

Relatório 4 - Prática: Principais Bibliotecas e Ferramentas Python para Aprendizado de Máquina (I)

João Pedro Gomes

1. Introdução

O objetivo do card 4, foi estudar sobre, Jupyter Notebook, Python para análise com numpy e pandas e realizar exercícios que estavam no docs que tinha as mesmas aulas e também realizar a prática dos dois temas e fazer um relatório.

2. Desenvolvimento

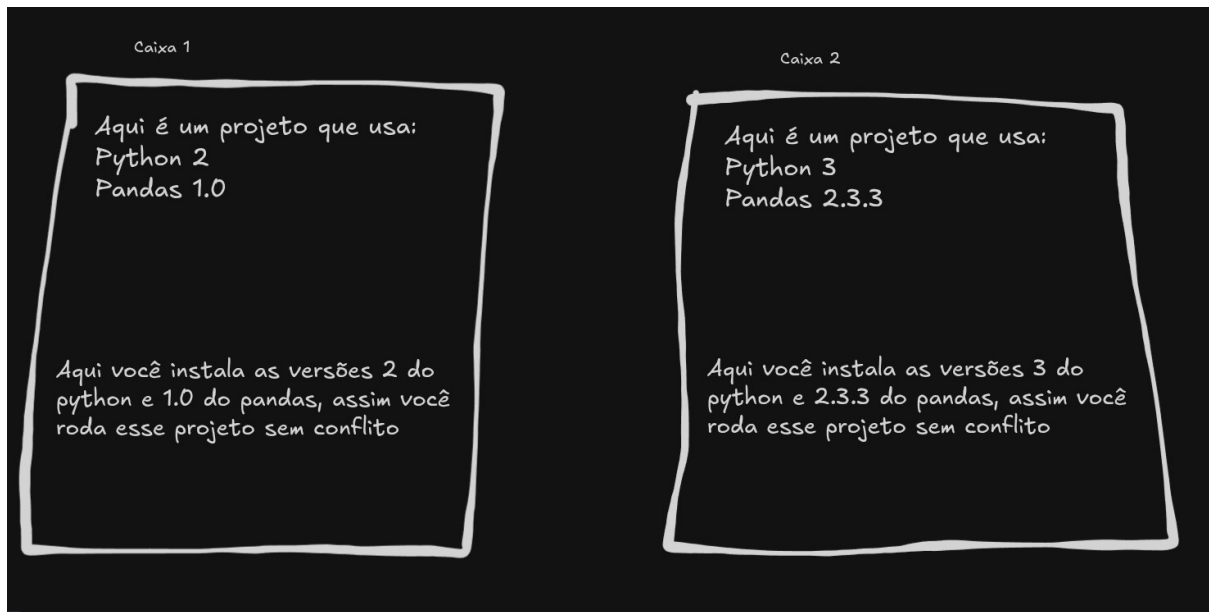
O card é dividido em 3 seções para estudo, Jupyter Notebook, Numpy e Pandas, eu vou falar sobre as 3 e depois explicar a minha prática com o pandas.

Jupyter Notebook:

O jupyter notebook é uma ferramenta que você baixa no pc pra rodar código python através do navegador, ele funciona com um esquema de célula ao invés de ter várias linhas igual em um vscode, você tem uma célula que pode ter vários códigos dentro dela e depois rodar ela sozinha podendo ter várias células com diferentes códigos em cada uma, ao rodar ele exibe o resultado do código abaixo da célula, indo por partes por assim dizer.

Ainda na seção do jupyter ele fala sobre ambientes virtuais, é tipo como se fosse uma caixa separada onde fica cada projeto com suas versões de bibliotecas, python e etc..., exemplo você tem um projeto que usa pandas 1.0, e no outro pandas 2.3.3, aí na caixa que usa pandas 1.0 você instala o pandas 1.0 e na outra o pandas 2.3.3, assim não dá conflito de versão.

Um exemplo visual:

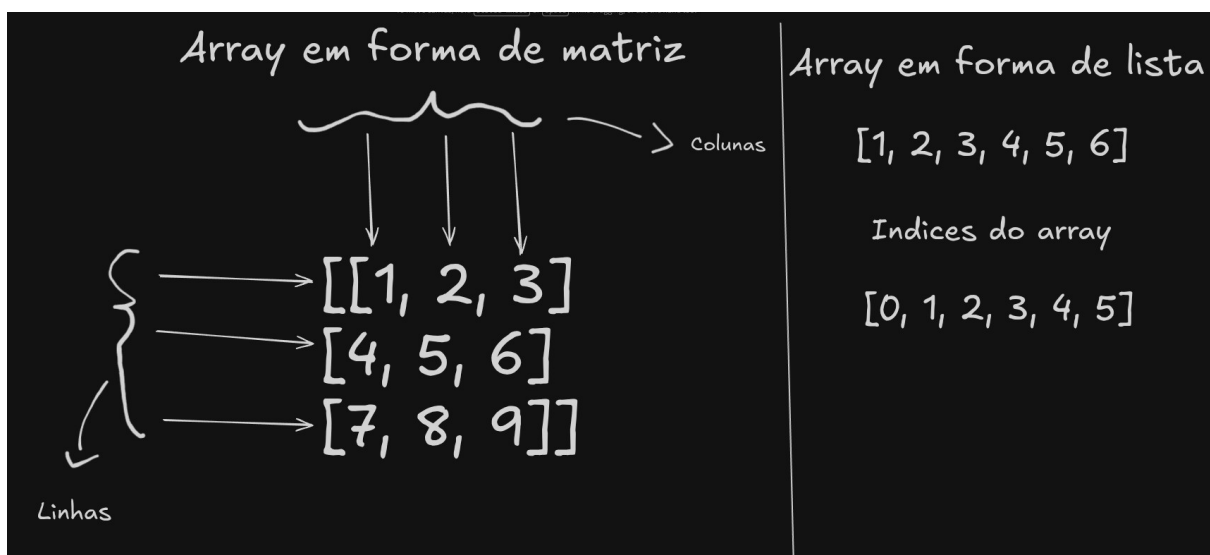


Numpy:

É a biblioteca de algebra linear do python, ela é muito rápida porque ela é escrita em C, com ela a gente mexe com arrays e ela tem funções pra manipular os dados desses arrays.

Um array vem de duas formas, em forma de vetor ou lista, e em matriz.

Um exemplo visual:



O array não tem limite de tamanho, ele pode ser até mesmo uma representação de um objeto 3D como um cubo. Pra criar um array usando o numpy é só importar o

numpy e dar um apelido geralmente é “np” e existem várias funções e formas de criar vou passar por algumas.

Se você cria uma lista tem como usar o np.array(e passar sua lista aqui como parametro e ela vira um array numpy), isso funciona se você tiver mais de uma lista ai ela vira uma matriz usando o mesmo modo.

```
np.arange(0, 10)
Saída: [0, 1, 2, 3, 4, 5, 6, 7, 8, 9]

np.arange(0, 10, 2)
Saída: [0, 2, 4, 6, 8]

np.zeros(3)
Saída: [0, 0, 0]

np.ones((3, 3))
Saída: [[1., 1., 1.]
        [1., 1., 1.]
        [1., 1., 1.]]

np.eye(4)
Saída: [[1, 0, 0, 0]
        [0, 1, 0, 0]
        [0, 0, 1, 0]
        [0, 0, 0, 1]]

np.linspace(0, 10, 3)
Saída: [0, 5, 10]
```

Esses são alguns dos principais, o arange cria um array de 10 indices indo de 0 a 9, o terceiro parametro que é colocado no arange é o espaçamento se você quiser, no exemplo ele vai somando 2 em 2, o np.zeros() cria um array apenas com zeros do tamanho que quiser, mas se você coloca 4 paranteses e dentro coloca 2 valores ele joga pra uma matriz no exemplo eu usei a função ones pra criar uma matriz 3x3 só de uns, a função eye cria a matriz identidade no exemplo eu fiz uma 4x4 por último o linspace faz um array de 0 a 10 e pega 3 números igualmente espaçados.

Existem as funções random que podem ser usado pra criar arrays, existe a random.rand(x) que pega a quantidade x de números aleatorios entre 0 e 1, a random.randn(x) faz a mesma coisa que a rand só que pega qualquer número aleatoriamente e não só de 0 a 1 e a random.randint(onde começa, onde para, quantidade de número que ira gerar) ex: random.randint(10, 60, 2) pega 2 numeros entre 10 e 60.

E para manipular os arrays existem:

O reshape que faz um array em forma de linha virar uma matriz e somente shape pega o formato da matriz exemplo: (5,4)

Tem funções que retornam os maiores e menores valores e também funções que dão o número do indice de onde esses valores estão. Existem as operações básicas como soma, subtração, divisão e multiplicação, as que pegam a raiz quadrada de cada valor, os valores exponenciais de cada valor, calcula a media do array todo, o desvio padrão e o seno de cada valor.

Também tem como fatiar os arrays da seguinte forma:

Mat Arr

[[1, 2, 3] [4, 5, 6] [7, 8, 9]]	[1, 2, 3, 4, 5, 6]
---------------------------------------	--------------------

→ Nossos arrays

arr[2:5] → Pega do índice 2 a 5, mas deixa o índice 5 de fora
Saída: [3, 4, 5]

arr[2:] → Pega do índice 2 adiante
Saída: [3, 4, 5, 6]

Mat[:2][:] → Pega até a 2 linha da matriz e todas as colunas
Saída: [[1, 2, 3]
[4, 5, 6]]

Pandas:

O pandas é outra biblioteca do python mais voltada pra análise de dados, ela tem estruturas como Series, Dataframes e outras funcionalidades que irei falar sobre.

Uma serie é uma estrutura muito parecida com um dicionário em python, ela tem chave – valor, onde a chave seria o index dessa estrutura pra você poder acessar os valores.

Exemplo visual:

Indice Valor

Serie =

Taylor Swift	10
The Beatles	8
Power Rangers	10
John Lennon	10

Series['Taylor Swift'] → Acessa o valor que o indice Taylor Swift tem
Saída: 10

Existe as operações normais com as series, um outro exemplo:

Serie =

Taylor Swift	10
The Beatles	8
Power Rangers	10
John Lennon	10

SerieDois =

Taylor Swift	10
The Beatles	8
Power Rangers	10
John Lennon	10
Kamen Rider	8,5

Se a gente faz a soma de Serie + SerieDois Resulta nessa serie:

Taylor Swift	20
The Beatles	16
Power Rangers	20
John Lennon	20
Kamen Rider	NaN

Pois ele nao achou um indice igual a Kamen Rider na primeira serie, logo nao deu pra fazer a soma

DataFrames:

Um dataframe é uma estrutura que é formada por várias séries, ela já é uma tabela muito mais completa, podendo ter colunas onde cada coluna é uma series.

Exemplo visual:

Df =

	W	Z
A	23	21
B	13	11
C	53	41

Um DataFrame normal

Df['W']

Saída:

	W
A	23
B	13
C	53

Pega apenas a coluna W

E tem como criar colunas novas, apenas fingindo que elas já existem fazendo assim definindo o nome da coluna e os dados. E pra deletar ela só usar o comando drop e o nome da coluna.

Pra fazer buscas por algum valor exato se usa o comando loc e passa a linha e a coluna da onde o valor está e usando o iloc você passa o número da linha e coluna.

Existe também uma forma de setar um index novo com alguma das colunas, exemplo:

To move canvas, hold `Scroll wheel` or `Space` while dragging, or use the hand tool

Este é o nosso DataFrame

	index	W	X	Y	Z	Estado
0	A	2.706850	0.628133	0.907969	0.503826	RS
1	B	0.651118	-0.319318	-0.848077	0.605965	RJ
2	C	-2.018168	0.740122	0.528813	-0.589001	SP
3	D	0.188695	-0.758872	-0.933237	0.955057	AM
4	E	0.190794	1.978757	2.605967	0.683509	SC

Usando o comando `set_index()` na coluna Estado

	index	W	X	Y	Z
Estado					
RS	A	2.706850	0.628133	0.907969	0.503826
RJ	B	0.651118	-0.319318	-0.848077	0.605965
SP	C	-2.018168	0.740122	0.528813	-0.589001
AM	D	0.188695	-0.758872	-0.933237	0.955057
SC	E	0.190794	1.978757	2.605967	0.683509

E ainda existe uma forma de colocar dois indexes em um dataframe só, podendo ter duas formas de acessar valores usando a função `xs` no `df`.

Uma forma de tratar dados que o curso mostrou é de preencher valores `NaN`, usando funções `fillna`, `dropna` que excluem ou aplicam um valor novo.

Também dá pra juntar dataframes de várias formas, concatenando que junta eles mas sem somar os dados, a mescla que junta os dataframes por algum objeto que eles tem iguais como parametro, e os métodos são iguais ao do SQL e a junção usando o índice como a referência pra mescla e usa o `left` como padrão, é tipo um `left join` no `sql`.

As operações com os dataframes são bem parecidas com as da series, mas normalmente você escolhe se vai ser de coluna a coluna ou linha por linha.

O pandas tem a opção de ler e exportar arquivos de `.csv` ou arquivo excel pra dataframe e vice e versa, também é possível ler páginas html e jogar elas pra um dataframe.

3. Conclusão

Em resumo, fazendo o card 4 deu pra aprender como usar o jupyter notebook, o numpy e o pandas, que são ferramentas importantes pra análise de dados em python. O aprendizado foi rápido e eficaz e a utilização do jupyter notebook foi essencial pra organização do código, o numpy foi muito bom para entender como mexer com arrays e matrizes e o pandas ajudou na organização de tabelas e dados estruturados.