

Centro Universitário Instituto de
Educação Superior de Brasília
Departamento de Pós graduação
Especialização em Big Data, BI e Analytics Aplicados aos
Negócios

Projeto Integrador I

Aluna: Mariana Borges de Sampaio

Professor: William de Almeida Silva

Disciplina: Fundamentos de Big Data e Conhecimentos de Dados

Agosto
2023

Centro Universitário Instituto de
Educação Superior de Brasília
Departamento de Pós graduação
Especialização em Big Data, BI e Analytics Aplicados aos
Negócios

Projeto Integrador I

Primeiro relatório do Projeto Integrador 1 da disciplina de Fundamentos de Big Data e Conhecimentos de dados.

Aluna: Mariana Borges de Sampaio

Professor: William de Almeida Silva

Agosto
2023

Conteúdo

1	Introdução	1
2	Dicionário de dados e Análise de dados	2
2.1	Dicionário de dados	2
2.2	Análise de dados	2
3	Considerações finais	5
	Bibliografia	6
	Anexo	7

1 Introdução

Ao longo dos anos as empresas de todos os países e de todos os continentes foram se desenvolvendo, cada um em sua área específica, sejam empresas de corretoras de seguro, empreiteiras, empresas de tecnologia, empresas de saúde, entre os mais diversos tipos e variados setores em que uma empresa pode se aplicar. Apesar dessas empresas terem conceitos e fundamentos diferentes visto que elas são pertencentes a mundos diferentes, tem-se detém um conjunto de dados, esse dado muitas vezes pode não ser estruturado, sendo ele um dado bruto, um dado que só teria sentido se fosse tratado a fim de se tornar uma informação.

Sendo assim, cada empresa existente até os dias de hoje inevitavelmente possui dado, logo possui informação. Esse dado pode ser ele armazenado em um banco de dados estruturado ou pode ser um big data, sendo este um dado não estruturado. O big data tem por definição ser um conjunto de dados que possuem uma maior variedade, maior volume e estes dados têm uma maior velocidade. Sendo assim, o seu armazenamento é diferente do que um dado estruturado. Para isso conforme foram-se passando os anos foi sendo necessário nos meios do trabalho ter profissionais que são voltados para essa área.

2 Dicionário de dados e Análise de dados

2.1 Dicionário de dados

A fonte de dados em questão foi extraída do Kaggle, conforme o que foi passado na descrição da atividade ativa.

Ao analisar o documento pode-se perceber a existência das seguintes colunas:

- id;
- roomid/id;
- noteddate;
- temp;
- out/in.

A seguir, tem-se a coluna com seu respectivo significado.

- id - corresponde à um valor alfanumérico que indica uma identificação única, sendo assim cada linha é um momento em que foi registrada a temperatura do edifício;
- roomid/id - corresponde ao nome da sala do edifício em que foi medida a temperatura. Para a identificação das medidas, tem-se o id que é um alfanumérico, visto que esses casos foram todos registrados em momentos aleatórios.
- noteddate - corresponde a data de anotação que foi identificada a temperatura no edifício, sendo está registrada em intervalos aleatórios;
- temp - corresponde a temperatura que foi atingida no edifício;
- out/in - indica se a classificação de cada se a temperatura medida foi no exterior ou no interior do edifício.

2.2 Análise de dados

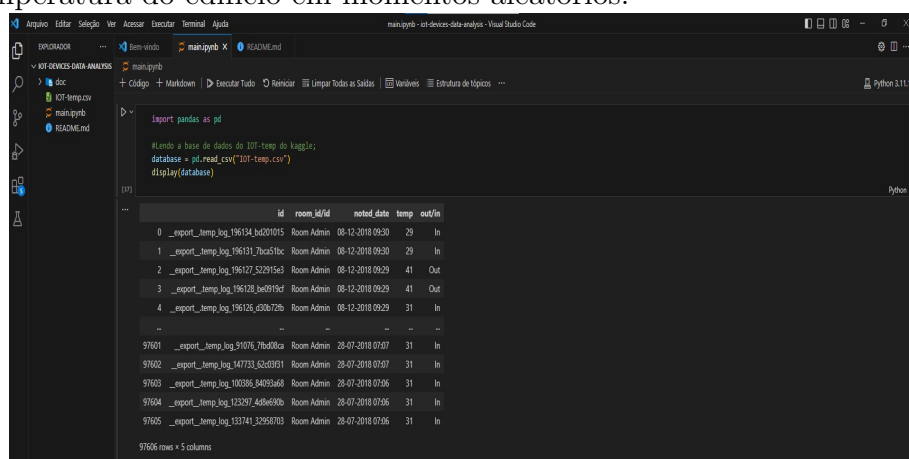
Para realizar a análise de dados, foi necessário entender um pouco mais sobre o contexto do arquivo que foi extraído do kaggle. Os dados apresentados correspondem a dados oriundos de um sistema de monitoramento de temperatura da sala do edifício empresarial (admin), tanto no exterior como no interior do edifício.

Para realizar a análise dos dados foi feita a leitura do que foi apresentado a fim de entender o contexto da análise. A partir disso, foi desenvolvido um código no arquivo (ipynb) jupyter notebook que possui todo o código desenvolvido que foi criado com sessão do aplicativo jupyter notebook. Para isso, utilizei o vscode studio com a extensão referente ao jupyter notebook. Para manter o versionamento, utilizei o github, dessa forma, conforme avançava na análise e no código realizada um update no github que foi destinado a esse trabalho, <https://github.com/sampaioariana/iot-devices-data-analysis>. Para o desenvolvimento do arquivo de relatório, utilizei o overleaf, que também conforme avançava foi sendo atualizado no github.

Na análise dos dados, foram seguidos cinco passos durante o desenvolvimento, sendo estes, os seguintes:

- Importar os dados;
- Visualizar a base;
- Tratamento de erros;
- Análise inicial dos dados;
- Análise profunda dos dados.

Para realizar todos esses passos, foram utilizadas as bibliotecas pandas e statistics. Para realizar a primeira análise e entender o os dados, importei os dados utilizando a biblioteca pandas, dessa forma consegui ver as colunas e o conteúdo das colunas, a partir dessa primeira análise, foi possível ver que existem 97606 linhas e 5 colunas, sendo assim, tem-se 97606 registros de temperatura do edifício em momentos aleatórios.



```

import pandas as pd

#Lendo a base de dados do IOT-temp do Kaggle;
database = pd.read_csv("IoT-temp.csv")
display(database)

```

	id	room_id	outed_date	temp	out/in
0	__report__temp_log_196134_b4d01015	Room Admin	08-12-2018 09:30	29	In
1	__report__temp_log_196131_7bca5fbc	Room Admin	08-12-2018 09:30	29	In
2	__report__temp_log_196127_522919e3	Room Admin	08-12-2018 09:29	41	Out
3	__report__temp_log_196128_ba0919cf	Room Admin	08-12-2018 09:29	41	Out
4	__report__temp_log_196126_630b770b	Room Admin	08-12-2018 09:29	31	In
...
97601	__report__temp_log_916076_7b0d0ca	Room Admin	28-07-2018 07:07	31	In
97602	__report__temp_log_147733_62d0351	Room Admin	28-07-2018 07:07	31	In
97603	__report__temp_log_100366_84093a68	Room Admin	28-07-2018 07:06	31	In
97604	__report__temp_log_123297_4a9b690b	Room Admin	28-07-2018 07:06	31	In
97605	__report__temp_log_133741_32950703	Room Admin	28-07-2018 07:06	31	In

97606 rows x 5 columns

A fim de entender qual era o tipo de conteúdo em cada linha, foi feito um comando com uma descrição do conteúdo de cada linha.

```
EXPLORADOR
  NOT DEVICIES DATA ANALYST
  C:\temp\cv
  main.pyb
  README.md

main.pyb
  README.md

D>
#Para visualizar os tipos de dados:
#Como a base é um pandas observar que todas as linhas estão preenchidas, não sendo necessário retirar dados nulos;
#Verificar a coluna 'out/in' é inteiro, enquanto as demais são object;
display(database.info())

#Como a coluna 'id' indica apenas a quantidade de lts registrados e essa contagem já será feita, essa informação pode ser retirada.
# lista >> axis = 0
# columns >> axis = 1
database = database.drop("id", axis = 1)
display(database.info())

[14]

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 97686 entries, 0 to 97685
Data columns (total 5 columns):
 # Column Non-Null Count  Dtype
---  --
 0 id      97686 non-null object
 1 row_id  97686 non-null object
 2 noted_date  97686 non-null object
 3 temp    97686 non-null int64
 4 out/in   97686 non-null object
dtypes: int64(1), object(4)
memory usage: 3.7+ MB

None

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 97686 entries, 0 to 97685
Data columns (total 4 columns):
 # Column Non-Null Count  Dtype
---  --
 0 row_id  97686 non-null object
 1 noted_date  97686 non-null object
 2 temp    97686 non-null int64
 3 out/in   97686 non-null object
dtypes: int64(1), object(3)
memory usage: 3.4v MB

None
```

3 Considerações finais

Bibliografia

AGUIRRE, L. A. Introdução à Identificação de Sistemas, Técnicas Lineares e Não lineares Aplicadas a Sistemas Reais. Belo Horizonte, Brasil, EDUFMG. 2004.

Anexo