NÃO PODE FALTAR

INTRODUÇÃO A MANIPULAÇÃO DE DADOS EM F

Vanessa Cadan Scheffer

MÉTODOS PARA MANIPULAÇÃO DE DADOS

Além dos métodos para carregar e salvar dados, a biblioteca pandas pos para a transformação dos dados e a extração de informação.



Fonte: Shutterstock.

(sistemas de armazenamento em cloud) ou em bancos de dado origens (e até mais), a biblioteca possui métodos capazes de fa dados e carregar em um DataFrame.

Todos os métodos capazes de fazer a leitura dos dados estrutuprefixo pd.read_xxx, onde pd é o apelido dado no momento da biblioteca e xxx é o restante da sintaxe do método. Além de faz biblioteca possui diversos métodos capazes de escrever o Data arquivo, em um banco ou ainda simplesmente copiar para a ár do sistema operacional. O Quadro 4.2, apresenta todos os métodos capazes de escrita. Veja que são suportados tanto a leitura de arquivos de binários e de bancos.

Quadro 4.2 - Métodos para leitura e escrita de dados estrut

Tipo de dado	Descrição do dado	Método para leitura	N
texto	CSV	read_csv	t
texto	Fixed-Width texto	read_fwf	
texto	J <u>SON</u>	read_json	t

Tipo de dado	Descrição do dado	Método para leitura	N e
binário	<u>Feather Format</u>	read_feather	to
binário	<u>Parquet Format</u>	read_parquet	to
binário	ORC Format	read_orc	
binário	<u>Msgpack</u>	read_msgpack	to
binário	<u>Stata</u>	read_stata	to
binário	<u>SAS</u>	read_sas	
binário	<u>SPSS</u>	read_spss	
binário	<u>Python Pickle Form</u> <u>at</u>	read_pickle	te
SQL	<u>SQL</u>	read_sql	to
SQL	Google BigQuery	read_gbq	to

Fonte: adaptado de https://pandas.pydata.org/pandas-docs/stable/us

Dentre todos os possíveis métodos para leitura, nessa aula var

CSV (comma-separated values - valores separados por vírgula) arquivo, nos quais os dados são separados por um delimitador esse delimitador é uma vírgula (por isso o nome), mas na práti pode ser criado com qualquer delimitador, por exemplo, por pipe (|), dentre outros. Por ser um arquivo de texto, é fácil de sistema, por isso se tornou tão democrático.

LEITURA DE ISON E CSV COM PANDAS

Agora que vamos começar a ver as implementações, vamos faz biblioteca, como já sabemos, só precisamos importar uma únio ou no script .py.

In [1]: import pandas as pd

A leitura de um arquivo JSON deve ser feita com o método: pandas.read_json(path_or_buf=None, orient=None, typ='frame', convert_axes=None, convert_dates=True, keep_default_dates=True precise_float=False, date_unit=None, encoding=None, lines=Falcompression='infer'). Os detalhes de cada parâmetro podem s documentação oficial: https://pandas.pydata.org/pandas-docs/i/pandas.read_json.html. O único parâmetro que é obrigatório dados é o "path_or_buf", no qual deve ser passado um caminho

20/08/2020 npf_ldkls202_u4s2_lin_pro

In [2]:

pd.read_json("https://api.bcb.gov.br/dados/serie/ formato=json").head()

Out[2]:

	data	valor
0	04/06/1986	0.065041
1	05/06/1986	0.067397
2	06/06/1986	0.066740
3	09/06/1986	0.068247
4	10/06/1986	0.067041

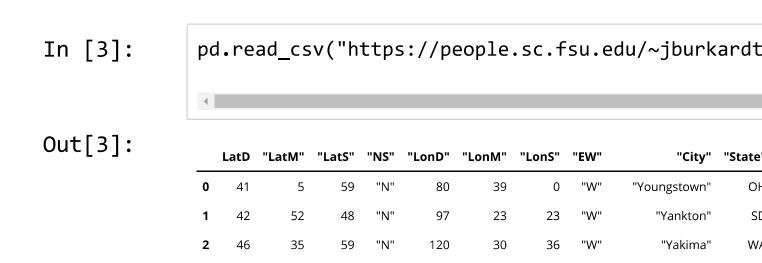
A leitura de um arquivo CSV deve ser feita com o método:

pandas.read_csv(filepath_or_buffer: Union[str, pathlib.Path, sep=',', delimiter=None, header='infer', names=None, index_color usecols=None, squeeze=False, prefix=None, mangle dupe cols=Tr engine=None, converters=None, true_values=None, false_values= skipinitialspace=False, skiprows=None, skipfooter=0, nrows=No keep default na=True, na filter=True, verbose=False, skip bla parse dates=False, infer datetime format=False, keep date col date_parser=None, dayfirst=False, cache_dates=True, iterator= chunksize=None, compression='infer', thousands=None, decimal: lineterminator=None, quotechar='"', quoting=0, doublequote=Tr escapechar=None, comment=None, encoding=None, dialect=None, error_bad_lines=True, warn_bad_lines=True, delim_whitespace=F low memory=True, memory map=False, float precision=None). São

o que proporciona uma versatilidade incrível para esse método

que estão separados por vírgula. O parâmetro header, tem com 'infer', que significa que o método realiza a inferência para os r partir da primeira linha de dados do arquivo.

Na entrada 3, estamos fazendo a leitura de uma fonte CSV, cuj separados por vírgula, logo não foi preciso especificar um delir



"N"

"N"

71

89

0

11

48

46

"W"

"W"

"Worcester"

W

"Wisconsin Dells"

MANIPULAÇÃO DE DADOS COM PANDAS

12

48

3

42

43

16

37

Além de vários métodos para carregar e salvar os dados, a bibl possui uma diversidade de métodos para a transformação dos de informação para áreas de negócio. Nessa seção vamos conl

O trabalho com dados: capturar os dados em suas origens, faz nos dados a fim de padronizá-los, aplicar técnicas estatísticas o algoritmos de machine/deep learning feito por engenheiros e

por exemplo, se você pretende fazer um financiamento o investimento, precisa olhar a taxa Selic, pois ela influenci pago ou ganho. "Selic é a sigla para Sistema Especial de l Custódia, um programa totalmente virtual em que os títu Nacional são comprados e vendidos diariamente por ins financeiras" (NUBANK, p. 1, 2020). Quem decide, de fato, Copom (Comitê de Política Monetária do Banco Central), se reúne para determinar se a taxa aumenta ou diminui. taxa está alta, os financiados podem ficar mais caros e o se a taxa está mais baixa, então os financiamentos ficam Resumindo, quando a taxa Selic aumenta o economia de ela abaixa a economia aquece, isso é preciso para contro

Agora que já conhecemos que a taxa Selic influencia, nossos finivestimentos e até mesmo o que compramos no mercado, va informações disponibilizadas pelo governo e fazer algumas ancomeçar pela etapa de extração e transformação dos dados.

ETAPA DE CAPTURA E TRANSFORMAÇÃO/ PADRONIZAÇÃO DO

A extração dos dados pode ser realizada por meio do método guardando em um DataFrame (DF) pandas. Ao carregar os dado podemos visualizar quantas linhas o colunas, hom como os tir

linhas, então não existem valores faltantes. Quanto ao tipo de object, ou seja, são todos do tipo strings ou existe mistura de t float.

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8552 entries, 0 to 8551
Data columns (total 2 columns):
data 8552 non-null object
valor 8552 non-null float64
dtypes: float64(1), object(1)
memory usage: 133.7+ KB

None

REMOVER LINHAS DUPLICADAS

devemos fazer é remover os dados duplicados. Certamente, qu depende da área de negócio e do problema a ser resolvido. Po queremos manter o registro da compra atual, ou queremos m compra. Um DataFrame da bilioteca pandas possui o método meu_df.drop_duplicates() que permite fazer essa remoção de d Observe a entrada 5 da linha de código a seguir, usamos o mé as linhas duplicadas, pedindo para manter o último registro (ke do parâmetro inplace=True, estamos fazendo com que a trans do DataFrame, na prática estamos sobrescrevendo o objeto na inplace não seja passado, a transformação é aplicada, mas seja, o DF continua da mesma forma anterior a transforma parâmetro interessante do método é o subset, que permite qu

Para o carregamento de uma base de dados, um dos primeiros

In [5]: df_selic.drop_duplicates(keep='last', inplace=Tru

dados duplicado seja feita com base em uma ou mais colunas.

CRIAR NOVAS COLUNAS

tipo de dados das datas, embora cada valor seja do tipo "date" ainda obtemos uma coluna object, para que de fato, a bibliotecum tipo data, vamos ter que utilizar o método da própria biblioconversão.

```
In [6]:
```

```
from datetime import date
from datetime import datetime as dt

data_extracao = date.today()

df_selic['data_extracao'] = data_extracao
 df_selic['responsavel'] = "Autora"

print(df_selic.info())
 df_selic.head()
```

memory usage: 334.1+ KB

None

Out[6]:

	data	valor	data_extracao	responsavel
0	04/06/1986	0.065041	2020-07-19	Autora
1	05/06/1986	0.067397	2020-07-19	Autora
2	06/06/1986	0.066740	2020-07-19	Autora
3	09/06/1986	0.068247	2020-07-19	Autora
4	10/06/1986	0.067041	2020-07-19	Autora

guardamos dentro da própria coluna, dessa forma os valores s linha 1, usamos a notação pd.to_datetime(), porque é um méto não do DF.

Na entrada 7 (linha 1), o método recebe a coluna com os valore (df_selic['data']) e um segundo parâmetro, indicando que no (antes da conversão) o dia está primeiro (dayfirst=True). Em se como a coluna "data_extracao" foi criada com o método today correto para a conversão. Nessa conversão usamos o método transforma os dados de uma coluna (que é uma Series) em um nesse caso, o tipo datetime especificado. Com astype() podem valores das colunas, por exemplo, transformando todos em flo outro tipo. Veja que agora, ao usar o método info(), temos que do tipo datetime (datetime da biblioteca pandas). O formato re dia é um padrão do datetime64[ns], que segue o padrão interr ano vem primeiro, seguido do mês e por último o dia. Podería: para transformar o traço em barra (/), mas aí o resultado seria datas.

In [7]:

```
df_selic['data'] = pd.to_datetime(df_selic['data'], dayfirst
df_selic['data_extracao'] = df_selic['data_extracao'].astype
print(df selic.info())
```

	data	valor	data_extracao	responsavel
0	1986-06-04	0.065041	2020-07-19	Autora
1	1986-06-05	0.067397	2020-07-19	Autora
2	1986-06-06	0.066740	2020-07-19	Autora
3	1986-06-09	0.068247	2020-07-19	Autora
4	1986-06-10	0.067041	2020-07-19	Autora

SERIES.STR

Muitas vezes precisamos padronizar os valores em colunas, po queremos ter certeza que a coluna "responsável" possui todas padronizadas em letras maiúsculas. Quando selecionamos um sabemos que o resultado é uma Series e esse objeto tem um r permite aplicar as funções de string para todos os valores da S

Observe o trecho de código seguir, selecionamos a coluna responsor o recurso str e aplicamos o método upper(). Dessa forma, a bible "entende" que queremos converter todos os valores dessa columaiúsculas. Como atribuímos o resultado na própria coluna, o substituído.

Out[8]:

	data	valor	data_extracao	responsavel
0	1986-06-04	0.065041	2020-07-19	AUTORA

MÉTODO SORT_VALUES()

No código a seguir, estamos usando o método sort_values() que DF, de acordo com os valores de uma coluna. Esse método é d isso a notação meu_df.metodo(). Utilizamos três parâmetros do o primeiro informando qual coluna deve ser usada para orden que seja feito em ordem decrescente (do maior para o menor) (inplace=True) significa que queremos modificar o próprio objectamos sobrescrevendo o DF.

Out[9]:

		data	valor	data_extracao	responsavel
855	51	2020-07-17	0.008442	2020-07-19	AUTORA
855	60	2020-07-16	0.008442	2020-07-19	AUTORA
854	19	2020-07-15	0.008442	2020-07-19	AUTORA
854	18	2020-07-14	0.008442	2020-07-19	AUTORA
854	! 7	2020-07-13	0.008442	2020-07-19	AUTORA

MÉTODO RESET_INDEX() E SET_INDEX()

Ao fazermos a ordenação dos dados com o método sort_value índices dos cinco primeiros registros é 8551, 8550...85XX. Nenh

em uma nova coluna e inplace, informa para gravar as alteraçõ objeto. Veja na saída que agora os cinco primeiros registros, po a 4.

Out[10]:

	data	valor	data_extracao	responsavel
0	2020-07-17	0.008442	2020-07-19	AUTORA
1	2020-07-16	0.008442	2020-07-19	AUTORA
2	2020-07-15	0.008442	2020-07-19	AUTORA
3	2020-07-14	0.008442	2020-07-19	AUTORA
4	2020-07-13	0.008442	2020-07-19	AUTORA

Durante a transformação dos dados, pode ser necessário defir para os índices, ao invés de usar o range numérico. Essa transfeita usando o método meu_df.set_index(). O método permite valores usando uma coluna já existente ou então passando un igual a quantidade de linhas.

Observe os códigos nas entradas 11 e 12. Na entrada 11, estan lista que usada os índices do DF, adicionando um prefixo para 2 são impressos os cinco primeiros itens da nova lista. Na entradefinindo o novo índice com base na lista criada. Veja que o parecebe como parâmetro uma lista de lista e o segundo parâme

20/08/2020 npf_ldkls202_u4s2_lin_pro

Out[12]:

	data	valor	data_extracao	responsavel
selic_0	2020-07-17	0.008442	2020-07-19	AUTORA
selic_1	2020-07-16	0.008442	2020-07-19	AUTORA
selic_2	2020-07-15	0.008442	2020-07-19	AUTORA
selic_3	2020-07-14	0.008442	2020-07-19	AUTORA
selic_4	2020-07-13	0.008442	2020-07-19	AUTORA

Ao especificar um índice com valor conceitual, podemos usar a idxmax() para descobrir qual o índice do menor e do maior de os exemplos a seguir.

ETAPA DE EXTRAÇÃO DE INFORMAÇÕES

Agora que já fizemos as transformações que gostaríamos, pod de extração de informações.

FILTROS COM LOC

Um dos recursos mais utilizados por equipes das áreas de dad

DataFrames da biblioteca pandas possuem uma propriedade o pandas.pydata.org/pandas-docs/stable/reference/api/pandas.l Essa propriedade permite acessar um conjunto de linhas (filtra do índice ou por um vetor booleano (vetor de True ou False).

Vamos começar explorando o filtro pelos índices. Entrada 14 e (loc), o registro que possui índice 'selic_0', como resultado obté entrada 15, foram filtrados três registros, para isso foi necessá contendo os índices, como resultado obtivemos um novo DF. Ne fizemos um fatiamento (slice), procurando um intervalo de índ

In [14]:	df_selic.loc['selic_0']
Out[14]:	data 2020-07-17 00:00:00
	valor 0.008442
	data_extracao 2020-07-19 00:00:00
	responsavel AUTORA
	Name: selic_0, dtype: object
In [15]:	<pre>df_selic.loc[['selic_0', 'selic_4', 'selic_200']]</pre>
Out[15]:	data valor data_extracao responsavel

			_	<u> </u>
selic_0	2020-07-17	0.008442	2020-07-19	AUTORA
selic_4	2020-07-13	0.008442	2020-07-19	AUTORA
selic_200	2019-09-30	0.020872	2020-07-19	AUTORA

In [16]: df_selic.loc[:'selic_5']

Out[16]:

	data	valor	data_extracao	responsavel
selic_0	2020-07-17	0.008442	2020-07-19	AUTORA
selic_1	2020-07-16	0.008442	2020-07-19	AUTORA
selic 2	2020-07-15	0.008442	2020-07-19	AUTORA

```
• df_selic.loc[['selic_0', 'selic_4', 'selic_200']]['valor'
```

• df_selic.loc[['selic_0', 'selic_4', 'selic_200']][['valor

Out[18]:

	valor	data_extracao
selic_0	0.008442	2020-07-19
selic_4	0.008442	2020-07-19
selic_200	0.020872	2020-07-19

selic_1 2020-07-16 0.008442

Antes de vermos a criação de filtros para o loc com condições le mencionar que existe também a propriedade iloc, a qual filtra considerando a posição que ocupam no objeto. Veja no exempusando o **iloc** para filtrar os 5 primeiros registros, usando a mentio de listas. Essa propriedade não possui a opção de tocolunas. Veja um exemplo a seguir.

2020-07-19

AUTORA

Observe a partir da entrada 20 a seguir. Estamos utilizando um relacional, para testar se os valores da coluna 'valor', são meno Armazenamos o resultado em uma variável chamada teste. Ve teste é uma Series, e na linha 4 teste[:5], estamos imprimindo resultados do teste lógico.

```
In [20]:
             teste = df_selic['valor'] < 0.01</pre>
              print(type(teste))
             teste[:5]
             <class 'pandas.core.series.Series'>
Out[20]:
             selic_0
                         True
             selic_1
                         True
             selic_2
                         True
             selic_3
                         True
             selic_4
                         True
             Name: valor, dtype: bool
```

No código, entrada 21 a seguir, realizamos mais um teste lógic da taxa é do ano de 2020. Para isso, utilizamos o método to_da converter a string para data e então fazer a comparação.

parênteses, senão ocorre um erro. Observe o código a seguir, linha seguinte, temos a construção de dois novos testes, o prin operação AND e o segundo usando OR.

```
In [22]: teste3 = (df_selic['valor'] < 0.01) & (df_selic['
pd.to_datetime('2020-01-01'))

teste4 = (df_selic['valor'] < 0.01) | (df_selic['
pd.to_datetime('2020-01-01'))

print("Resultado do AND:\n")
print(teste3[:3])

print("Resultado do OR:\n")
print(teste4[:3])</pre>
```

Resultado do AND:

```
selic_0 True
selic_1 True
selic_2 True
dtype: bool
Resultado do OR:
selic_0 True
selic_1 True
selic_1 True
```

dtype: bool

Agora que já sabemos criar as condições, basta aplicá-las no D o filtro. A construção é feita passando a condição para a propr o código a seguir. Na linha 1 estamos criando a condição do fil

Series hooleana) e na entrada 23 nassamos como narâmetro i

20/08/2020 npf_ldkls202_u4s2_lin_pro

	data	valor	data_extracao	responsavel
selic_5	2020-07-10	0.008442	2020-07-19	AUTORA
selic_6	2020-07-09	0.008442	2020-07-19	AUTORA
selic_7	2020-07-08	0.008442	2020-07-19	AUTORA
selic_8	2020-07-07	0.008442	2020-07-19	AUTORA
selic_9	2020-07-06	0.008442	2020-07-19	AUTORA
selic_10	2020-07-03	0.008442	2020-07-19	AUTORA
selic_11	2020-07-02	0.008442	2020-07-19	AUTORA
selic_12	2020-07-01	0.008442	2020-07-19	AUTORA
selic_13	2020-06-30	0.008442	2020-07-19	AUTORA
selic_14	2020-06-29	0.008442	2020-07-19	AUTORA
selic_15	2020-06-26	0.008442	2020-07-19	AUTORA
selic_16	2020-06-25	0.008442	2020-07-19	AUTORA
selic_17	2020-06-24	0.008442	2020-07-19	AUTORA
selic_18	2020-06-23	0.008442	2020-07-19	AUTORA
selic_19	2020-06-22	0.008442	2020-07-19	AUTORA
selic_20	2020-06-19	0.008442	2020-07-19	AUTORA
selic_21	2020-06-18	0.008442	2020-07-19	AUTORA
selic_7606	1990-03-16	0.000000	2020-07-19	AUTORA
selic_7607	1990-03-15	0.000000	2020-07-19	AUTORA
selic_7608	1990-03-14	0.000000	2020-07-19	AUTORA

Na entrada 23, criamos primeiro a condição, guardamos na va aplicamos, mas poderíamos ter passado a condição direta: df_selic.loc[df_selic['valor'] < 0.01]. Cabe ao desenvolvedo que se sente mais a vontade. Nesse livro vamos adotar a sintar e guardar em variáveis por questões didáticas e de legibilidade

Na entrada 24 (linha 1), criamos uma condição para exibir o repartamento de 2020. Primeiro criamos duas vari armazenar as datas, na linha 4 criamos o filtro e na linha 6 o aprilemento de 2020.

20/08/2020 npf_ldkls202_u4s2_lin_pro

(df selic['data'] <= data2)</pre>

In [24]:

```
data2 = pd.to_datetime('2020-01-31')
filtro_janeiro_2020 = (df_selic['data'] >= data1)
```

data1 = pd.to_datetime('2020-01-01')

df janeiro 2020 = df selic.loc[filtro janeiro 202 df janeiro 2020.head()

Out[24]:

	data	valor	data_extracao	responsavel
selic_114	2020-01-31	0.017089	2020-07-19	AUTORA
selic_115	2020-01-30	0.017089	2020-07-19	AUTORA
selic_116	2020-01-29	0.017089	2020-07-19	AUTORA
selic_117	2020-01-28	0.017089	2020-07-19	AUTORA
selic_118	2020-01-27	0.017089	2020-07-19	AUTORA

Vamos criar mais um filtro e um novo DF para que possamos v filtros. No código a seguir, vamos criar um novo DF contendo a mês de janeiro do ano de 2019.

```
data1 = pd.to_datetime('2019-01-01')
data2 = pd.to_datetime('2019-01-31')
filtro_janeiro_2019 = (df_selic['data'] >= data1)
(df_selic['data'] <= data2)</pre>
df_janeiro_2019 = df_selic.loc[filtro_janeiro_201
```

Out[25]:

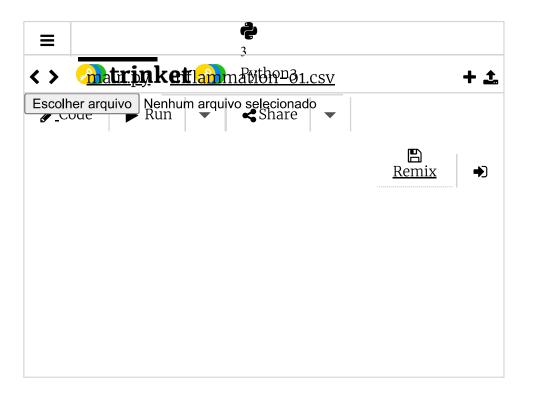
	data	valor	data_extracao	responsavel
selic_367	2019-01-31	0.02462	2020-07-19	AUTORA
selic_368	2019-01-30	0.02462	2020-07-19	AUTORA
selic_369	2019-01-29	0.02462	2020-07-19	AUTORA
selic 370	2019-01-28	0.02462	2020-07-19	AUTORA

df janeiro 2019.head()

20/08/2020 npf_ldkls202_u4s2_lin_pro

```
print('Mínimo geral = ', df_selic['valor'].min())
In [26]:
             print('Mínimo janeiro de 2019 = ', df_janeiro_201
             print('Mínimo janeiro de 2020 = ',df_janeiro_2020
             '\n')
             print('Máximo geral = ', df_selic['valor'].max())
             print('Máximo janeiro de 2019 = ', df_janeiro_201
             print('Máximo janeiro de 2020 = ',df_janeiro_2020
             '\n')
             print('Média geral = ', df_selic['valor'].mean())
             print('Média janeiro de 2019 = ', df_janeiro_2019
             print('Média janeiro de 2020 = ',df_janeiro_2020[
             '\n')
             Mínimo geral = 0.0
             Mínimo janeiro de 2019 = 0.02462
             Mínimo janeiro de 2020 = 0.017089
```

Veja como os filtros permitem começar a tirar respostas para a No ano de 2019 tanto a mínima quanto a máxima foram super 2020. A máxima geral é bem superior a máxima desses meses, média geral, que é bem superior, ou seja, nesses meses a taxa média geral, sendo que em janeiro de 2020 foi ainda pior.



Poderíamos extrair diversas outras informações dos dados. To parte do cotidiado nos engenheiros, cientistas e analistas de de engenheiros de dados mais focados na preparação e disponible os cientistas focados em responder questões de negócio, inclumodelos de machine learning e deep learning e os analistas, ta a perguntas de negócios e apresentando resultados. Se você g fizemos quem sabe não deva investir mais nessa área?!

SAIBA MAIS

Existem diversos portais que disponibilizam dados, por el kaggle, os portais brasileiros https://www.dados.gov.br/dataset, o portal https://archive.asets.php, ou o https://wincentarelbundock.github.io/Rdataset

ml Enfim ção inúmeros os respecitórios que podem ser

20/08/2020 npf_ldkls202_u4s2_lin_pro

pandas.read_sql(sql, con, index_col=None, coerce_float=Tr
parse_dates=None, columns=None, chunksize=None)

pandas.read_sql_query(sql, con, index_col=None, coerce_fl params=None, parse_dates=None, chunksize=None)

O mínimo de parâmetros que ambos métodos exigem é a instruccionexão com um banco de dados (con). A conexão com o banco ser feita usando uma outra biblioteca, por exemplo, sqlalchem diversos bancos), pyodbc para SQL Server, cx_Oracle para Orac Postgresql, dentre outras. Seguindo as recomendações da PEP bibliotecas precisam fornecer um método "connect", o qual reconexão. A síntaxe da string de conexão depende da biblioteca dados.

A seguir apresentamos de maneira didática dois exemplos, o proconexão com um banco postgresql e outro com um banco myst diferença entre eles é a importação da biblioteca específica e a Dessa forma, ao estabelecer conexão com um banco de dados instância em uma variável, basta passá-la como parâmetro do biblioteca pandas.

import psycopg2

```
conn = psycopg2.connect(conn_str)
query = "select * from XXX.YYYY"
df = pd.read_sql(query, conn)import mysql.connector
host = 'XXXXX'
port = 'XXXXX'
database = 'XXXXX'
username = 'XXXXX'
password = 'XXXXX'
conn_str = fr"host={host}, user={username}, passwd={password}
{database}"
conn = mysql.connector.connect(host="localhost", user="root";
database="bd")
query = "select * from XXX.YYYY"
df = pd.read_sql(query, conn)
```

REFERÊNCIAS E LINKS ÚTEIS

Kaggle. Titanic: Machine Learning from Disaster. Disponível em <u>e.com/c/titanic</u>). Acesso em: 20 jun. 2020.

Pandas Team. pandas documentation. Disponível em: https://pandas-docs/stable/index.html. Acesso em: 17 jun. 2020.

Pandas Team. Series. Disponível em: https://pandas.pydata.org
e/reference/api/pandas.Series.html. Acesso em: 17 jun. 2020.