


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

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

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

```
In [ ]: dtl["Nova coluna"] = 10
dtl
```

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

```
In [ ]: del dtl["peça1"]
dtl
```

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 [7]: dtl.sort_values(by="peça2")

Out[7]:
```

	peça1	peça2	peça3
1	2	2	3
2	3	3	4
0	1	5	2
3	4	5	3

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 parametros (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\\e-shop clothing 2008.csv"
dt = pd.read_csv(caminho, sep = ";")
```

Alguns repositórios de dados disponibilizam os mesmos 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 = ",")
```

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\\db_addre
dt_excel = pd.read_excel(caminho_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"
dt.to_csv(caminho)
```

```
In [ ]: dtl
```

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 ';', os dados ficam tabulares no Excel.
3. `encoding = "utf-8-sig"`: Permite exportar acentos.

Assim, o código melhorado fica:

```
In [37]: caminho = r"G:\\Meu Drive\\Arquivos\\UFFPR\\Disciplinas\\Arquivo_exportado.csv"
dtl.to_csv(caminho, sep = ";", index = False, encoding = "utf-8-sig")

Out[37]:
```

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 [45]: #dtl.shape
#dtl.info()
dtl.describe()
```

```
Out[45]:
```

	peça1	peça2	peça3
count	4.000000	4.00	4.000000
mean	2.500000	3.75	3.000000
std	1.290994	1.50	0.816497
min	1.000000	2.00	2.000000
25%	1.750000	2.75	2.750000
50%	2.500000	4.00	3.000000
75%	3.250000	5.00	3.250000
max	4.000000	5.00	4.000000

Exercícios II

1. Considerando 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?

1. 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 de processamento das atividades. Repnda 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*?