


```
KeyError
File ~\AppData\Roaming\Python\Python310\site-packages\pandas\core\indexes\base.py:3621, in Index.get_loc(self, key, method, tolerance)
--> 3621     return self._engine.get_loc(casted_key)
3622 except KeyError as err:
File ~\AppData\Roaming\Python\Python310\site-packages\pandas\core\indexes\base.py:3621, in Index.get_loc(self, key, method, tolerance)
Traceback (most recent call last)
Input In [37], in <cell line: 1>()
----> 1 print(dt2.iloc[0])

File ~\AppData\Roaming\Python\Python310\site-packages\pandas\core\indexes\py:967, in _LocationIndexer._getitem_m_(self, key)
366 axis = self.axis or 0
366 maybe_callable = com.apply_if_callable(key, self.obj)
--> 967 return self._getitem_axis(maybe_callable, axis=axis)

File ~\AppData\Roaming\Python\Python310\site-packages\pandas\core\indexes\py:1202, in _iLocIndexer._getitem_axis(self, key, axis)
1200 # fall thru to straight lookup
1201 self._validate_key(key, axis)
--> 1202 return self._get_label(key, axis=axis)

File ~\AppData\Roaming\Python\Python310\site-packages\pandas\core\indexes\py:1153, in _iLocIndexer._get_label(self, label, axis)
1151 def _get_label(self, label, axis: int):
1152     # GH#567 this will fail if the label is not present in the axis.
--> 1153     return self._obj_xs(label, axis=axis)

File ~\AppData\Roaming\Python\Python310\site-packages\pandas\core\generic.py:3864, in NDFrame.xs(self, key, axis, level, drop_level)
3862         new_index = index[loc]
3863     else:
--> 3864         loc = self.index.get_loc(key)
3866         if isinstance(loc, np.ndarray):
3867             if loc.dtype == np.bool_:

File ~\AppData\Roaming\Python\Python310\site-packages\pandas\core\indexes\base.py:3623, in Index.get_loc(self, key, method, tolerance)
3621     return self._engine.get_loc(casted_key)
3622 except KeyError as err:
--> 3623     raise KeyError(key) from err
3624 except TypeError:
3625     # If we have a listlike key, _check_indexing_error will raise
3626     # InvalidIndexError. Otherwise we fall through and re-raise
3627     # the TypeError.
3628     self._check_indexing_error(key)

KeyError: 0
```

Ainda, podemos usar índices sequenciais inteiros para acessar as linhas (mesmo que as mesmas tenham outros nomes em seus índices) usando o método `iloc` (que também retorna uma `Series`):

```
In [ ]: dt2.iloc[0]
```

Usando o método `columns` obtemos um objeto do tipo `Index` com as colunas do `DataFrame`.

```
In [ ]: type(dt1.columns)
```

Podemos adicionar uma nova coluna no `DataFrame` usando as chaves com o nome da coluna:

```
In [ ]: # Todos os elementos da nova coluna são preenchidos com o valor 10
dt1["Nova coluna"] = 10
dt1
```

Da mesma forma podemos remover colunas usando o método `del`.

```
In [ ]: del dt1["peça1"]
dt1
```

Ordenando

Podemos ordenar todo um `dataframe` com base nos dados de uma coluna usando a função `sort_values()`, passando o argumento `by=` com o nome da coluna que queremos ordenar. Considerando o banco de dados das peças, o código abaixo ordena a `dataframe` de Peças pelos valores de `peça2` (note que que os índices das linhas foram alterados):

```
In [ ]: dt1.sort_values(by="peça2")
```

Carregando dados em um DataFrame

A maior utilidade dos `DataFrames` é a manipulação de dados. Dessa forma, o `Pandas` contém inúmeras maneiras para se carregar dados externos, e a estrutura de dados padrão gerada é um `DataFrame`. Inicialmente, faremos a leitura de dados no formato `csv` do próprio computador. Para isso usamos o método `pd.read_csv()`. Esse método possui diversos parâmetros (mais de 50), porém os dois principais são: o caminho do arquivo a ser lido e o delimitador dos dados. Considere o exemplo abaixo que carrega os dados "e-shop clothing 2008.csv", contido na pasta `Data`.

```
In [ ]: caminho = "G:\Meu Drive\Arquivos\UFFPR\Disciplinas\2 - Intro Mineração de Dados\Python\Datasets\le-shop
dt = pd.read_csv(caminho, sep = ";")
dt
```

Alguns repositórios de dados disponibilizam os dados diretamente da internet, de forma que podemos carregar os dados sem mesmo baixá-los no computador. Para isso só precisamos do URL dos dados. Por exemplo: https://raw.githubusercontent.com/cs109/2014_data/master/countries.csv. Lendo esses dados em um `DataFrame` temos:

```
In [ ]: caminho_url = "https://raw.githubusercontent.com/cs109/2014_data/master/countries.csv"
dt_url = pd.read_csv(caminho_url, sep = ",")
dt_url
```

Também podemos ler dados tabulares direto de uma planilha de excel com o método `read_excel`. OBS: Para isso o `pandas` requer a instalação do pacote `openpyxl`. Assim, abra um terminal e instale o pacote pelo `pip install`:

```
pip install openpyxl
```

```
In [ ]: # Podemos usar a string pura do caminho (sem barras invertidas), usando a letra 'r' antes de começar o caminho
caminho_excel = r"G:\Meu Drive\Arquivos\UFFPR\Disciplinas\2 - Intro Mineração de Dados\Python\Datasets\dt_excel
dt_excel = pd.read_excel(caminho_excel)
dt_excel
```

Exportando dados de um DataFrame

Podemos exportar os dados de um `DataFrame` usando o método `to_csv()`, em seu modo mais simples com o único argumento do caminho do arquivo.

```
In [ ]: caminho = r"G:\Meu Drive\Arquivos\UFFPR\Disciplinas\Arquivo_exportado.csv"
dt1.to_csv(caminho)
dt
```

```
In [ ]: dt1
```

Exportando dessa forma surgem 3 problemas (ou melhorias possíveis):

1. As índices das linhas foram exportados também.
2. O arquivo não fica tabulado ao abri-lo com o Excel.
3. Os nomes não estão com a acentuação correta.

Para melhorar a exportação usamos os seguintes argumentos:

```
1. index = False: Não exporta o índice das linhas.
2. sep = ";": Adicionando o separador ";" para os dados ficarem tabulares no Excel.
3. encoding = "utf-8-sig": Permite exportar acentos.
```

Assim, o código melhorado fica:

```
In [ ]: caminho = r"G:\Meu Drive\Arquivos\UFFPR\Disciplinas\Arquivo_exportado.csv"
dt1.to_csv(caminho, sep = ";", index = False, encoding = "utf-8-sig")
dt
```

Assim que carregamos um conjunto de dados, podemos obter algumas informações superficiais e rápidas sobre eles, por exemplo:

```
1. .shape: Retorna uma tupla com o número de linhas e colunas do DataFrame.
2. .info(): Mostra o nome das colunas e seus tipos de dados associados.
3. .describe(): Retorna um DataFrame com várias estatísticas descritivas sobre as colunas.
```

```
In [ ]: #dt1.shape
#dt1.info()
dt1.describe()
```

Filtros e indexação booleana

Os `DataFrames` são muito utilizados para realizarmos filtros no banco de dados. A mesma lógica da indexação booleana e do fatiamento usados nas listas e `ndarrays` pode ser utilizada aqui. Considere o conjunto de dados `Production_Data.csv`:

```
In [ ]: dt_production = pd.read_csv(r"G:\Meu Drive\Arquivos\UFFPR\Disciplinas\2 - Intro Mineração de Dados\Python\Data
dt
```

Podemos aplicar uma condição booleana em alguma das colunas, para obtermos um array de `True/False`. Por exemplo, todas as linhas em que a coluna "Activity" é igual a "Turning & Milling - Machine 4":

```
In [ ]: cond = dt_production["Activity"] == "Turning & Milling - Machine 4"
cond
```

Se atribuirmos esse vetor ao `dataframe`, teremos somente as linhas em que a `cond.` é verdadeira:

```
In [ ]: dt_production[cond]
```

Podemos escrever a mesma coisa de forma direta, ou seja, colocamos a condição diretamente no `dataframe`:

```
In [ ]: dt_production[dt_production["Activity"] == "Turning & Milling - Machine 4"]
```

Usando o método `unique()` (das `Series` em um determinado atributo (coluna), conseguimos encontrar os valores sem repetição que ocorrem nos registros dessa coluna). Por exemplo, quais são os tipos de atividade desempenhadas nesse banco de dados?

```
In [ ]: dt_production["Activity"].unique()
```

Exercícios II

1. Considerando o exercício dos pães da seção anterior: salve os dados em um `dataframe` adequado para se realizar operações. Crie 2 `dataframes` iguais usando métodos diferentes: um a partir de um dicionário e um a partir de uma lista de tuplas.
 - A. Usando o `dataframe`, encontre as somas de pesos por dias da semana e por tipo de pão.
 - B. Usando o `dataframe`, encontre a média de pesos por dias da semana e por tipo de pão.
2. Leia os dados "clientes_shopping.csv", exclua a coluna "Genre" e exporte os dados novamente em um arquivo `csv`.
3. Considere o conjunto de dados `MateriaisConstrução.xlsx`. Este conjunto contém dados referente a compra em uma loja de construção. Cada linha representa um pedido, sendo que as colunas contém os itens comprados e as células as quantidades adquiridas. Responda às seguintes questões:
 - A. Quantos registros de compra existem?
 - B. Quantos e quais os itens vendidos pela loja?
 - C. Quais as médias de vendas dos itens?
 - D. Qual é o item com a maior média de vendas?
 - E. Qual é o item que está presente na maioria das compras? Em quantas?
4. Considere o conjunto de dados `Production_Data.csv`. Este conjunto contém dados de produção, a coluna `Case_ID` indica as ordens de produção, uma ordem de produção passa por diversas atividades (`Activity`), portanto existem diversas linhas para cada `Case_ID`. A coluna `worker_ID` indica o número de identificação do funcionário que realizou a atividade. `Qty_Rejected` indica quantas peças foram perdidas na atividade executada naquela linha. `Start_Timestamp` e `Complete Timestamp` indicam as datas e horas de início e fim do processamento das atividades. Responda às seguintes questões:
 - A. Quantos trabalhadores existem nesse db?
 - B. Qual o total de peças rejeitadas no db?
 - C. Quantas ordens de produção foram processados no total?
 - D. Quais são as datas mais cedo e mais tarde de início de processamento de OPs?
 - E. O db compreende um período de quantos dias de produção?
 - F. Quais as médias de peças rejeitadas/dia e ordens de produção/dia no período todo?
 - G. Quais as médias de peças rejeitadas/dia e ordens de produção/dia nos seguintes períodos:
 - a. [2012/01/02,2012/02/01] -> janeiro
 - b. [2012/02/01,2012/03/01] -> fevereiro
 - c. [2012/03/01,2012/03/30] -> março
 - H. Qual é o tempo médio, em minutos, de processamento da atividade "Turning & Milling - Machine 4"?
5. Considere o conjunto de dados "ProducaoGrega.csv", com os dados de uma produção de cerâmica na Grécia, incluindo as informações do dia da semana, temperatura, medida do diametro e se houve defeitos ou não. Quais informações você pode extrair dos dados? Faça uma análise com o que já aprendeu.