 1 - 50 de 1151956

Histórico de jobs

ATUALIZAR

_GOOGLE COLAB

Após código SQL pronto e tabela OK no Google Cloud, foi transferido para o Google Colab, por meio de uma conexão python. Obs.: essa transferência é automática.

```
[ ] ##Setup
# @title Setup
from google.colab import auth
from google.cloud import bigquery
from google.colab import data_table

project = 'postecprojetotres' # Project ID inserted based on the query results selected to explore
location = 'southamerica-east1' # Location inserted based on the query results selected to explore
client = bigquery.Client(project=project, location=location)
data_table.enable_dataframe_formatter()
auth.authenticate_user()

# Reference SQL syntax from the original job

Use the jobs.query method to return the SQL syntax from the job. This can be copied from the output cell below to edit the query now or in the future. Alternatively, you can use this link back to BigQuery to edit the query within the BigQuery user interface.

[ ] # Running this code will display the query used to generate your previous job

job = client.get_job('bquxjob_2ea294bb_18dbdb03e20') # Job ID inserted based on the query results selected to explore
print(job.query)

SELECT
CASE UF
WHEN 11 THEN 'Rondônia'
WHEN 12 THEN 'Acre'
```

_GOOGLE COLAB

Imagem da tabela de dados no Google Colab e a quantidade de dados.

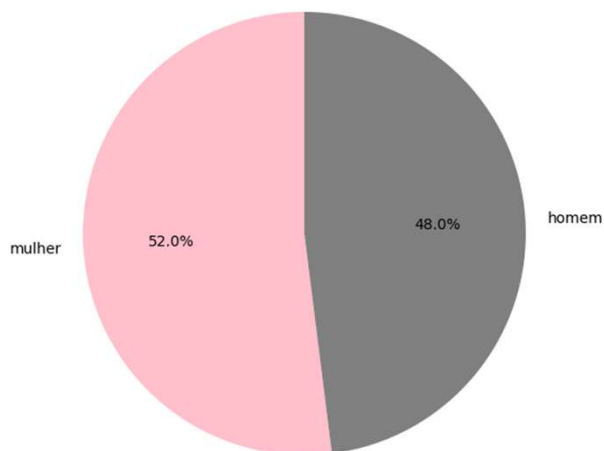
	Estado	Mes	Zona	Idade	Sexo	Raça	Escolaridade	Estabelecimento saúde	Providência - casa	Providência - profissional	Atendimento PS	Atendimento hospital	Internação	Plano de saúde	Afastamento	Trabalho remoto	PQ não procurou trabalho	Bolsa Família	Seguro desemprego	Domicílio
0	Rondônia	junho	urbana	3	homem	Parda	Sem instrução	Não	Sim	Não	Não aplicável	Não aplicável	Não aplicável	Não	Não aplicável	Não aplicável	Não aplicável	não	não	Próprio - ainda pagando
1	Rondônia	junho	urbana	7	homem	Branca	Fundamental incompleto	Não aplicável	Não aplicável	Não aplicável	Não aplicável	Não aplicável	Não aplicável	Sim	Não aplicável	Não aplicável	Não aplicável	não	não	Alugado
2	Rondônia	junho	urbana	8	homem	Parda	Fundamental incompleto	Não	Sim	Não	Não aplicável	Não aplicável	Não aplicável	Não	Não aplicável	Não aplicável	Não aplicável	não	não	Cedido por familiar
3	Rondônia	junho	urbana	4	mulher	Parda	Sem instrução	Não	Sim	Não	Não aplicável	Não aplicável	Não aplicável	Não	Não aplicável	Não aplicável	Não aplicável	sim	não	Alugado
4	Rondônia	junho	urbana	2	mulher	Parda	Sem instrução	Não aplicável	Não aplicável	Não aplicável	Não aplicável	Não aplicável	Não aplicável	Não	Não aplicável	Não aplicável	Não aplicável	sim	não	Cedido por familiar
...
1151951	São Paulo	agosto	urbana	29	homem	Parda	Fundamental completa	Não aplicável	Não aplicável	Não aplicável	Não aplicável	Não aplicável	Não aplicável	Não	Estava em quarentena, isolamento, distanciamen...	Não aplicável	Não aplicável	não	não	Alugado
1151952	Paraná	agosto	urbana	39	mulher	Branca	Médio completo	Não aplicável	Não aplicável	Não aplicável	Não aplicável	Não aplicável	Não aplicável	Sim	Estava em quarentena, isolamento, distanciamen...	Não aplicável	Não aplicável	não	sim	Próprio - já pago
1151953	Paraná	agosto	rural	51	mulher	Parda	Fundamental incompleto	Não aplicável	Não aplicável	Não aplicável	Não aplicável	Não aplicável	Não aplicável	Sim	Licença remunerada por motivo de saúde ou acid...	Não aplicável	Por problemas de saúde ou gravidez	não	não	Cedido por empregador
1151954	Santa Catarina	agosto	urbana	47	mulher	Branca	Médio completo	Não aplicável	Não aplicável	Não aplicável	Não aplicável	Não aplicável	Não aplicável	Não	Estava em quarentena, isolamento, distanciamen...	Não aplicável	Não quer trabalhar ou é aposentado	não	não	Próprio - já pago
1151955	Distrito Federal	agosto	urbana	20	homem	Parda	Superior incompleto	Não aplicável	Não aplicável	Não aplicável	Não aplicável	Não aplicável	Não aplicável	Não	Estava em quarentena, isolamento, distanciamen...	Não aplicável	Devido à pandemia (isolamento, quarentena ou d...	não	não	Próprio - já pago

1151956 rows × 20 columns

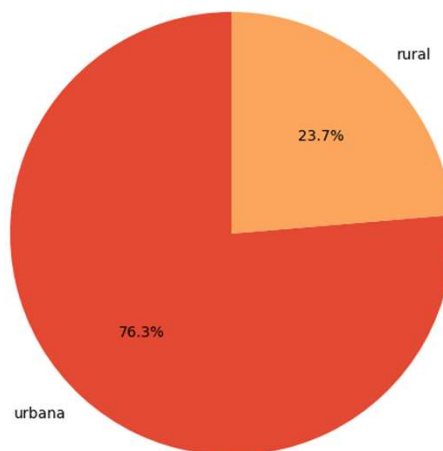
Criado todos os gráficos a seguir em python no Google Colab.

_ CARACTERÍSTICAS GERAIS

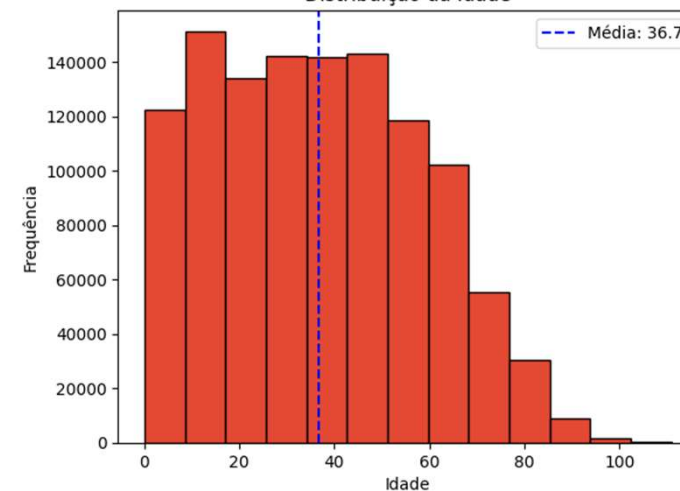
Distribuição do Sexo



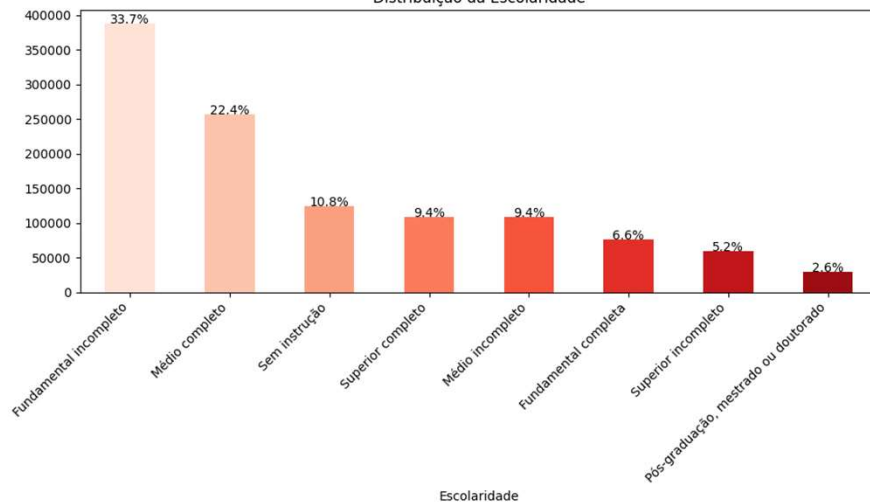
Situação do domicílio



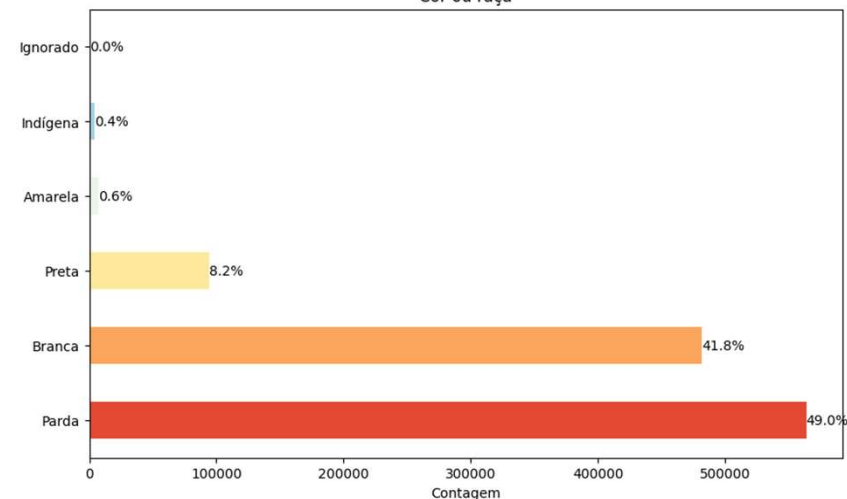
Distribuição da Idade



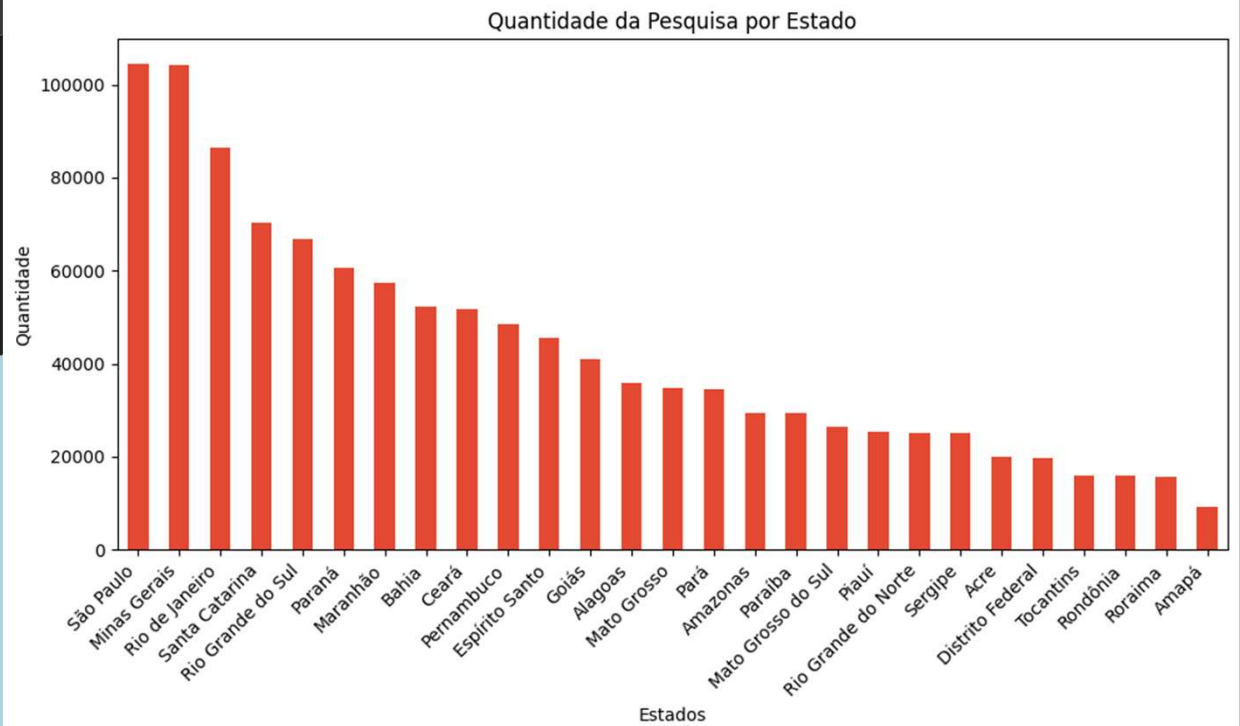
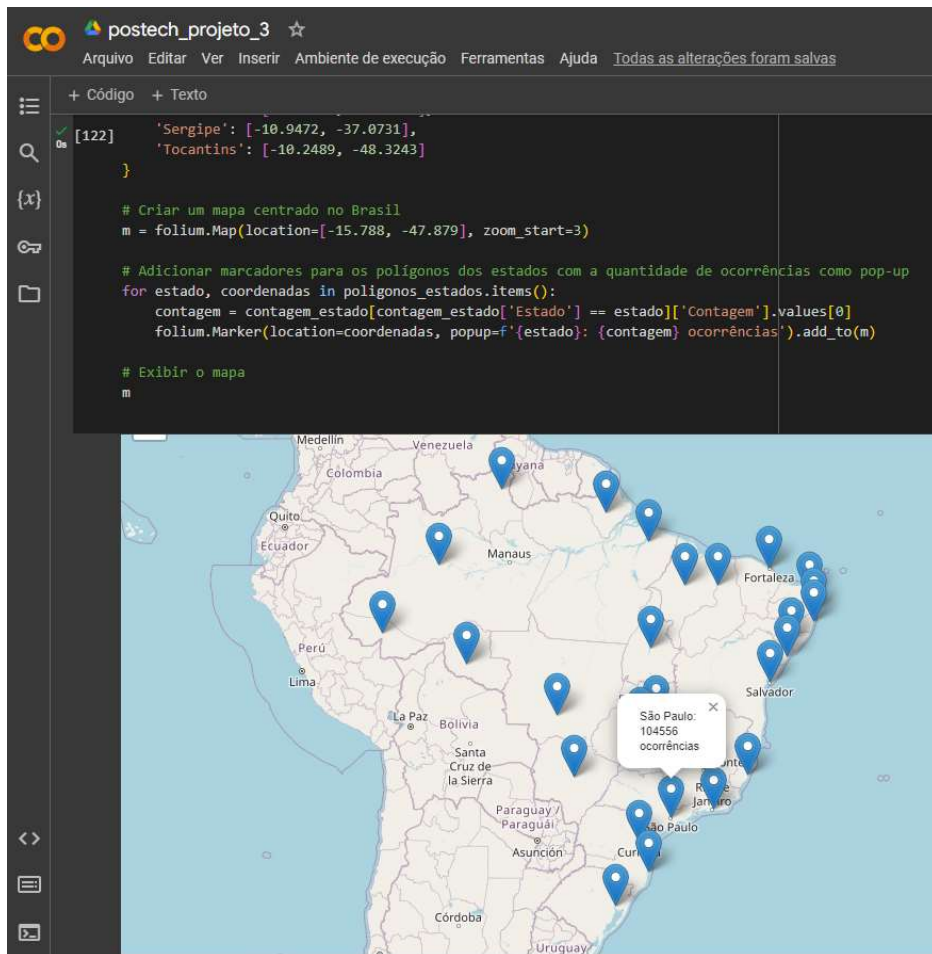
Distribuição da Escolaridade



Cor ou raça

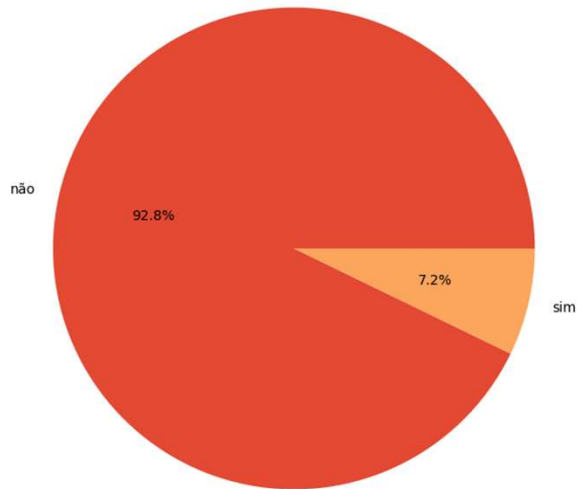


_ CARACTERÍSTICAS DE LOCALIZAÇÃO

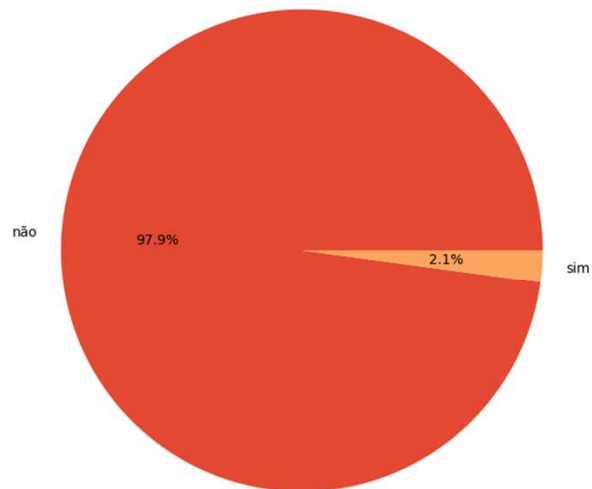


CARACTERÍSTICAS ECONÔMICAS

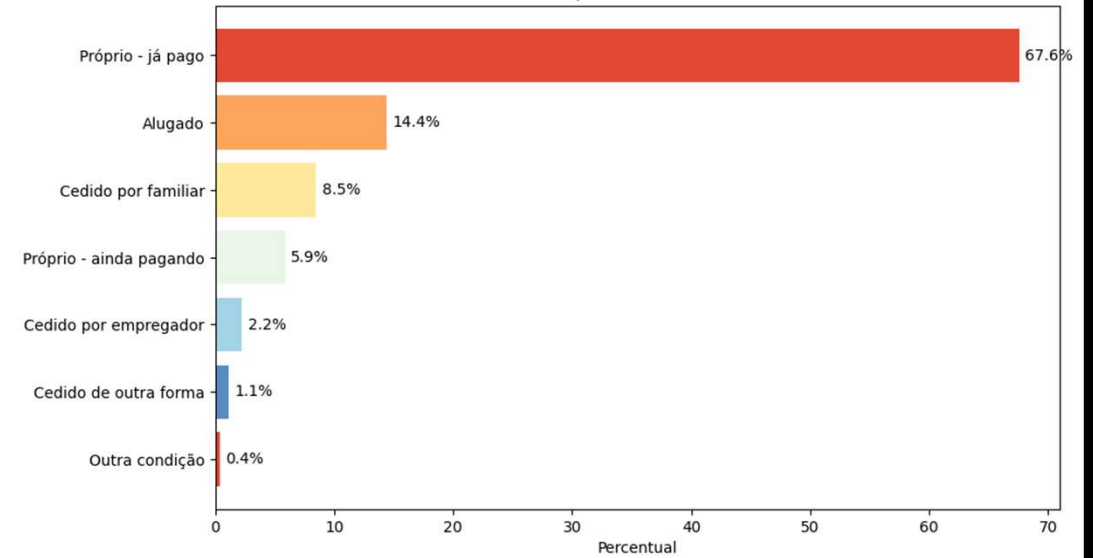
Bolsa Família



Distribuição de Seguro Desemprego

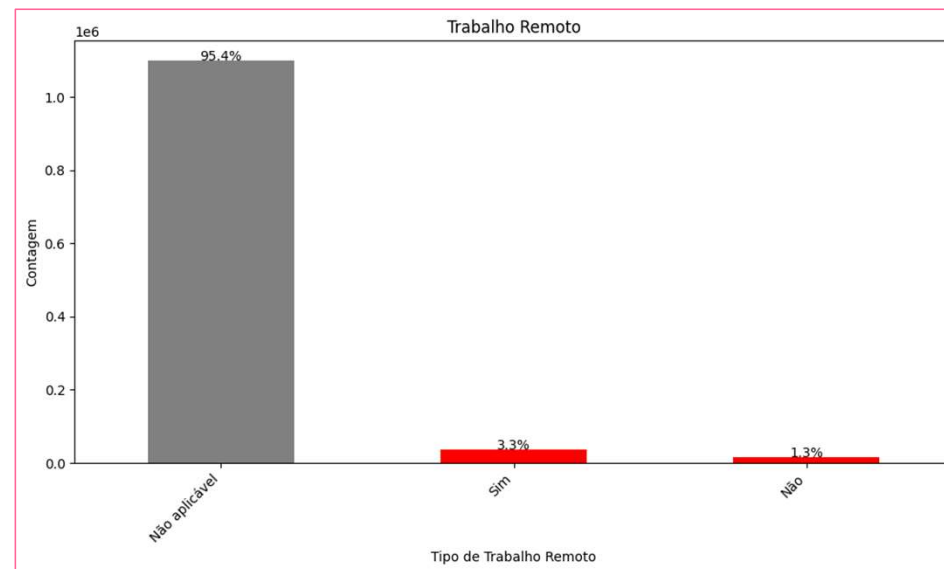
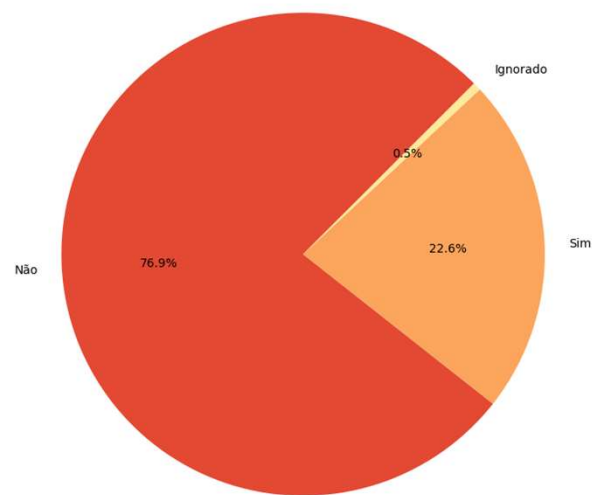


Tipo de Domicílio

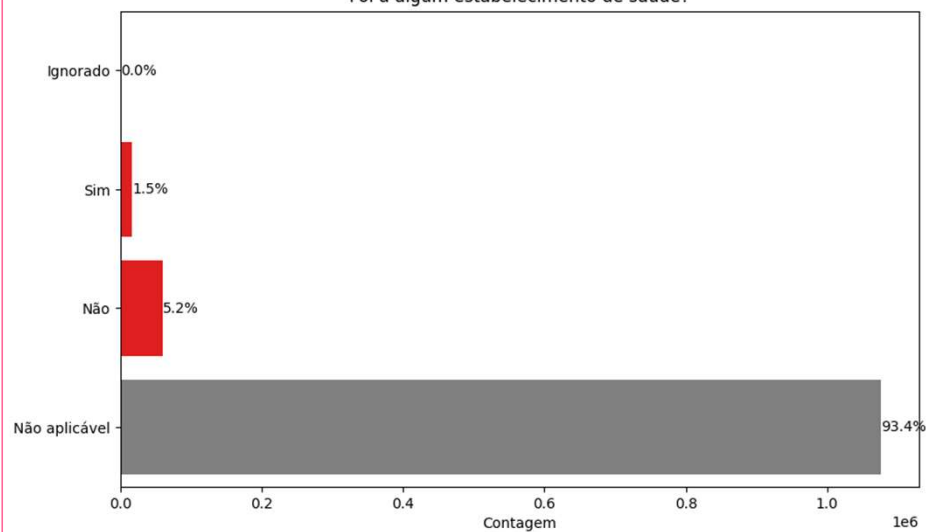


_ CARACTERÍSTICAS GERAIS

Tem algum plano de saúde médico, seja particular, de empresa ou de órgão público



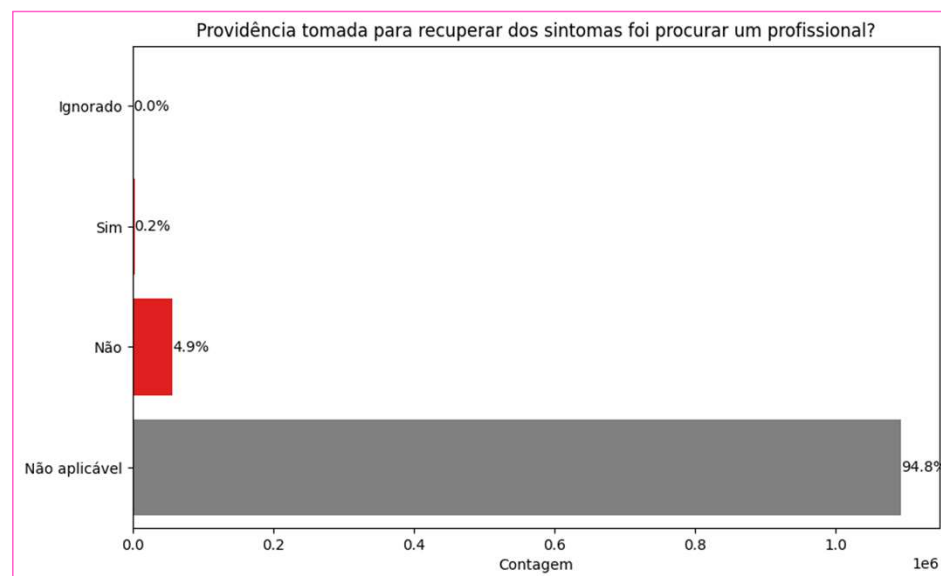
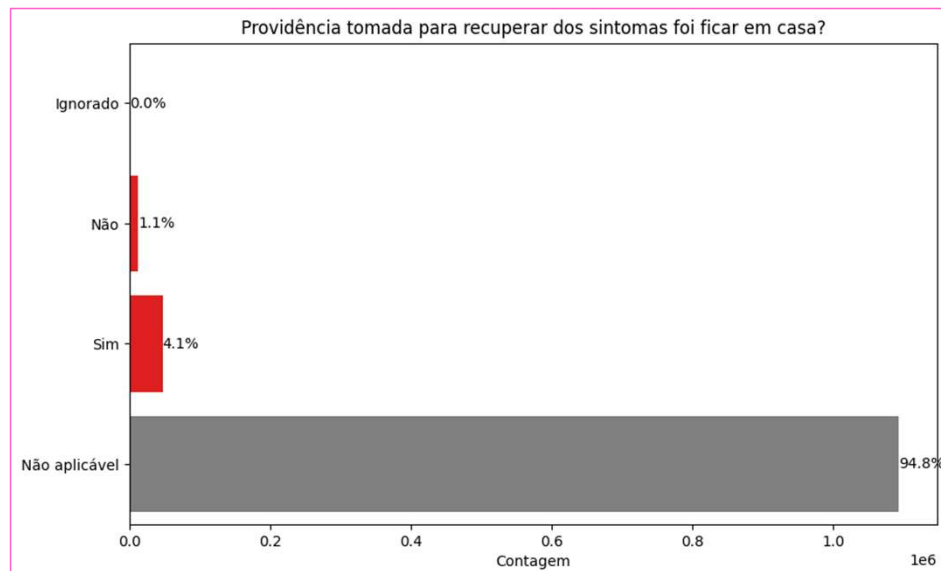
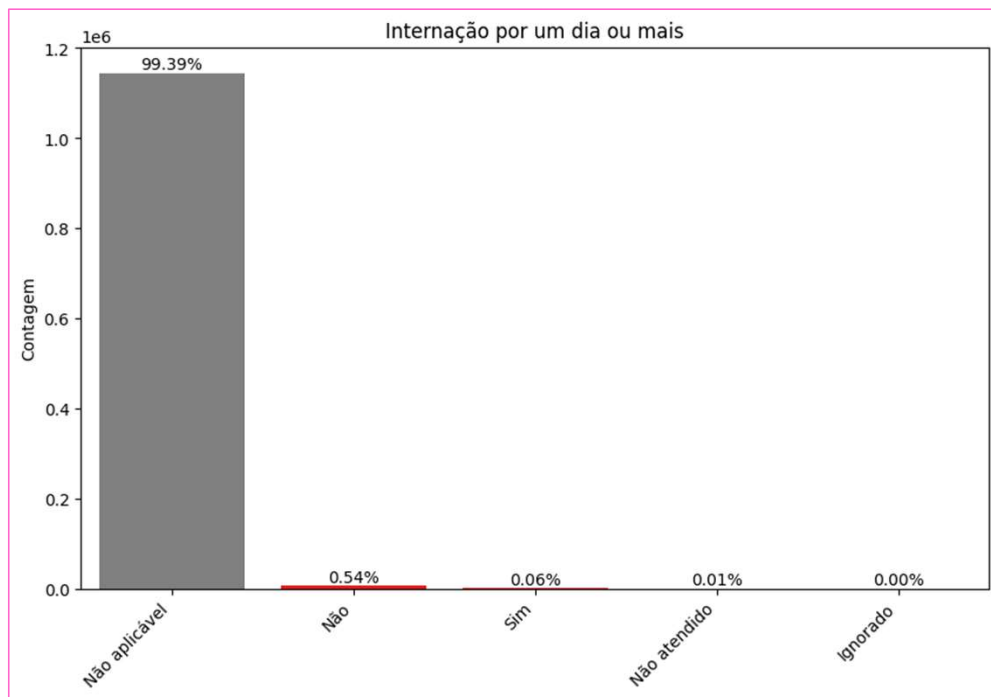
Foi a algum estabelecimento de saúde?



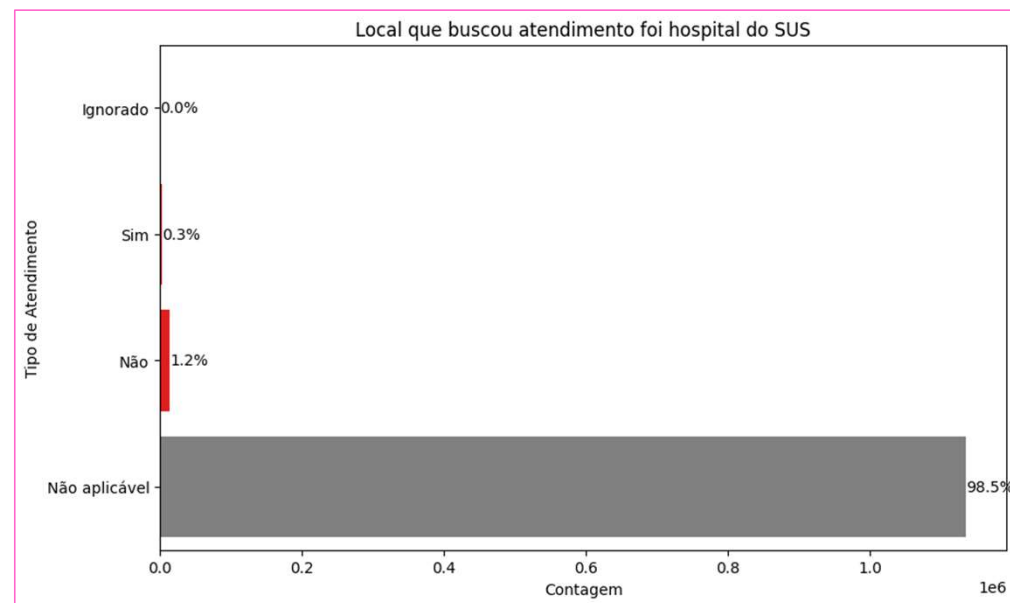
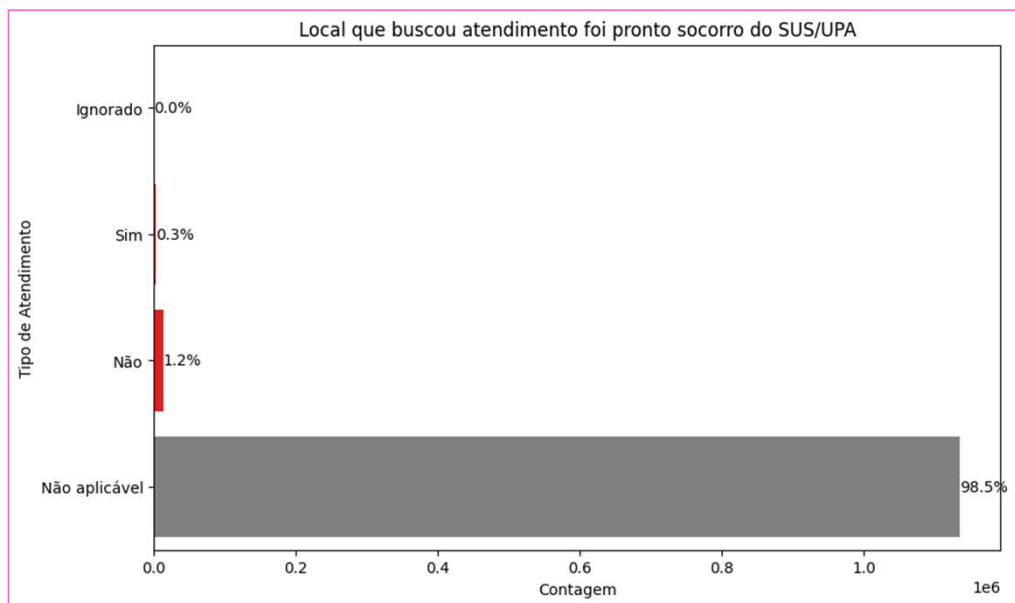
Qual o principal motivo deste afastamento temporário?



_ CARACTERÍSTICAS CLÍNICAS



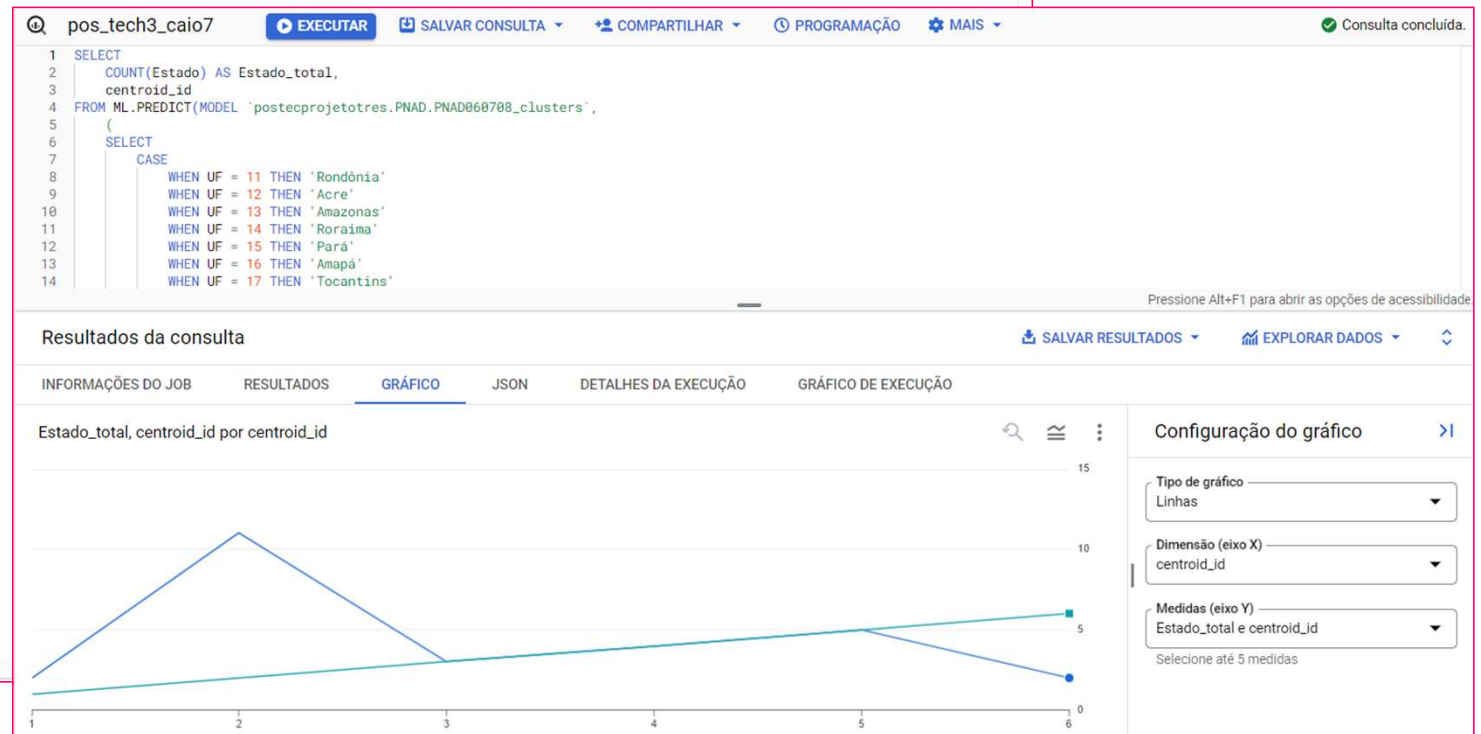
_ CARACTERÍSTICAS CLÍNICAS



_ CLUSTERING-ML

Mostrado agrupamento dos estados e o total de internação mostrando a similaridade dos dados entre os 3 meses.

```
1 SELECT
2   COUNT(Estado) AS Estado_total,
3   centroid_id
4 FROM ML.PREDICT(MODEL 'postecprojetotres.PNAD.PNAD060708_clusters',
5 (
6   SELECT
7     CASE
8       WHEN UF = 11 THEN 'Rondônia'
9       WHEN UF = 12 THEN 'Acre'
10      WHEN UF = 13 THEN 'Amazonas'
11      WHEN UF = 14 THEN 'Roraima'
12      WHEN UF = 15 THEN 'Pará'
13      WHEN UF = 16 THEN 'Amapá'
14      WHEN UF = 17 THEN 'Tocantins'
15      WHEN UF = 21 THEN 'Maranhão'
16      WHEN UF = 22 THEN 'Piauí'
17      WHEN UF = 23 THEN 'Ceará'
18      WHEN UF = 24 THEN 'Rio Grande do Norte'
19      WHEN UF = 25 THEN 'Paraíba'
20      WHEN UF = 26 THEN 'Pernambuco'
21      WHEN UF = 27 THEN 'Alagoas'
22      WHEN UF = 28 THEN 'Sergipe'
23      WHEN UF = 29 THEN 'Bahia'
24      WHEN UF = 31 THEN 'Minas Gerais'
25      WHEN UF = 32 THEN 'Espírito Santo'
26      WHEN UF = 33 THEN 'Rio de Janeiro'
27      WHEN UF = 35 THEN 'São Paulo'
28      WHEN UF = 41 THEN 'Paraná'
29      WHEN UF = 42 THEN 'Santa Catarina'
30      WHEN UF = 43 THEN 'Rio Grande do Sul'
31      WHEN UF = 50 THEN 'Mato Grosso do Sul'
32      WHEN UF = 51 THEN 'Mato Grosso'
33      WHEN UF = 52 THEN 'Goiás'
34      WHEN UF = 53 THEN 'Distrito Federal'
35      ELSE CAST(UF AS STRING)
36    END AS Estado,
37    SUM(CASE WHEN B005 = 1 THEN 1 ELSE 0 END) AS total_internacao
38 FROM (
39   SELECT UF, B005 FROM postecprojetotres.PNAD.PNAD06
40   UNION ALL
41   SELECT UF, B005 FROM postecprojetotres.PNAD.PNAD07
42   UNION ALL
43   SELECT UF, B005 FROM postecprojetotres.PNAD.PNAD08
44 ) AS data_merged
45 GROUP BY Estado
46 )
47 )
48 GROUP BY centroid_id
49 ORDER BY centroid_id;
50
```





CONSIDERAÇÕES FINAIS

04

_ CONSIDERAÇÕES FINAIS

Após uma análise detalhada do comportamento da população na época de pandemia da COVID-19, recomendamos os seguintes indicadores para uma melhor gestão de um grande hospital para se planejar de forma mais adequada em novos surtos:

Alguns dos atributos de maior peso para monitoração do Público em atendimento no hospital:

idade, sexo, etnia, status socioeconômico, sintomas e comorbidades

% atendimento diário no ambiente/ moradia urbana

% público com menor grau de escolaridade sendo sem instrução, médio completo ou fundamental incompleto

% público sem plano de saúde

_ CONSIDERAÇÕES FINAIS

Além destes indicadores combinados com outros dados para atuar de forma preventiva:

1. Uso de sistema de saúde. Observar as tendências como visitas no pronto socorro e internações devido ao vírus pode apoiar no entendimento da gravidade da doença;
- 2- Dados de mobilidade para entender padrões de deslocamento da população e permitir que o hospital aloque recursos;
- 3- Monitoramento de mídias sociais na adesão a medidas preventivas para apoiar em ideias de conscientização e ajustes na comunicação;
- 4- Monitoramento de indicadores de saúde pública como taxa de positividade de testes.

De maneira simplificada, um modelo para apoiar o hospital deverá se concentrar uma grande quantidade de dados que, cruzados, indicam padrões e tendências futuras sobre os mais diferentes problemas, como a disseminação de doenças.

O modelo proposto pode contribuir para a diminuição da subjetividade na tomada de decisão, para que o profissional de saúde possa tomar suas decisões baseadas em dados de forma mais assertiva.



REFERÊNCIAS

05

REFERÊNCIAS

- ❖ <https://covid19.ibge.gov.br/pnad-covid/>
- ❖ <https://www.fiocruzbrasil.fiocruz.br/covid-19-balanco-de-dois-anos-da-pandemia-aponta-vacinacao-como-prioridade/>
- ❖ https://www.ibge.gov.br/estatisticas/investigacoes-experimentais/estatisticas-experimentais/27946-divulgacao-semanal-pnadcovid1?t=o-que-e&utm_source=covid19&utm_medium=hotsite&utm_campaign=covid_19
- ❖ <https://www.gov.br/saude/pt-br/assuntos/coronavirus>
- ❖ https://www.ibge.gov.br/estatisticas/investigacoes-experimentais/estatisticas-experimentais/27946-divulgacao-semanal-pnadcovid1?t=downloads&utm_source=covid19&utm_medium=hotsite&utm_campaign=covid_19