

# Telecomunicações e conectividade móvel no Brasil em 2020 e 2021

Projeto de Conclusão de Curso – Engenharia de Dados BC17





# DESENVOLVEDORES



Gabriel  
Janjacomio



Jum  
Sahek  
i



Luma  
Guimarães



Victor  
Ribeir  
o



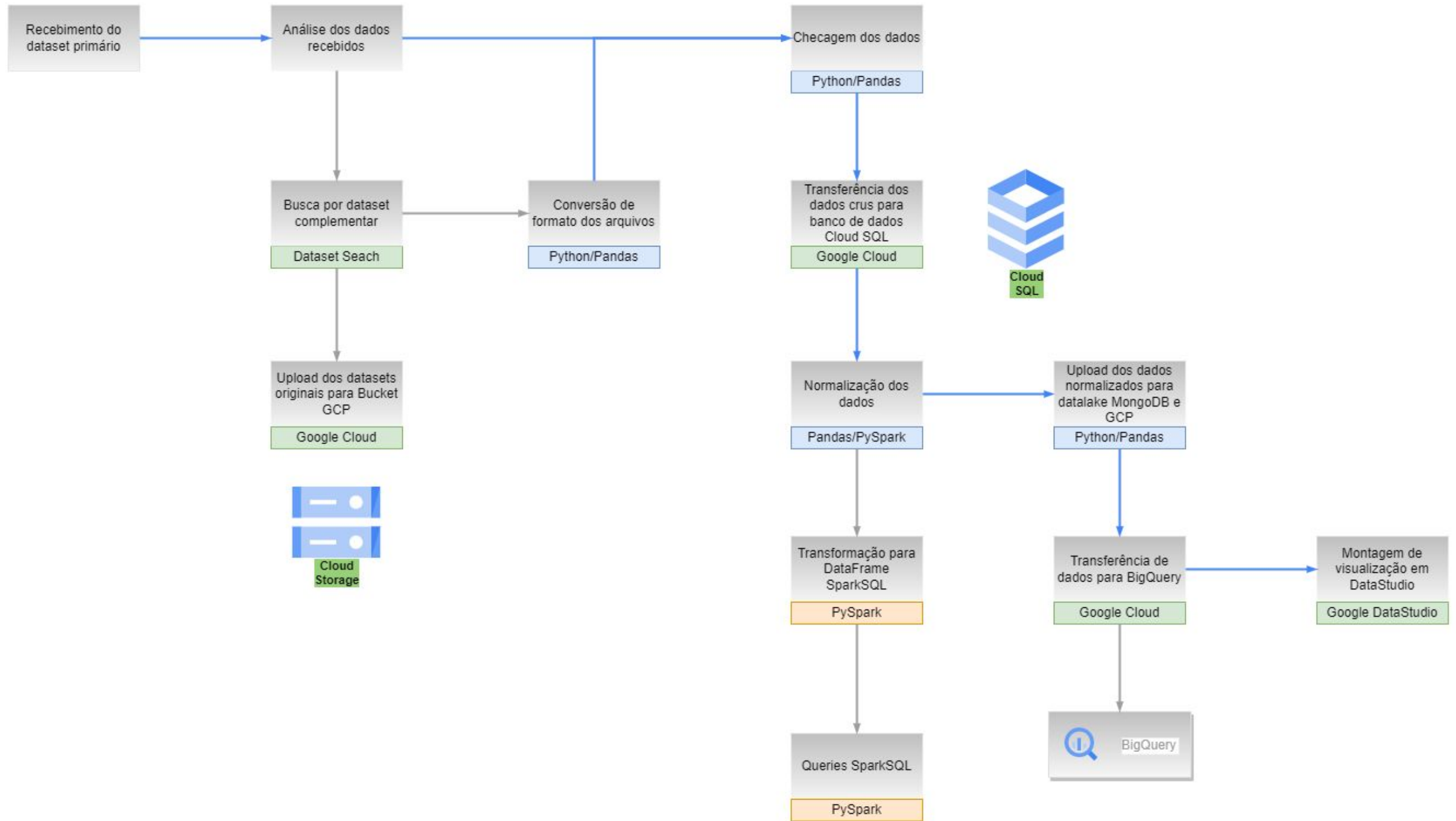
# SUMÁRIO

Esse projeto teve por objetivo exercitar e demonstra o aprendizado do conjunto de tecnologias e infraestruturas que foram passados ao longo do bootcamp.

Utilizando uma base fornecida pelos professores, fomos encarregados de analisar, somar, normalizar e tratar dados a fim de gerar insights e finalmente visualizações gráficas.

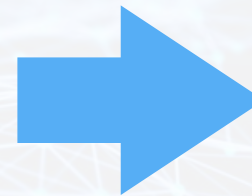
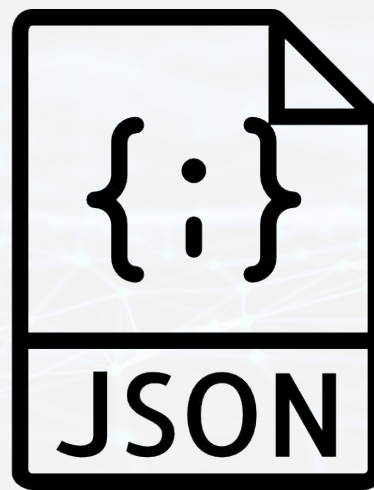
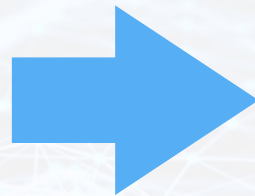
# TECNOLOGIAS UTILIZADAS







# TRANSFORMAÇÕES NOS FORMATOS DE ARQUIVO



# TRANSFORMAÇÕES NOS FORMATOS DE ARQUIVO

```
1 '''
2 PUXA ARQUIVO CSV DO BUCKET GCP CONVERTE PARA JSON E SALVA LOCAL
3 '''
4
5 listaNomes = ['Acessos_Telefonia_Movel_202107-202112.csv',
6               'Acessos_Telefonia_Movel_202107-202112_Colunas.csv', 'Acessos_Telefonia_Movel_202201-202106_Colunas.csv',
7               'Acessos_Telefonia_Movel_202201-202206.csv', 'Acessos_Telefonia_Movel_Pre_Pos_Total.csv',
8               'Acessos_Telefonia_Movel_Total.csv', 'Densidade_Telefonia_Movel.csv']
9
10 for nomeArquivo in listaNomes:
11
12     GCpath = "gs://soulcode-bc17-telecom/" + nomeArquivo
13
14     if nomeArquivo == 'Acessos_Telefonia_Movel_202107-202112.csv':
15         colsList = [0, 1, 3, 6, 7, 11, 12, 13, 14, 15, 16]
16         with pd.read_csv(GCpath, sep=';', usecols=colsList, chunksize=500000) as leitor:
17             chunkNO = 1
18             for chunk in leitor:
19                 nomeArquivo = listaNomes[0].split('.')
20                 nomeProcc = (str(nomeArquivo[0]) + '_chunk' + str(chunkNO) + '.json')
21                 localpath = r"/content/ArquivosOut/" + nomeProcc
22                 chunk.to_json(localpath)
23                 chunkNO = chunkNO + 1
24
25     else:
26         dfCSV = pd.read_csv(GCpath, sep=';')
27
28         nomeProcc = nomeArquivo.split('.')
29         nomeArquivo = nomeProcc[0] + '.json'
30         localpath = r"/content/ArquivosOut/" + nomeArquivo
31         dfCSV.to_json(localpath)
32
33
```

# INFRAESTRUTURA DE ARMAZENAMENTO EM NUVEM

Cloud Storage

Browser

Monitoring

Settings

Marketplace

Release Notes

<|

Bucket details

REFRESH

HELP ASSISTANT

LEARN

soulcode-bc17-telecom

Location

us (multiple regions in United States)

Storage class

Standard

Public access

Public to internet

Protection

None

OBJECTS

CONFIGURATION

PERMISSIONS

PROTECTION

LIFECYCLE

Buckets > soulcode-bc17-telecom

UPLOAD FILES

UPLOAD FOLDER

CREATE FOLDER

MANAGE HOLDS

DOWNLOAD

DELETE

Sort and filter

Filter

Filter objects and folders

?

Show deleted data

|||

<input type="checkbox"/>	Name ↑	Size	Type	Created ?	Storage class	Last modified	Publi	
<input type="checkbox"/>	.ipynb_checkpoints/	—	Folder	—	—	—	—	⋮
<input type="checkbox"/>	Acessos_Telefonia_Movel_202107-2021...	3 GB	text/csv	Jun 8, 2022, 10...	Standard	Jun 8, 2022, 10:...	⚠ F	⬇ ⋮
<input type="checkbox"/>	Acessos_Telefonia_Movel_202107-2021...	977.3 MB	application/json	Jun 9, 2022, 2:...	Standard	Jun 9, 2022, 2:4...	⚠ F	⬇ ⋮
<input type="checkbox"/>	Acessos_Telefonia_Movel_202107-2021...	987 MB	application/json	Jun 9, 2022, 2:...	Standard	Jun 9, 2022, 2:3...	⚠ F	⬇ ⋮
<input type="checkbox"/>	Acessos_Telefonia_Movel_202107-2021...	1 GB	application/json	Jun 9, 2022, 2:...	Standard	Jun 9, 2022, 2:3...	⚠ F	⬇ ⋮
<input type="checkbox"/>	Acessos_Telefonia_Movel_202107-2021...	1 GB	application/json	Jun 9, 2022, 2:...	Standard	Jun 9, 2022, 2:2...	⚠ F	⬇ ⋮



# INFRAESTRUTURA DE ARMAZENAMENTO EM NUVEM

The screenshot displays the ClusterO web interface, which is used for managing cloud storage infrastructure. The interface is divided into several sections:

- Header:** The ClusterO logo is on the left. On the right, it shows the **VERSION** as 5.0.9 and the **REGION** as AWS Sao Paulo (sa-east-1).
- Navigation Tabs:** A row of tabs includes Overview, Real Time, Metrics, Collections (which is the active tab), Search, Profiler, Performance Advisor, and Online Archive.
- Database Overview:** Below the tabs, it states "DATABASES: 1" and "COLLECTIONS: 6". There are buttons for "VISUALIZE YOUR DATA" and "REFRESH".
- Left Sidebar:** Contains a "+ Create Database" button, a "Search Namespaces" search bar, and a list of namespaces under the "soulcode" database. The namespaces listed are: Acesso\_Banda\_Larga..., Acesso\_Telefonia\_Movel, Acessos\_Municipio, Cobertura\_Municipios, Cobertura\_Setor, and Velocidade\_Contratada.
- Main Content Area:** Displays details for the "soulcode.Acesso\_Banda\_Larga\_Fixa" collection. It shows statistics: "STORAGE SIZE: 106.82MB", "TOTAL DOCUMENTS: 1323219", and "INDEXES TOTAL SIZE: 37.61MB". Below these are tabs for Find (active), Indexes, Schema Anti-Patterns (1), Aggregation, and Search Indexes (0). There is an "INSERT DOCUMENT" button.
- Filter and Document View:** A "FILTER" section shows a query: `{ field: 'value' }`. Below this, a document is displayed with the following fields: `_id: ObjectId("62a7b2058256b4c1be4a29b9")`, `Ano: 2022`, `Mes: 3`, `Grupo_Economico: "OI"`, `Empresa: "OI"`, and `CNDI: 76535764000143`. There are "Apply" and "Reset" buttons for the filter.
- Bottom Right:** A circular chat icon with two people silhouettes.

# CRIAÇÃO DOS TRIGGERS EM POSTGRESQL

```
CREATE OR REPLACE FUNCTION public.reg_acessos()
    RETURNS trigger
    LANGUAGE 'plpgsql'
    COST 100
    VOLATILE NOT LEAKPROOF
AS $BODY$
BEGIN
    --Trigger operation inserção
    IF (TG_OP = 'INSERT') THEN
        INSERT INTO Log_changes (log_id, data_criacao, operacao)
        VALUES (new.acessos, current_timestamp, 'Operação de inserção.
        A linha: ' || NEW.acessos || ' foi inserida');
        RETURN NEW;
    -- Trigger operation atualização
    ELIF (TG_OP = 'UPDATE') THEN
        INSERT INTO Log_changes (log_id, data_criacao, operacao)
        VALUES (NEW.acessos, current_timestamp, 'Operação de atualização.
        A linha: ' || NEW.acessos || ' foi atualizada'
        || OLD || ' com ' || NEW.* || '.');
        RETURN NEW;
    -- Trigger operation exclusão
    ELIF (TG_OP = 'DELETE') THEN
        INSERT INTO Log_changes (log_id, data_criacao, operacao)
        VALUES (OLD.acessos, current_timestamp, 'Operação de exclusão. A linha: '
        || OLD.acessos || ' foi excluída');
        RETURN OLD;
    END IF;
    RETURN NULL;
END;
$BODY$;

ALTER FUNCTION public.reg_acessos()
    OWNER TO postgres;
```

```
CREATE TRIGGER trigger_acessos
    AFTER INSERT OR DELETE OR UPDATE
    ON public."Acessos_Telefonia_Movel_Pre_Pos_Total"
    FOR EACH ROW
    EXECUTE FUNCTION public.reg_acessos();
```



# CONECTIVIDADE VIA PYTHON

## Google Cloud Platform

```
1  # Conectando ao drive para acessar a chave
2  drive.mount('/content/drive')
3
4  # Utilizando a chave para obter as informações de acesso
5  serviceAccount = '/content/drive/MyDrive/SoulCode/Datasets/macro-mercury-349020-d9ed9a670580.json'
6  os.environ['GOOGLE_APPLICATION_CREDENTIALS'] = serviceAccount
7
8  # Conectando ao bucket
9  client_bucket = storage.Client()
10 bucket = client_bucket.get_bucket('soulcode-bc17-telecom')
```

# CONECTIVIDADE VIA PYTHON

## MongoDB

```
1 # Credenciais de acesso
2 path = 'mongodb://soulcode:a1b2c3@cluster0-shard-00-00.ap3gr.mongodb.net/test?replicaSet=atlas-fg5u6j-shard-0&ssl=true&authSource=admin'
3 # Configurações do banco
4 client = MongoClient(path)
5 db = client['soulcode']
```



# ANÁLISE PRELIMINAR DOS DATASETS

```
1  # Quantidade de linhas
2  df_blf2021.shape
3  # Tipo de dado de cada coluna
4  df_blf2021.dtypes
5  # Valores únicos por coluna
6  df_blf2021.nunique()
7  # Valores vazios
8  df_blf2021.isna().sum()
9  # Determinando o número de duplicatas
10 df_blf2021.duplicated().sum()
11 pd.unique(df_blf2021['Mês'])
```

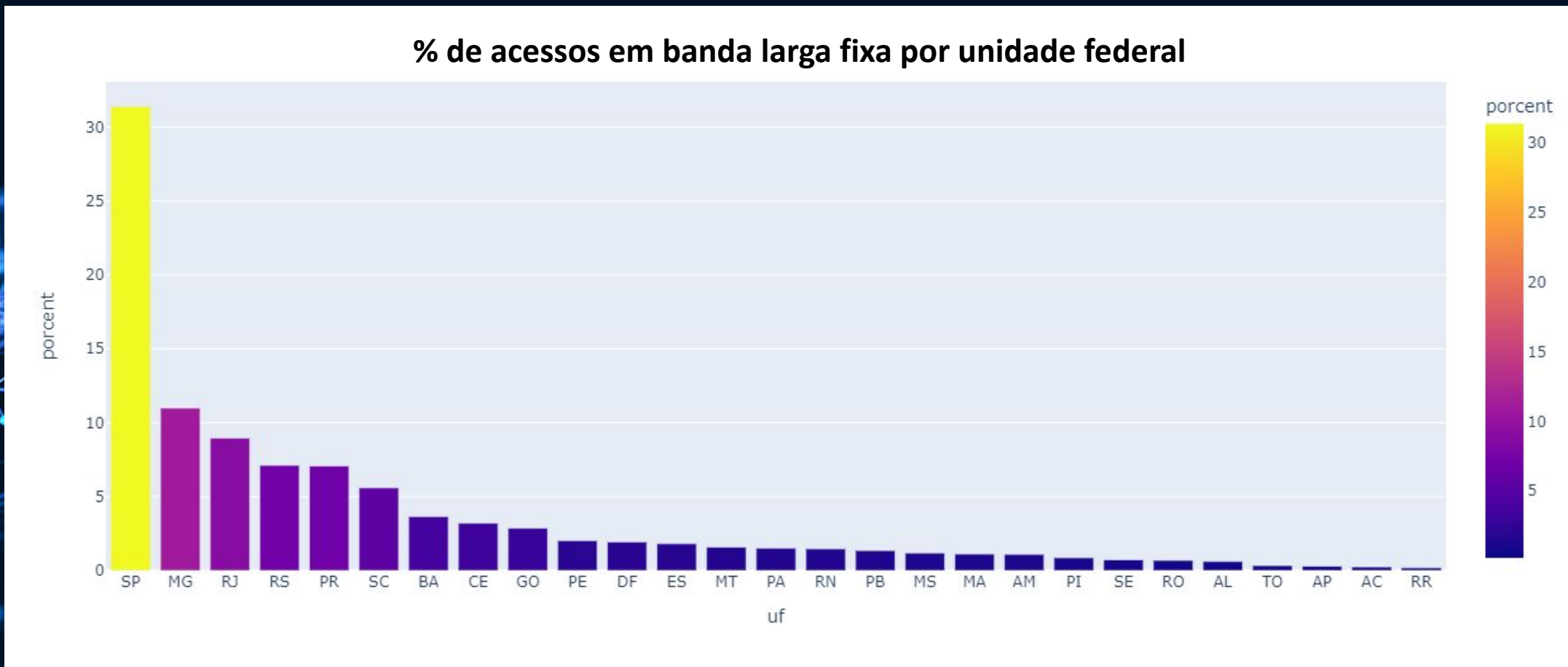
```
1  # Quantidade de linhas
2  df_blf2022.shape
3  # Tipo de dado de cada coluna
4  df_blf2022.dtypes
5  # Valores únicos por coluna
6  df_blf2022.nunique()
7  # Valores vazios
8  df_blf2022.isna().sum()
9  # Determinando o número de duplicatas
10 df_blf2022.duplicated().sum()
```

# ANÁLISE COM PANDAS

```
1 df2_extrai=pd.read_csv('gs://telecomunicacao2/Velocidade_Contratada_SCM.csv',sep=";")
2 # Copy
3 df2_pd = df2_extrai.copy()
4
5 df2_pd.head(2)
6 # Unique
7 pd.unique(df2_pd['Mês'])
8 # Info
9 df2_pd.info()
10 # Duplicated
11 df2_pd.duplicated().sum()
12 df2_pd=df2_pd[df2_pd.duplicated()]
13
14 # Drop duplicates
15 df2_pd = df2_pd.drop_duplicates()
16 pd.isnull(df2_pd).sum()
17 # Drop
18 df2_pd=df2_pd.drop(columns=['municipio_uf'])
19
```



# VISUALIZAÇÃO COM PANDAS



# TRANSFORMAÇÃO PARA SPARK

```
1 df_spark_blf = spark.createDataFrame(df_blf)
2 df_spark_bfl = (spark.read.format("parquet")
3                 .option("header", True)
4                 .option("delimiter", ",")
5                 .load("https://storage.googleapis.com/soulcode-bc17-telecom/Tratados/Acessos_Banda_Larga_Fixa.parquet"))
6 )
7 df_spark_bfl=spark.read.parquet("https://storage.googleapis.com/soulcode-bc17-telecom/Tratados/Acessos_Banda_Larga_Fixa.parquet")
```



# CHECAGEM DOS DADOS (PYSPARK)

```
1  # Número de linhas e colunas
2  print((df_spark_cob.count(), len(df_spark_cob.columns)))
3  # Tipo de dados
4  df_spark_cob.printSchema()
5  # Contagem de valores únicos por coluna
6  df_spark_cob.select([F.countDistinct(c).alias(c) for c in df_spark_cob.columns]).show()
7  # Valores nulos
8  df_spark_cob.select([F.count(F.when(F.col(c).contains('None') | \
9                        F.col(c).contains('NULL') | \
10                       (F.col(c) == '' ) | \
11                       F.col(c).isNull() | \
12                       F.isnan(c), c
13                       )).alias(c)
14                        for c in df_spark_cob.columns]).show()
15  # Valores duplicados
16  total = (df_spark_cob.count()) - (df_spark_cob.dropDuplicates().count())
17  print(total)
```

```
1  # Tipo de dado de cada coluna
2  df_blf2021.dtypes
```

Ano	int64
Mês	int64
Grupo Econômico	object
Empresa	object
CNPJ	int64
Porte da Prestadora	object
UF	object
Município	object
Código IBGE Município	int64
Faixa de Velocidade	object
Tecnologia	object
Meio de Acesso	object
Tipo de Pessoa	object
Acessos	int64
dtype:	object

# TRANSFORMAÇÃO DE DADOS (PYSPARK)

```
1 # Trocando o tipo dos dados de float para inteiro
2 df_mov_col['2021-07'] = df_mov_col['2021-07'].astype('Int64')
3 df_mov_col['2021-08'] = df_mov_col['2021-08'].astype('Int64')
4 df_mov_col['2021-09'] = df_mov_col['2021-09'].astype('Int64')
5 df_mov_col['2021-10'] = df_mov_col['2021-10'].astype('Int64')
6 df_mov_col['2021-11'] = df_mov_col['2021-11'].astype('Int64')
7 df_mov_col['2021-12'] = df_mov_col['2021-12'].astype('Int64')
8 # Removendo colunas
9 df_mov_col = df_mov_col.drop('Código Nacional', axis=1)
10 df_mov_col = df_mov_col.drop('Grupo Econômico', axis=1)
11 df_mov_col = df_mov_col.drop('Código IBGE Município', axis=1)
12 # Trocando o tipo dos dados de float para inteiro
13 df_mov_col2['2022-01'] = df_mov_col2['2022-01'].astype('Int64')
14 df_mov_col2['2022-02'] = df_mov_col2['2022-02'].astype('Int64')
15 df_mov_col2['2022-03'] = df_mov_col2['2022-03'].astype('Int64')
```

```
1 # Removendo colunas
2 df_mov_col2 = df_mov_col2.drop('Código Nacional', axis=1)
3 df_mov_col2 = df_mov_col2.drop('Grupo Econômico', axis=1)
4 df_mov_col2 = df_mov_col2.drop('Código IBGE Município', axis=1)
5 # Renomeando colunas devido à compatibilidade com a Bigquery
6 (df_mov_col.rename(columns = {'2021-07': '_2021-07',
7                               '2021-08': '_2021-08',
8                               '2021-09': '_2021-09',
9                               '2021-10': '_2021-10',
10                              '2021-11': '_2021-11',
11                              '2021-12': '_2021-12'}, inplace = True)
12 )
13
14 (df_mov_col2.rename(columns = {'2022-01': '_2022-01',
15                               '2022-02': '_2022-02',
16                               '2022-03': '_2022-03'}, inplace = True)
17 )
```



# NORMALIZAÇÃO DOS DADOS (PYSPARK)

```
1  # Eliminando a coluna 'Operadora' pois ela possuiu somente um valor
2  df_spark_cob = df_spark_cob.drop('Operadora')
3  # Eliminando a coluna 'Localidade Agregadora' pois ela está vazia
4  df_spark_cob = df_spark_cob.drop('Localidade Agregadora')
5  # Substituindo os valores nulos por 'Não informado'
6  df_spark_cob = df_spark_cob.fillna(0)
7  # Renomeando as colunas devido à compatibilidade com Spark
8  df_spark_cob = (df_spark_cob.withColumnRenamed('Código Setor Censitário','Codigo_Setor_Censitario')
9                      .withColumnRenamed('Tipo Setor','Tipo_Setor')
10                     .withColumnRenamed('Código Localidade','Codigo_Localidade')
11                     .withColumnRenamed('Nome Localidade','Nome_Localidade')
12                     .withColumnRenamed('Categoria Localidade','Categoria_Localidade')
13                     .withColumnRenamed('Código Município','Codigo_Municipio')
14                     .withColumnRenamed('Município','Municipio')
15                     .withColumnRenamed('Região','Regiao')
16                     .withColumnRenamed('Área (km2)','Area_km2')
17                     .withColumnRenamed('Domicílios','Domicilios')
18                     .withColumnRenamed('Percentual Cobertura','Percentual_Cobertura')
19  )
```

# UPLOAD PARA BANCO POSTGRESQL/GCP

```
1 from sqlalchemy import create_engine
2
3 # Credenciais de acesso para o banco de dados Postgres sql criado na GCP
4 engine = create_engine("postgresql://postgres:xk5fMiHn74FDBw_@34.132.83.213/projeto_final")
5 df_cob.to_sql('Cobertura_Todas',engine,if_exists='replace',index=False)
```



# QUERY COM SPARK

```
1 # Acessos por Estado
2 spark.sql("select uf,sum(acessos) as acessos,round((sum(acessos)/122858439)*100,2) as percent From bandfixa group by uf order by percent desc ").show(3)
```

uf	acessos	percent
SP	186970496	31.47
MG	64044645	10.78
RJ	55161665	9.29
RS	42030644	7.07
PR	41986126	7.07
SC	32525189	5.47
BA	21765710	3.66
CE	19763136	3.33
GO	17719883	2.98
DF	11997008	2.02
PE	11742930	1.98
ES	10809381	1.82
MT	8893073	1.5
RN	8773989	1.48
PA	8796813	1.48
PB	8053971	1.36
MS	7077543	1.19
AM	6545058	1.1
MA	6436060	1.08



# QUERY COM SPARK

```
1 # Empresas responsáveis pelas maiores quantidades de acesso
2 spark.sql("select empresa,sum(acessos) as acessos,round((sum(acessos)/122858439)*100,2) as percent From bandfixa group by empresa order by percent desc ").show(4)
```

empresa	acessos	percent
CLARO	146286095	24.62
VIVO	94990779	15.99
OI	77473755	13.04
Brisanet Servicos...	11573482	1.95



# QUERY COM SPARK

```
1 # Tecnologia e meio de acessos
2 spark.sql("select Tecnologia,meioacesso, sum(acessos) as acessos,round((sum(acessos)/122858439)*100,2) as percent From bandfixa group by tecnologia, meioacesso order by percent desc ").show(7)
```

Tecnologia	Meio_Acesso	acessos	percent
FTTH	Fibra	314309777	52.91
HFC	Cabo Coaxial	137634016	23.17
ADSL2	Cabo Metálico	52803821	8.89
ETHERNET	Fibra	19616141	3.3
ETHERNET	Rádio	9599155	1.62
ETHERNET	Cabo Metálico	8967252	1.51
Wi-Fi	Rádio	8766433	1.48

# BIGQUERY

Query Editor interface showing three SQL queries for BigQuery. The first query is highlighted with a red box, and a red arrow points from it to a table of results.

```
1 SELECT Operadora, ROUND(SUM(Area_km2_Coberta),2) AS Cobertura_4g FROM `macro-mercury-349020.ProjetoFinal.Cobertura_Municipios`
2 WHERE Tecnologia_Cobertura = '4G'
3 GROUP BY Operadora HAVING
4 Operadora = 'VIVO' OR
5 Operadora = 'CLARO' OR
6 Operadora = 'TIM' OR
7 Operadora = 'OI' OR
8 Operadora = 'NEXTEL' ORDER BY Cobertura_4g DESC;
```

Row	Operadora	Cobertura_4g
1	CLARO	513038.52
2	VIVO	504152.35
3	TIM	424294.77
4	OI	131109.62
5	NEXTEL	69570.68



# LIDERANÇA DE COBERTURA MÓVEL NO BRASIL

Row	Operadora	Cobertura
1	VIVO	1560608.61
2	CLARO	1442949.09
3	TIM	885242.64
4	OI	529072.33
5	NEXTEL	279764.12



# BIGQUERY

```
1 SELECT Regiao, Tipo_Setor, ROUND(100*(SUM(Area_km2*(Percentual_Cobertura/100))/sum(Area_km2)),2) AS  
   porcentagem_cobertura FROM `macro-mercury-349020.ProjetoFinal.Cobertura_Setor` WHERE Tecnologia = '2G3G4G' GROUP  
   BY Regiao, Tipo_Setor order by Regiao;
```

Press Alt+F1 for Accessibility Options

## Query results

 SAVE RESULTS ▾

 EXPLORE DATA ▾



### JOB INFORMATION

### RESULTS

### JSON

### EXECUTION DETAILS

Row	Regiao	Tipo_Setor	porcentagem_cobertura
1	Centro-oeste	Rural	13.46
2	Centro-oeste	Urbano	83.25
3	Nordeste	Rural	30.85
4	Nordeste	Urbano	98.67
5	Norte	Rural	10.41
6	Norte	Urbano	93.18
7	Sudeste	Rural	36.96
8	Sudeste	Urbano	92.89
9	Sul	Rural	38.14
10	Sul	Urbano	96.4



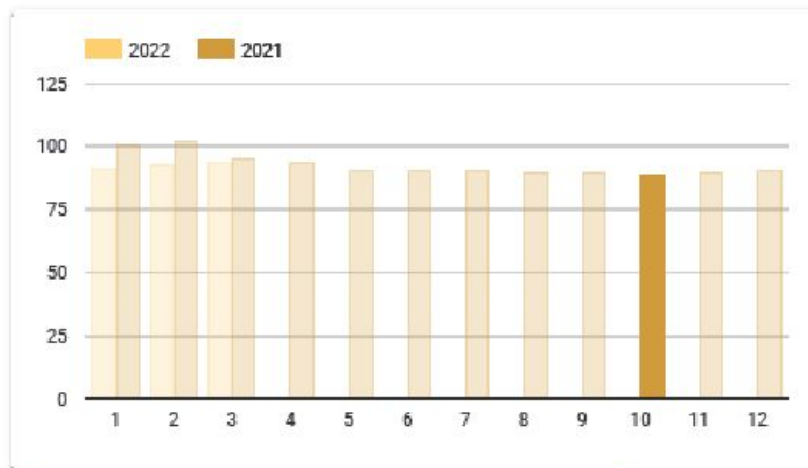
# DATA STUDIO

## Banda Larga fixa

Total de Acessos

40,8 mi

### Acessos por Mês

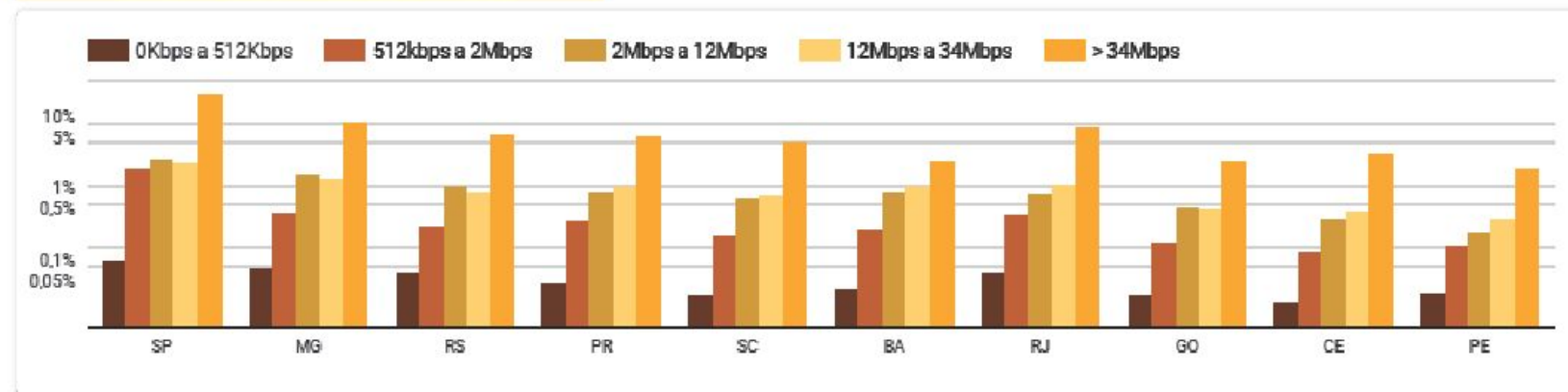


### UF

### Acessos

SP	12,8 mi
MG	4,4 mi
RJ	3,7 mi
RS	2,9 mi
PR	2,9 mi
SC	2,2 mi
BA	1,5 mi
CE	1,4 mi
GO	1,2 mi

### Faixa de velocidade por estado



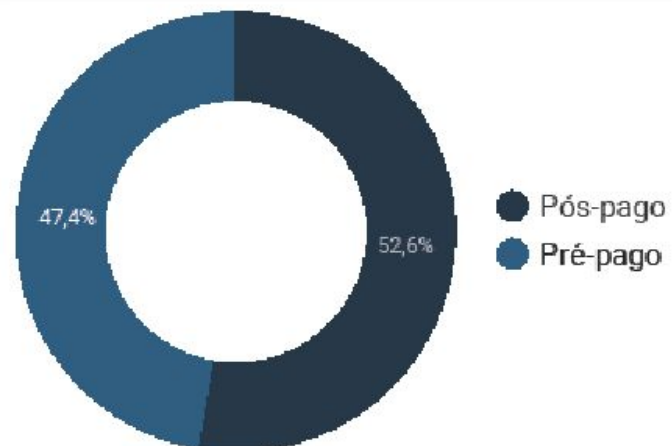
# DATA STUDIO

## Telefonia Móvel

Total de Acessos

3,4 bi

### Modalidade de cobrança

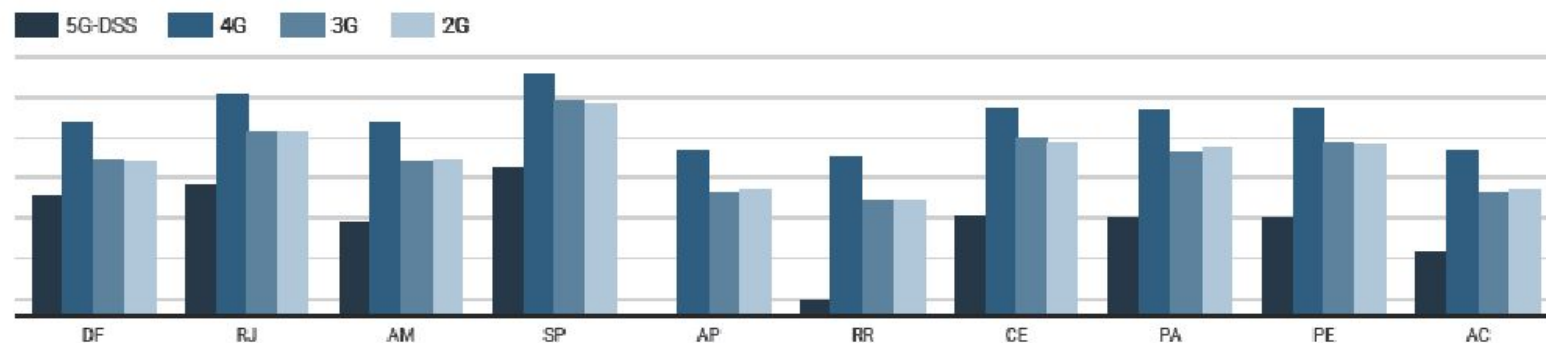


### UF

### Acessos

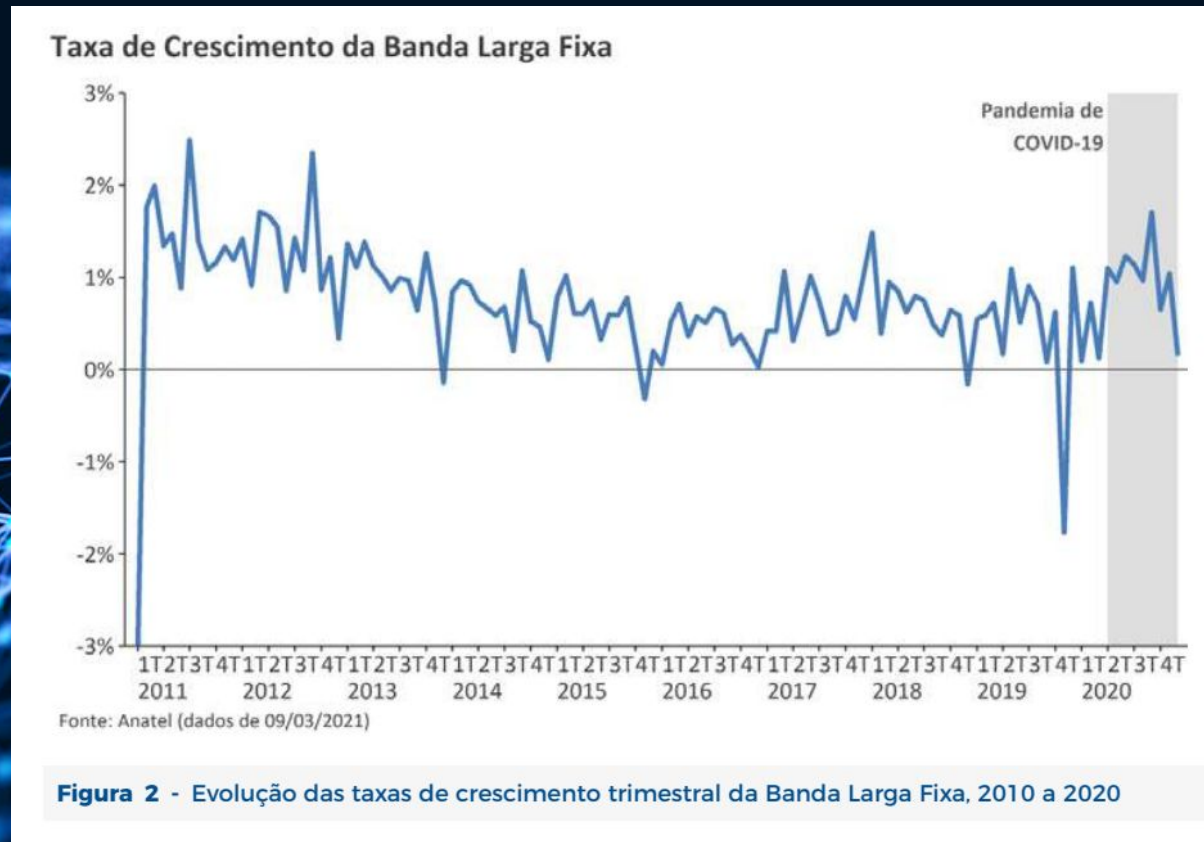
SP	1 bi
MG	326,8 mi
RJ	280,2 mi
BA	203,1 mi
RS	184,2 mi
PR	181,4 mi
PE	137,7 mi
CE	131,9 mi
SC	114,7 mi
GO	110,2 mi

### Tecnologia por estado





# RELATÓRIO DE ACOMPANHAMENTO DO SETOR DE TELECOMUNICAÇÕES / ANATEL



OBRIGADO

