

Aula 4 - NumPy e Pandas I

4.1 NumPy

O NumPy (NumPy Python) é um dos pacotes básicos mais importantes para o processamento numérico em Python. Isso pelo sua estrutura de dados principal o numpy array (array), que faz o processamento de dados de forma muito mais eficiente do que listas, por exemplo.

4.1.1 Array NumPy

O array numpy é uma estrutura para armazenar dados numéricos (em sua maioria), e tem seu funcionamento como um vetor ou lista. Existem diversas formas de se criar o array. Abaixo citamos array de 3 formas distintas: usando uma lista com valores, a partir da função range() e a função np.arange(). (equivalente ao range do NumPy).

```
In [3]: import numpy as np

lista = [1,2,3,4,5] # lista normal

# Array de lista
arr1 = np.array(lista)

# Array de range
arr2 = np.array(range(10))

# Array de arange
arr3 = np.arange(10)
```

Os arrays em NumPy podem ser processados de forma vetorizada, o que aumenta a eficiência dos cálculos. Isso quer dizer que podemos realizar operações matemáticas em todos os elementos do array sem usar laços for (sempre vai existir um laço, porém ele é realizado em funções pré-compiladas em C++ ou Fortran, imbutidas no pacote NumPy). Considere um vetor de 10000 elementos representado por uma lista e por um array, que deve ter seus elementos individuais multiplicados por 2. O código abaixo faz esses cálculos e coleta o tempo de processamento de cada um, usando uma lista e um array (com a função so Notebook %time).

```
In [2]: l1 = list(range(100000))
l2 = np.arange(100000)

%time for i in range(len(l1)): l1[i] = l1[i]*2
%time l2 = l2 * 2

#time for i in range(100)
```

4.1.2 Inicialização de arrays

np.arange

Existem outras formas de inicializar arrays. Usando np.arange() cria um array com valores internos. np.arange() possui vários argumentos que podem ser utilizados, algumas construções são mostradas abaixo:

```
In [3]: # Valores entre 0 e 5
arr1 = np.arange(6,10)
arr1

# Valores entre 5 e 14
arr2 = np.arange(5,15)
arr2

# Valores entre 5 e 14 com passo de 0.5
arr3 = np.arange(5,15, 0.5)
arr3

# Valores entre -3 e 9 com passo de 0,5
arr4 = np.arange(-3, 10)
arr4

Out[3]: array([-3., -2., -1., 0., 1., 2., 3., 4., 5., 6., 7., 8., 9.])
```

np.zeros() np.ones()

Podemos ainda inicializar arrays com valores nulos ou com valores unitários usando as funções np.zeros() e np.ones().

```
In [4]: # Array com 10 elementos nulos
arr0 = np.zeros(10)
arr0

# Array com 10 elementos iguais a 1
arr0 = np.ones(10)
arr0

Out[4]: array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.])
```

np.random()

np.random: fornece diversas ferramentas para a geração de dados aleatórios em arrays. Abaixo algumas opções (extraídas de <https://numpy.org/doc/1.16/reference/routines.random.html>)

rand(d0, d1, ..., dn)	Random values in a given shape.
randn(d0, d1, ..., dn)	Return a sample (or samples) from the "standard normal" distribution.
randint(low[, high, size, dtype])	Return random integers from <i>low</i> (inclusive) to <i>high</i> (exclusive).
random_integers(low[, high, size])	Random integers of type np.int between <i>low</i> and <i>high</i> , inclusive.
random_sample(size)	Return random floats in the half-open interval [0.0, 1.0).
randin[size]	Return random floats in the half-open interval [0.0, 1.0).
randf[size]	Return random floats in the half-open interval [0.0, 1.0).
choice(a, size, replace, p)	Generates a random sample from a given 1-D array
bytes(length)	Return random bytes.

O código abaixo cria arrays de números aleatórios de diversas formas:

```
In [5]: # Mostra de 10 números aleatórios gerados pela distribuição Normal Padrão
rand_arr1 = np.random.randn(10)
rand_arr1

# Mostra de 10 números aleatórios gerados uniformemente entre 0 e 5
rand_arr2 = np.random.randint(5, size = 10)
rand_arr2

# Mostra de 10 números aleatórios gerados uniformemente entre 100 e 200
rand_arr3 = np.random.randn(100,200, size = 10)
rand_arr3

Out[5]: array([177, 188, 155, 138, 188, 199, 128, 123, 192, 199])
```

4.1.3 Arrays multidimensionais (N-dimensional array)

Arrays multidimensionais podem ser pensados como matrizes. Podemos criar arrays multidimensionais (ndarrays) das mesmas formas vistas acima, porém especificamos as suas dimensões. Abaixo alguns exemplos.

```
In [6]: # A partir de uma lista de listas
lista_lista = [[1,2,3], [4,5,6]]
nd_arr1 = np.array(lista_lista)
nd_arr1

# Matriz 2x3 de aleatórios
nd_arr2 = np.random.randn(2,3)
nd_arr2

# Matriz 2x3 de zeros - passamos uma tupla com as dimensões
nd_arr3 = np.zeros((2,3))
nd_arr3

# Matriz 2x3 de 1s - passamos uma tupla com as dimensões
nd_arr3 = np.ones((2,3))
nd_arr3

#Criando uma matriz identidade 5x5
iden = np.identity(5)
iden

Out[6]: array([[1., 0., 0., 0., 0.],
               [0., 1., 0., 0., 0.],
               [0., 0., 1., 0., 0.],
               [0., 0., 0., 1., 0.],
               [0., 0., 0., 0., 1.]])
```

Podemos verificar o tamanho dos arrays usando o método .shape(). Este método retorna uma tupla com o número de elementos referente ao número de dimensões do array, e para cada dimensão, o número representa a quantidade de elementos que existe nela. Considere o exemplo:

```
In [7]: # Matriz 2x3 de zeros - passamos uma tupla com as dimensões
nd_arr3 = np.zeros((2,3))
print(nd_arr3)

print(nd_arr3.shape)

print("Número de linhas: \n", nd_arr3.shape[0])
print("Número de colunas: \n", nd_arr3.shape[1])

[[0. 0. 0.]
 [0. 0. 0.]]
(2, 3)
Número de linhas :
2
Número de colunas :
3
```

4.1.4 Aritmética com arrays

Como dissemos, a grande vantagem de usar arrays está no processamento vetorizado, o que permite expressar operações matemáticas em lotes sem usar laços "for". Qualquer operação matemática aplicada em um array faz a operação ser aplicada a todos os seus elementos. Considere os exemplos abaixo:

```
In [8]: # Gera uma matriz 3x3 com dados aleatórios entre 2-100
arr4 = np.random.randint(0,6, size=(4,4))
print("Aleatorios :\n", arr4)

# Multiplica a linha 0 por 2:
arr4[0] = arr4[0]*2
print("Multiplica linha 0 por 2 : \n", arr4)

# Linha 0 - 1
arr4[0] = arr4[0] - 1
print("Linha 0 - 1 : \n", arr4)

# Eleva todos os elementos ao quadrado:
arr4 = arr4**2
print("Todos os elementos*2 : \n", arr4)

# Linhas:
arr4[1:] = arr4[1:] - arr4[0]
print("Linha 1 = linha 1 - linha 0 : \n", arr4)

Aleatorios :
[[ 2.5  5.]
 [ 4.2  5.]
 [ 5.2  5.]]
Multiplica linha 0 por 2 :
[[ 6. 4 10. 6]
 [ 2. 5 3.]
 [ 4. 2 4. 2.]]
Linha 0 - 1:
[[ 5. 3 9. 5.]
 [ 2. 5 3.]
 [ 4. 2 4.]]
Todos os elementos*2 :
[[10. 9 6 10.]
 [ 8 10 6 10.]
 [16 4 16 4.]
 [20 4 20 20]]
Linha 1 = linha 1 - linha 0 :
[[ 25. 9 81 25.]
 [-21. -06 -16]
 [16 4 16 4.]
 [20 4 25 25]]

Perceba-se que as operações algébricas ficam muito facilitadas com os arrays. Considere o código abaixo, que encontra a inversa da seguinte matriz:
```

```
M = [[4, 3, 3, 4],
      [4, 3, 3, 2],
      [8, 3, 5, 5],
      [5, 6, 3, 4]]

In [9]: M = np.array([[10, .3, .3, .4],[2, .3, .3, .2],[8, .3, .5, .5],[5, .6, .3, .4]])
M1 = np.identity(4)
Mprint("Inversa = np.linalg.inv(M)")
#Print "\n", M)
for i in range(M.shape[0]):
    piv = M[i,i]
    M[i] = M[i] / piv
    for j in range(M.shape[1]):
        if i != j:
            M[j] = M[j] - M[i] * M[j,i]
            M[j] = M[j] - M[i] * M[j,i]
print("Inversa : \n",M1)
```

Inversa :

[[0.28125	0.15625	-0.1875	-0.125
[0.13541667	0.26041667	-0.3125	0.125
[0.09375	0.71875	-0.0625	-0.375
[-0.625	-1.125	0.75	0.5

Por sorte, podemos conferir o resultado pelo próprio NumPy...

```
In [10]: M = np.array([[10, .3, .3, .4],[2, .3, .3, .2],[8, .3, .5, .5],[5, .6, .3, .4]])
print("Inversa pelo NumPy : \n",np.linalg.inv(M))

Inversa pelo NumPy :
[[ 0.28125 0.15625 -0.1875 -0.125]
 [ 0.13541667 0.26041667 -0.3125 0.125]
 [ 0.09375 0.71875 -0.0625 -0.375]
 [-0.625 -1.125 0.75 0.5]]
```

Uma observação importante é em relação ao tipo numérico dos arrays. Considere o seguinte caso, em que uma matriz é criada e a primeira linha substituída por ela /10.

```
In [11]: M = np.array([[10, .3, .3, .4],[2, .3, .3, .2],[8, .3, .5, .5],[5, .6, .3, .4]])
M1 = M[0]/10
print(M[0])

[[1. 0. 0. 0.]
```

O resultado não é como o esperado, pois o tipo dos dados foi inferido como inteiro. Podemos verificar o tipo de dados usando o dtype (no caso abaixo, int32).

```
In [12]: print(M.dtype)

int32

O problema pode ser corrigido ao se inicializar os valores da matriz, colocando um ponto após os números, indicando que são reais:
```

```
In [13]: M = np.array([[10, .3, .3, .4],[2, .3, .3, .2],[8, .3, .5, .5],[5, .6, .3, .4]])
M[0] = M[0]/10
print(M[0])

[[ 1. 0.3 0.3 0.4]
 float64
```

Ou ainda especificando o próprio tipo dos dados:

```
In [14]: M = np.array([[10, .3, .3, .4],[2, .3, .3, .2],[8, .3, .5, .5],[5, .6, .3, .4]], dtype=np.float64)
M[0] = M[0]/10
print(M[0])

[[ 1. 0.3 0.3 0.4]
 float64
```

4.1.5 Fatiamento de arrays

O fatiamento de arrays permite visualizar partes do mesmo. Para arrays unidimensionais a sintaxe é muito parecida com o fatiamento de listas. Considere os exemplos abaixo.

```
In [15]: # Gera 10 valores extraídos da normal padrão
arr = np.random.randn(10)
print(arr)

# Imprime os 5 primeiros valores (de 0 a 4)
print(arr[:5])

# Imprime os últimos valores, a partir do índice 5
print(arr[5:])

# Imprime os elementos de índices 2-5
print(arr[2:6])

[-0.74513883 -1.13877339 -0.04052013 0.48931399 -0.50095004 -0.17217431
 0.45768785 2.3528207 -0.02000582 0.2703764 ]
[-0.74513883 -1.13877339 -0.04052013 0.48931399 -0.50095004
 0.47174331 0.45768785 2.3528207 -0.02000582 0.2703764 ]
[-0.04052013 0.48931399 -0.50095004 -0.17217431]
```

Uma diferença importante entre o fatiamento de listas e de arrays, é que estes últimos são visualizações (views) do próprio array, ou seja, alterando a matriz também altera o array. Considere o exemplo:

```
In [16]: arr = np.zeros(10, dtype = np.float64)
print(arr)

arr[5] = 10
print("Alterando os valores por fatiamento :",arr)

[0. 0. 0. 0. 0. 10. 0. 0. 0. 0.]
Alterando os valores por fatiamento : [10. 10. 10. 10. 10. 0. 0. 0. 0. 0.]

Se quisermos uma cópia do fatiamento precisamos dizer explicitamente, usando o método .copy()
```

```
In [17]: arr = np.zeros(10, dtype = np.float64)
print(arr)

copla = arr[5:].copy()
copla = 10
print("Copiando não altera os valores :",arr)

[0. 0. 0. 0. 0. 0. 0. 0. 0. 0.]
Copiando não altera os valores : [0. 0. 0. 0. 0. 0. 0. 0. 0. 0.]

Em arrays multidimensionais os fatiamentos de cada índice não são mais escalares, mas arrays unidimensionais. Considere o caso 2d:
```

```
In [18]: M = np.array([[10, .3, .3, .4],[2, .3, .3, .2],[8, .3, .5, .5],[5, .6, .3, .4]])
print("Matriz original : \n", M)

# Imprime todas as linhas a partir do índice 1
print("Linhas a partir do índice 1 :\n",M[1:])

# De todas as linhas a partir do índice 1 (igual anterior), seleciona as colunas até o índice 2
print("Colunas até o índice 2, das linhas a partir do índice 1 :\n",M[1:,1:3])

Matriz original :
[[10. 3. 3. 4.]
 [ 2. 3. 3. 2.]
 [ 8. 3. 5. 5.]
 [ 5. 6. 3. 4.]]
Linhas a partir das índice 1 :
[[ 2. 3. 3. 2.]
 [ 8. 3. 5. 5.]
 [ 5. 6. 3. 4.]]
Colunas até o índice 2, das linhas a partir do índice 1 :
[[2. 3. 3.]
 [ 8. 3. 5.]
 [ 5. 6. 3.]]
```

4.1.6 Indexação booleana

Também podemos realizar operações booleanas em arrays, de forma que o resultado será um novo array de valores booleanos, de acordo com a condição. Considere o exemplo:

```
In [19]: arr_string = np.array(["Dwight", "Michael", "Angela", "Oscar", "Michael", "Angela"])
arr_condico = quais elementos do array são iguais a "Michael"?
arr_bool = arr_string == "Michael"
print(arr_bool)

# Condico : quais elementos do array são iguais a "Michael" OU "Angela"
arr_bool = arr_string == "Michael" | arr_string == "Angela"
print(arr_bool)

[False True False False True False]
[False True False True True]

Também podemos fazer o processo reverso: passamos um array de booleanos para um array, e ele retorna somente os elementos (ou arrays) em que a condição é verdadeira:
```

```
In [20]: arr_string = np.array(["Dwight", "Michael", "Angela", "Oscar", "Michael", "Angela"])
arr_booleano = np.array([True,False,False,True,False,False])
arr_booleano

# Seleciona somente os elementos em que arr_booleano == True
print(arr_string[arr_booleano])

['Dwight' 'Oscar']

Também podemos fazer a indexação booleana em arrays multidimensionais. Nesses casos as condições verdadeiras retornam arrays de dimensões menores. Considere o seguinte caso:
```

```
In [21]: # Gerando uma matriz 3x4 de aleatórios entre 5 e 9
ndarray = np.random.randint(5,10, size=(3,4))
ndarray

# Gerando um array de booleanos com o mesmo número de elementos da primeira dimensão da Matriz (3)
arr_bool = np.array([True,False,False])

# Imprimindo somente as linhas de ndarray que satisfizerem as condições de arr_bool
print(ndarray[arr_bool])

[[8 6 9 8]]

Combinando as duas indexações nos fornece uma poderosa ferramenta para a análise de dados. Considere o seguinte cenário: temos os dados de produção de uma indústria de pães, em que a cada vez que um lote é produzido, uma amostra de 5 pães é verificada pela qualidade, aferindo o peso total. Os tipos de pães são armazenados em um array chamado arr_paes: e as coletas dos pesos em uma ndarray chamado arr_pesos. Os valores são os seguintes: arr_paes = np.array(["frances","italiano","sirio","frances","sirio"])
```

```
arr_pesos = np.array([3.0,2.8,3.1,3.0,3.23],
                     [5.0,5.3,4.95,4.9,5.23],
                     [3.0,2.8,3.1,3.0,3.23],
                     [6.0,6.8,6.1,6.0,6.23],
                     [3.0,2.8,3.1,3.0,3.23]])
```

Podemos realizar filtros na matriz de pesos com base nos pães que desejamos verificar. Considere os exemplos:

```
In [22]: arr_paes = np.array(["frances","italiano","sirio","frances","sirio"])
arr_pesos = np.array([[3.0,2.8,3.1,3.0,3.23],
                     [5.0,5.3,4.95,4.9,5.23],
                     [3.0,2.8,3.1,3.0,3.23],
                     [6.0,6.8,6.1,6.0,6.23],
                     [3.0,2.8,3.1,3.0,3.23]])

# Filtrando todas as linhas que contém medidas do pão frances
arr_frances = arr_pesos[arr_paes == "frances"]
print("Linhas pao frances \n", arr_frances)

# Filtrando todas as linhas que contém medidas do pão sirio
arr_frances = arr_pesos[arr_paes == "sirio"] | arr_paes == "frances")
print("Linhas pao sirio ou frances \n", arr_frances)

# Filtrando todas as linhas que contém medidas do pão sirio OU frances
arr_frances = arr_pesos[(arr_paes == "sirio") | (arr_paes == "frances")]
print("Linhas pao sirio ou frances \n", arr_frances)
```

Linhas pao frances

[3. 2.8 3.1 3. 3.23]
[6. 6.8 6.1 6. 6.23]

Linhas pao sirio

[3. 2.8 3.1 3. 3.23]
[3. 2.8 3.1 3. 3.23]

Linhas pao sirio ou frances

[3. 2.8 3.1 3. 3.23]
[3. 2.8 3.1 3. 3.23]
[5. 5.3 4.95 4.9 5.23]
[6. 6.8 6.1 6. 6.23]
[3. 2.8 3.1 3. 3.23]

Note que a indexação booleana, diferentemente do fatiamento, não produz uma view do array, mas sim uma cópia! Ou seja, alterar o resultado de uma indexação booleana não altera os valores originais. Considere o exemplo abaixo:

```
In [23]: arr_paes = np.array(["frances","italiano","sirio","frances","sirio"])
arr_pesos = np.array([[3.0,2.8,3.1,3.0,3.23],
                     [5.0,5.3,4.95,4.9,5.23],
                     [3.0,2.8,3.1,3.0,3.23],
                     [6.0,6.8,6.1,6.0,6.23],
                     [3.0,2.8,3.1,3.0,3.23]])

arr_frances = arr_pesos[arr_paes == "frances"]
print(arr_frances)

arr_frances[0] = 99
print("Alterando arr_frances \n",arr_frances)

print("Não altera arr_pesos \n",arr_pesos)

[[3. 2.8 3.1 3. 3.23]
 [6. 6.8 6.1 6. 6.23]]
Alterando arr_frances
[[99. 99. 99. 99. 99. ]
 [ 6. 6.8 6.1 6. 6.23]]
Não altera arr_pesos
[[3. 2.8 3.1 3. 3.23]
 [5. 5.3 4.95 4.9 5.23]
 [3. 2.8 3.1 3. 3.23]
 [6. 6.8 6.1 6. 6.23]
 [3. 2.8 3.1 3. 3.23]]
```

4.1.7 Métodos matemáticos e estatísticos

Os arrays do NumPy possuem muitos métodos que matemáticos que facilitam o processamento. Alguns deles são: sum() mean() std() var() cumsum() min() max() argmin() argmax()

```
In [24]: # Gera 20 elementos aleatórios (entre 10 e 19)
arr_rand = np.random.randint(10,20, size=(20))
print("Valores : \n", arr_rand)

# calcula a soma
print("Soma : \n", arr_rand.sum())

# calcula a média
print("Media : \n", arr_rand.mean())

# calcula o desvio padrão
print("Desvio padrão : \n", arr_rand.std())

# calcula a variância
print("Variancia : \n", arr_rand.var())

# Máximo
print("Máximo : \n",arr_rand.max())

# Índice do Máximo
print("Índice do Máximo : \n",arr_rand.argmax())

# Soma cumulativa dos elementos começando em 0
print("Soma cumulativa : \n", arr_rand.cumsum())

Valores :
[14 17 15 12 14 19 12 18 11 15 19 16 19 13 11 19 10 13 17 19]
Soma :
300
Média :
15.15
Desvio padrão :
3.021175268004159
Variância :
9.127500000000001
Máximo :
19
Índice do Máximo :
5
Soma cumulativa :
[ 14 31 46 58 72 91 103 121 132 147 166 182 201 214 225 244 254 267 284 303]
```

Em arrays multidimensionais podemos escolher em relação a qual eixo que desejamos coletar as informações (não todas):

```
In [25]: arr_m = np.array([[1,1,1,1],
                       [4,5,6,6]],
                       [[1,1,1,1],
                       [2,2,2,2]])

print("Média por colunas", arr_m.mean(axis=0))
print("Média por linhas", arr_m.mean(axis=1))

print("Maior elemento", arr_m.max())

Média por colunas [5. 3.33333333 3.33333333 3. ]
Média por linhas [1. 5.25 4.75]
Maior elemento 10
```

Exercícios I

1. Escreva os seguintes vetores como arrays numpy.

v1 = [10,20,30,20,10,1,0,2,5,0,20,1,4,0,20,20,30,40,13,44,55]

v2 = [1,25,50,41,5,20,10,23,5,10,20,13,4,20,100,20,50,35,40,4,55,55]

2. Considere as seguintes sequências matemáticas, e para cada uma delas escreva um algoritmo que armazene os elementos em um array, e calcule a soma e o desvio padrão dos valores.

A. $(n), n = 1, \dots, 100$

B. $\left\{ \frac{n}{n+1} \right\}, n = 1, \dots, 100 = \left\{ \frac{1}{2}, \frac{2}{3}, \dots \right\}$

C. $\left\{ \frac{(-1)^n (n+1)}{3^n} \right\}, n = 1, \dots, 100 = \left\{ -\frac{2}{3}, \frac{4}{9}, \frac{1}{27}, \dots \right\}$

3. Crie um array arrN20 com 20 dados aleatórios extraídos da distribuição Normal padrão.

4. Crie um array arrU20 com 20 dados aleatórios extraídos de uma distribuição uniforme, com valores entre -10 e 10.

5. Imprima a multiplicação de arrN20 por 10.

6. Imprima a multiplicação de arrN20 por arrU20, esse é o resultado esperado de uma multiplicação vetorial?

7. Gere uma ndarray MW 5x10 com dados extraídos de uma Uniforme(-10,60).

8. Gere uma ndarray MW 5x10 com valores somente os valores pares.

9. Considerando o array arrN20, imprima os elementos entre os índices 5 e 10 incluindo o 10 (usando fatiamento).

11. Considerando o array arrU20, imprima todos os valores, exceto o último (usando fatiamento).

12. Considerando o array arrU20, imprima todos os valores, exceto o primeiro (usando fatiamento).

13. Considere o seguinte ndarray:

```
M = np.array([[4, 3, 3, 4],
              [4, 3, 3, 2],
              [5, 3, 5, 5],
              [5, 3, 3, 4]])
```

Use fatiamento para imprimir os números do array, de acordo com a imagem abaixo:

1. Ainda considerando o ndarray, do exemplo anterior, encontre A. O soma de os elementos por linha. B. O array com a soma das linhas em cada coluna. C. A soma e a média de todos os elementos.

1. Resolva os sistemas de equações lineares abaixo usando NumPy

<https://numpy.org/doc/stable/reference/generated/numpy.linalg.solve.html>:

A. $2x + 3y = 4$

B. $x - 5y = 2$

C. $x + y = 0$

D. $2x + y + z = 0$

E. $4x + 3y + z = 0$

2. Considere os seguintes dados de coleta de amostras de pesos de pães (como no exemplo):

```
arr_paes = np.array(["frances","italiano","sirio","frances","sirio"])
arr_pesos = np.array([[3.0,2.8,3.0,3.0,3.23,3.0,2.8,3.0],
                     [5.0,5.3,4.95,4.9,5.23,5.5,5.6],
                     [3.0,2.8,3.1,3.0,3.23,3.3,3.3,3.1],
                     [6.0,6.8,6.1,6.0,6.23,5.8,5.9,6.1],
                     [3.0,2.8,3.1,3.0,3.23,3.3,3.1,3.1]])
```

O controle de qualidade define que, se uma amostra tem variância maior do que metade da média, existe algo errado com os dados (muita variabilidade), que precisa de uma amostra deve ser coletada novamente. Crie um código (usando indexação booleana) que retorne o pão se existir algum) de forma que uma nova amostra seja coletada.

1. Ainda considerando os dados dos pães. A qualidade precisa saber a média dos pesos de todos os pães no primeiro e no último dia de coletas. Use fatiamentos e o máximo duas linhas para extrair as duas informações.

4.2 Pandas I

O pandas é um pacote essencial para se realizar análise de dados, muito disso se dá pelas suas duas estruturas de dados principais, a series e o Dataframe, usados em quase todas as aplicações de mineração de dados. Utilizaremos a importação do pacote com a seguinte convenção:

```
In [4]: import pandas as pd
```

4.2.1 Series

Uma serie é um objeto do tipo array unidimensional contendo uma sequência de valores (de tipos semelhantes aos do NumPy) e um array associado de rótulos (labels) de dados, chamado index. A series mais simples é composta de um array de dados:

```
In [27]: ser1 = pd.Series([4,3,4,5])
print(ser1)
print(type(ser1))

0    4
1    3
2    4
3    5
dtype: int64
<class 'pandas.core.series.Series'>
```

Podemos acessar tanto os valores quanto os índices de uma Series pelos métodos .values e .index:

```
In [28]: print(ser1.values)
print(ser1.index)

[4 3 4 5]
RangeIndex(start=0, stop=4, step=1)
```

Note que o tipo de estrutura de dados do values é justamente um array NumPy:

```
In [29]: print(type(ser1.values))

print('numpy.ndarray')

<class 'numpy.ndarray'>
```

Podemos criar uma Series e alterar os valores de index para o que quisermos, considere:

```
In [30]: ser2 = pd.Series([1,2,3,4], index=["a","b","c","d"])
print(ser2)

a    1
b    2
c    3
d    4
dtype: int64
```

Podemos usar os valores dos índices para acessar e alterar os elementos:

```
In [31]: print(ser2["a"])

# Alterando o elemento
ser2["a"] = 999
print(ser2.values)

1
999 2 3 4
```

4.2.2 Dataframe

Um dataframe representa uma tabela. O dataframe contém uma coleção ordenada de dados, em que cada uma é uma Series e pode ter um tipo de dado diferente. O dataframe tem um índice tanto para as linhas quanto para as colunas. Existem diversas formas para se criar Dataframes (embora na maioria dos casos ele será criado automaticamente ao carregarmos dados externos), algumas delas são mostradas abaixo:

Criação de Dataframes

```
In [5]: # Dataframe a partir de um dicionário de listas; as chaves são os nomes das colunas e as listas os valores
dic1 = {"peca1": [1,2,3,4],
        "peca2": [5,2,3,5],
        "peca3": [2,3,4,3]}
df1 = pd.DataFrame(dic1)
df1
```

peca1 peca2 peca3

0	1	5	2
1	2	2	3
2	3	3	5
3	4	3	3

