

Eléments De réponses

Sghuri Chaimaa

Introduction

- Ces slides accompagnent le code réponse pour le test Technique qui concerne :Traitement d'une base de donnée 'Market_data', il présente un document explicatif du code et des réponses
- Cette base de donnée expose les données financières de 6 entreprises :MSFT, IBM, AAPL, FB, GOOG et NDX
- Pour chacune des entreprises citées ci-haut, on peut tirer les informations suivantes :
 - Prix de l'action (stock price)
 - Id (de 0 à 1646)
 - Date



Step 01

- Importer les données et
- Organiser ces données sous la forme souhaitée afin de faciliter l'extraction

01

- Définir une fonction qui retourne les prix des stocks spécifiés sinon par défaut tous les stocks , pendant une période déterminée

02

- Analyse des données manquantes,
- Gérer les Datamissing

03

TABLE OF CONTENTS

04

- Visualisation des données

05

Calculer les daily Returns

06

Calculer la moyenne et l'écart type de ces daily returns

01

**Importer les données et
Organiser ces données sous la forme
souhaitée afin de faciliter l'extraction**



Importer les données:

- On a la base de donnée est sous forme un fichier json , Pour Faciliter l'import des données , On va utiliser la librairie Pandas .
- On a une dataframe dont les données ont la forme d'un dictionnaire d'ou la nécessité de reorganiser cette dataframe afin de faciliter l'extraction
- Data Initiale :

```
Entrée [28]: market_data=pd.read_json('market_data.json')
market_data
```

Out[28]:

	MSFT	IBM	AAPL	FB	GOOG	NDX
0	{'_id': 0, 'data': {'Date': '2012-05-18T00:00:00.000Z', 'Open': 28.5, 'High': 29.0, 'Low': 28.0, 'Close': 28.75, 'Volume': 150000000}}	{'_id': 0, 'data': {'Date': '2012-05-18T00:00:00.000Z', 'Open': 160.0, 'High': 165.0, 'Low': 155.0, 'Close': 162.5, 'Volume': 100000000}}	{'_id': 0, 'data': {'Date': '2012-05-18T00:00:00.000Z', 'Open': 110.0, 'High': 115.0, 'Low': 105.0, 'Close': 112.5, 'Volume': 80000000}}	{'_id': 0, 'data': {'Date': '2012-05-18T00:00:00.000Z', 'Open': 75.0, 'High': 80.0, 'Low': 70.0, 'Close': 77.5, 'Volume': 60000000}}	{'_id': 0, 'data': {'Date': '2012-05-18T00:00:00.000Z', 'Open': 550.0, 'High': 560.0, 'Low': 540.0, 'Close': 555.0, 'Volume': 40000000}}	{'_id': 0, 'data': {'Date': '2012-05-18T00:00:00.000Z', 'Open': 2800.0, 'High': 2850.0, 'Low': 2750.0, 'Close': 2820.0, 'Volume': 20000000000}}
1	{'_id': 1, 'data': {'Date': '2012-05-21T00:00:00.000Z', 'Open': 29.0, 'High': 29.5, 'Low': 28.5, 'Close': 29.25, 'Volume': 160000000}}	{'_id': 1, 'data': {'Date': '2012-05-21T00:00:00.000Z', 'Open': 162.0, 'High': 167.0, 'Low': 157.0, 'Close': 164.5, 'Volume': 110000000}}	{'_id': 1, 'data': {'Date': '2012-05-21T00:00:00.000Z', 'Open': 112.0, 'High': 117.0, 'Low': 107.0, 'Close': 114.5, 'Volume': 85000000}}	{'_id': 1, 'data': {'Date': '2012-05-21T00:00:00.000Z', 'Open': 77.0, 'High': 82.0, 'Low': 72.0, 'Close': 79.5, 'Volume': 65000000}}	{'_id': 1, 'data': {'Date': '2012-05-21T00:00:00.000Z', 'Open': 560.0, 'High': 570.0, 'Low': 550.0, 'Close': 565.0, 'Volume': 45000000}}	{'_id': 1, 'data': {'Date': '2012-05-21T00:00:00.000Z', 'Open': 2850.0, 'High': 2900.0, 'Low': 2800.0, 'Close': 2870.0, 'Volume': 21000000000}}
2	{'_id': 2, 'data': {'Date': '2012-05-22T00:00:00.000Z', 'Open': 28.75, 'High': 29.25, 'Low': 28.25, 'Close': 29.0, 'Volume': 155000000}}	{'_id': 2, 'data': {'Date': '2012-05-22T00:00:00.000Z', 'Open': 164.0, 'High': 169.0, 'Low': 159.0, 'Close': 166.5, 'Volume': 105000000}}	{'_id': 2, 'data': {'Date': '2012-05-22T00:00:00.000Z', 'Open': 114.0, 'High': 119.0, 'Low': 109.0, 'Close': 116.5, 'Volume': 82000000}}	{'_id': 2, 'data': {'Date': '2012-05-22T00:00:00.000Z', 'Open': 79.0, 'High': 84.0, 'Low': 74.0, 'Close': 81.5, 'Volume': 62000000}}	{'_id': 2, 'data': {'Date': '2012-05-22T00:00:00.000Z', 'Open': 565.0, 'High': 575.0, 'Low': 555.0, 'Close': 570.0, 'Volume': 42000000}}	{'_id': 2, 'data': {'Date': '2012-05-22T00:00:00.000Z', 'Open': 2870.0, 'High': 2920.0, 'Low': 2820.0, 'Close': 2890.0, 'Volume': 20500000000}}
3	{'_id': 3, 'data': {'Date': '2012-05-23T00:00:00.000Z', 'Open': 29.25, 'High': 29.75, 'Low': 28.75, 'Close': 29.5, 'Volume': 165000000}}	{'_id': 3, 'data': {'Date': '2012-05-23T00:00:00.000Z', 'Open': 166.0, 'High': 171.0, 'Low': 161.0, 'Close': 168.5, 'Volume': 115000000}}	{'_id': 3, 'data': {'Date': '2012-05-23T00:00:00.000Z', 'Open': 116.0, 'High': 121.0, 'Low': 111.0, 'Close': 118.5, 'Volume': 88000000}}	{'_id': 3, 'data': {'Date': '2012-05-23T00:00:00.000Z', 'Open': 81.0, 'High': 86.0, 'Low': 76.0, 'Close': 83.5, 'Volume': 68000000}}	{'_id': 3, 'data': {'Date': '2012-05-23T00:00:00.000Z', 'Open': 570.0, 'High': 580.0, 'Low': 560.0, 'Close': 575.0, 'Volume': 48000000}}	{'_id': 3, 'data': {'Date': '2012-05-23T00:00:00.000Z', 'Open': 2890.0, 'High': 2940.0, 'Low': 2840.0, 'Close': 2910.0, 'Volume': 21500000000}}
4	{'_id': 4, 'data': {'Date': '2012-05-24T00:00:00.000Z', 'Open': 29.5, 'High': 30.0, 'Low': 29.0, 'Close': 29.75, 'Volume': 170000000}}	{'_id': 4, 'data': {'Date': '2012-05-24T00:00:00.000Z', 'Open': 168.0, 'High': 173.0, 'Low': 163.0, 'Close': 170.5, 'Volume': 120000000}}	{'_id': 4, 'data': {'Date': '2012-05-24T00:00:00.000Z', 'Open': 118.0, 'High': 123.0, 'Low': 113.0, 'Close': 120.5, 'Volume': 90000000}}	{'_id': 4, 'data': {'Date': '2012-05-24T00:00:00.000Z', 'Open': 83.0, 'High': 88.0, 'Low': 78.0, 'Close': 85.5, 'Volume': 70000000}}	{'_id': 4, 'data': {'Date': '2012-05-24T00:00:00.000Z', 'Open': 575.0, 'High': 585.0, 'Low': 565.0, 'Close': 580.0, 'Volume': 50000000}}	{'_id': 4, 'data': {'Date': '2012-05-24T00:00:00.000Z', 'Open': 2910.0, 'High': 2960.0, 'Low': 2860.0, 'Close': 2930.0, 'Volume': 22000000000}}
...
1642	{'_id': 1642, 'data': {'Date': '2018-12-19T00:00:00.000Z', 'Open': 110.0, 'High': 115.0, 'Low': 105.0, 'Close': 112.5, 'Volume': 80000000}}	{'_id': 1642, 'data': {'Date': '2018-12-19T00:00:00.000Z', 'Open': 160.0, 'High': 165.0, 'Low': 155.0, 'Close': 162.5, 'Volume': 100000000}}	{'_id': 1642, 'data': {'Date': '2018-12-19T00:00:00.000Z', 'Open': 110.0, 'High': 115.0, 'Low': 105.0, 'Close': 112.5, 'Volume': 80000000}}	{'_id': 1642, 'data': {'Date': '2018-12-19T00:00:00.000Z', 'Open': 75.0, 'High': 80.0, 'Low': 70.0, 'Close': 77.5, 'Volume': 60000000}}	{'_id': 1642, 'data': {'Date': '2018-12-19T00:00:00.000Z', 'Open': 550.0, 'High': 560.0, 'Low': 540.0, 'Close': 555.0, 'Volume': 40000000}}	{'_id': 1642, 'data': {'Date': '2018-12-19T00:00:00.000Z', 'Open': 2800.0, 'High': 2850.0, 'Low': 2750.0, 'Close': 2820.0, 'Volume': 20000000000}}
1643	{'_id': 1643, 'data': {'Date': '2018-12-24T00:00:00.000Z', 'Open': 112.0, 'High': 117.0, 'Low': 107.0, 'Close': 114.5, 'Volume': 85000000}}	{'_id': 1643, 'data': {'Date': '2018-12-24T00:00:00.000Z', 'Open': 162.0, 'High': 167.0, 'Low': 157.0, 'Close': 164.5, 'Volume': 110000000}}	{'_id': 1643, 'data': {'Date': '2018-12-24T00:00:00.000Z', 'Open': 112.0, 'High': 117.0, 'Low': 107.0, 'Close': 114.5, 'Volume': 85000000}}	{'_id': 1643, 'data': {'Date': '2018-12-24T00:00:00.000Z', 'Open': 77.0, 'High': 82.0, 'Low': 72.0, 'Close': 79.5, 'Volume': 65000000}}	{'_id': 1643, 'data': {'Date': '2018-12-24T00:00:00.000Z', 'Open': 560.0, 'High': 570.0, 'Low': 550.0, 'Close': 565.0, 'Volume': 45000000}}	{'_id': 1643, 'data': {'Date': '2018-12-24T00:00:00.000Z', 'Open': 2850.0, 'High': 2900.0, 'Low': 2800.0, 'Close': 2870.0, 'Volume': 21000000000}}
...
1644	{'_id': 1644, 'data': {'Date': '2018-12-24T00:00:00.000Z', 'Open': 114.0, 'High': 119.0, 'Low': 109.0, 'Close': 116.5, 'Volume': 82000000}}	{'_id': 1644, 'data': {'Date': '2018-12-24T00:00:00.000Z', 'Open': 164.0, 'High': 169.0, 'Low': 159.0, 'Close': 166.5, 'Volume': 105000000}}	{'_id': 1644, 'data': {'Date': '2018-12-24T00:00:00.000Z', 'Open': 114.0, 'High': 119.0, 'Low': 109.0, 'Close': 116.5, 'Volume': 82000000}}	{'_id': 1644, 'data': {'Date': '2018-12-24T00:00:00.000Z', 'Open': 79.0, 'High': 84.0, 'Low': 74.0, 'Close': 81.5, 'Volume': 62000000}}	{'_id': 1644, 'data': {'Date': '2018-12-24T00:00:00.000Z', 'Open': 565.0, 'High': 575.0, 'Low': 555.0, 'Close': 570.0, 'Volume': 42000000}}	{'_id': 1644, 'data': {'Date': '2018-12-24T00:00:00.000Z', 'Open': 2870.0, 'High': 2920.0, 'Low': 2820.0, 'Close': 2890.0, 'Volume': 20500000000}}

Ps:On remarque que tous les enregistrements des prix des stocks sont effectués dans la même date d'où le résultat suivant :

```
market_data=pd.read_json('market_data.json')
L=list(market_data.columns)
#On remarque que pour tous les enregistrements de tous les stocks :ont le meme id et le meme stock , donc on procede comme suit
market_data['id']=market_data[L[0]].apply(lambda x:x['_id'])
market_data['Date']=market_data[L[0]].apply(lambda x:x['data']['Date'])
for t in L:
    market_data[t]=market_data[t].apply(lambda x: x['data']['price'])
market_data
```

Out[31]:

	MSFT	IBM	AAPL	FB	GOOG	NDX	id	Date
0	24.327055	146.358612	65.592659	38.230000	299.078979	2478.530029	0	2012-05-18T00:00:00
1	NaN	147.763336	69.414101	34.029999	305.908386	2545.429932	1	2012-05-21T00:00:00
2	24.734301	147.060989	68.881065	31.000000	299.278229	2539.199951	2	2012-05-22T00:00:00
3	24.194075	146.537994	70.561768	32.000000	303.592072	2547.080078	3	2012-05-23T00:00:00
4	24.160828	146.515564	69.913734	33.029999	300.702881	2531.350098	4	2012-05-24T00:00:00
...
1642	101.617836	108.498039	157.706192	133.240005	1023.010010	6342.970215	1642	2018-12-19T00:00:00
1643	92.248878	100.241623	143.924454	124.059998	976.219971	5899.350098	1643	2018-12-24T00:00:00
1644	99.158005	106.028572	153.059998	134.520004	1043.880005	6288.299805	1644	2018-12-27T00:00:00
1645	98.383781	105.329666	153.138428	133.199997	1037.079956	6285.270020	1645	2018-12-28T00:00:00



02

Définir une fonction qui retourne les prix
des stocks spécifiés sinon par défaut
tous les stocks , pendant une période
déterminée

on a défini une fonction sur Python qui permet d'afficher le prix des stocks sur la période comprise entre la date de début et la date de fin, ces dernières étant prises comme arguments au niveau de la fonction suivante:

```
Entrée [2]: #Step1#Reponse
#on va creer une fonction qui prend comme arguments positionnels datedebut et datefin ,et un argument par default , une liste
#qui contient tous les stock , C'est une fonction qui retourne les prix des stock specifiee sinon par default tous les stock
#pendant la periode specifiee :
T=['MSFT','IBM','AAPL','FB','GOOG','NDX']
def ExtractDonnee(datedebut,datefin,Stock=T):
    #Organisation de La dataframe
    market_data=pd.read_json('market_data.json')
    L=list(market_data.columns)
    #On remarque que pour tous les enregistrements de tous les stocks :ont le meme id et le meme stock , donc on procede comme suit
    market_data['id']=market_data[L[0]].apply(lambda x:x['_id'])
    market_data['Date']=market_data[L[0]].apply(lambda x:x['data']['Date'])
    for t in L:
        market_data[t]=market_data[t].apply(lambda x: x['data']['price'])
    #eT Voila LE RESultat , on obtient une dataframe bien organisee

    market_data.set_index('Date',inplace=True)
    datedebut=datedebut+'T00:00:00'
    datefin=datefin+'T00:00:00'
    N=market_data.copy()
    return N.loc[datedebut:datefin,Stock]
```


Testons cette fonction sur la période comprise du 01-01-2017 au 31-12-2018, notre fonction retourne le résultat suivant :

```
S=ExtractDonnee('2017-01-01','2018-12-31').reset_index()
```

```
S
```

Out[4]:

	Date	MSFT	IBM	AAPL	FB	GOOG	NDX
0	2017-01-03T00:00:00	58.969059	143.492233	110.392334	116.860001	786.140015	4911.330078
1	2017-01-04T00:00:00	58.705212	145.268768	110.268791	118.690002	786.900024	4937.200195
2	2017-01-05T00:00:00	58.705212	144.788177	110.829552	120.669998	794.020020	4964.950195
3	2017-01-06T00:00:00	59.214054	145.500534	112.065109	123.410004	806.150024	5007.080078
4	2017-01-09T00:00:00	59.025597	143.886993	113.091560	124.900002	806.650024	5024.899902
...
479	2018-12-19T00:00:00	101.617836	108.498039	157.706192	133.240005	1023.010010	6342.970215
480	2018-12-24T00:00:00	92.248878	100.241623	143.924454	124.059998	976.219971	5899.350098
481	2018-12-27T00:00:00	99.158005	106.028572	153.059998	134.520004	1043.880005	6288.299805
482	2018-12-28T00:00:00	98.383781	105.329666	153.138428	133.199997	1037.079956	6285.270020
483	2018-12-31T00:00:00	99.540192	105.926056	154.618546	131.089996	1035.609985	6329.959961

484 rows × 7 columns

Resultat step 1

```
Entrée [6]: #Resultat step1:  
#On appelle la fonction ExtractDonnee en specifiant les stock :  
S=ExtractDonnee('2017-01-01','2018-12-31',['GOOG','AAPL']).reset_index()  
S
```

Out[6]:

	Date	GOOG	AAPL
0	2017-01-03T00:00:00	786.140015	110.392334
1	2017-01-04T00:00:00	786.900024	110.268791
2	2017-01-05T00:00:00	794.020020	110.829552
3	2017-01-06T00:00:00	806.150024	112.065109
4	2017-01-09T00:00:00	806.650024	113.091560
...
479	2018-12-19T00:00:00	1023.010010	157.706192
480	2018-12-24T00:00:00	976.219971	143.924454
481	2018-12-27T00:00:00	1043.880005	153.059998
482	2018-12-28T00:00:00	1037.079956	153.138428
483	2018-12-31T00:00:00	1035.609985	154.618546

484 rows × 3 columns





03

- Analyse des données manquantes,
 - Gérer les Datamissing
-

A ce stade, nous nous intéressons particulièrement à deux actions uniquement, à savoir AAPL et GOOG. Nous avons utilisé la fonction définie : on spécifie cette information au niveau de l'argument de notre fonction, de la manière suivante :

```
Entrée [3]: #Resultat step1:
#On appelle la fonction ExtractDonnee en specifiant Les stock :
J=ExtractDonnee('2017-01-01','2018-12-31',['GOOG','AAPL']).reset_index()
#step2:
#Premiere Partie
#Pour gerer Les Datamissing , je propose deux methode , celui de l'interpolation lineaire ou supprimer Les donnees ;
#Methode 1 :Interpolation lineaire
#Mais avant , Checkons s'il y a Missing data dans notre base de donnee
J.isnull().describe()
```

Out[3]:

	Date	GOOG	AAPL
count	484	484	484
unique	1	1	2
top	False	False	False
freq	484	484	483

En vue de gérer les datamissing dans notre base de données, 2 cas de figure se présentent : - la suppression des enregistrements qui contiennent des datamissings - l'utilisation d'une interpolation linéaire sur la période déterminée J'ai choisi d'utiliser la seconde méthode qui peut être implémentée de la manière suivante

Entrée [4]: *#On remarque que stock appl a seulement 483 valeurs de false , donc y a surement une valeur True , ce qui signifie une valeur manquante pour le stock app .Pour gerer cela on va utiliser methode d'interpolation lineaire*
`J=J.interpolate(method ='linear', limit_direction ='forward')`
#Verification
`J.isnull().describe()`

Out[4]:

	Date	GOOG	AAPL
count	484	484	484
unique	1	1	1
top	False	False	False
freq	484	484	484

Entrée [11]: *#Deuxieme Partie*
#Dans notre portefeuille on detient 100 \$ de GooG et 100\$ de AAPL à la date precise dans l'enonce , donc on detient X1 de Google et X2 de Y, determinons X1 et X2 :
`X1=100/(J['GOOG'].iloc[0])`
`X2=100/(J['AAPL'].iloc[0])`
`X2,X1`

Out[11]: (0.9058600030519698, 0.12720380356763814)

Le graphe qui représente l'évolution des prix des stocks en fonction du temps est donné comme suit :

Entrée [13]:

```
#notre investissement evolue tant quele prix evolue :  
J['InvestGoog']=J['GOOG'].apply(lambda x: x*X1)  
J['InvestaPPLE']=J['AAPL'].apply(lambda x: x*X2)  
J
```

Out[13]:

	Date	GOOG	AAPL	InvestGoog	InvestaPPLE
0	2017-01-03T00:00:00	786.140015	110.392334	100.000000	100.000000
1	2017-01-04T00:00:00	786.900024	110.268791	100.096676	99.888088
2	2017-01-05T00:00:00	794.020020	110.829552	101.002367	100.396058
3	2017-01-06T00:00:00	806.150024	112.065109	102.545349	101.515300
4	2017-01-09T00:00:00	806.650024	113.091560	102.608951	102.445121
...
479	2018-12-19T00:00:00	1023.010010	157.706192	130.130764	142.859732
480	2018-12-24T00:00:00	976.219971	143.924454	124.178893	130.375406
481	2018-12-27T00:00:00	1043.880005	153.059998	132.785507	138.650930
482	2018-12-28T00:00:00	1037.079956	153.138428	131.920515	138.721977
483	2018-12-31T00:00:00	1035.609985	154.618546	131.733529	140.062756

484 rows × 5 columns



En ce qui concerne le rendement des deux stocks Aapl et Goog, la fonction " pct_change " est utilisée dans ce cas de la manière suivante :

```
Entrée [24]: #Calcul des rendements quotidiens :  
J1=J[['InvestGoog','InvestaPPLE']]  
Returns=J1.pct_change()[1:]  
Returns['Date']=J['Date']  
Returns
```

Out[24]:

	InvestGoog	InvestaPPLE	Date
1	0.000967	-0.001119	2017-01-04T00:00:00
2	0.009048	0.005085	2017-01-05T00:00:00
3	0.015277	0.011148	2017-01-06T00:00:00
4	0.000620	0.009159	2017-01-09T00:00:00
5	-0.002306	0.001009	2017-01-10T00:00:00
...
479	-0.005541	-0.031192	2018-12-19T00:00:00
480	-0.045738	-0.087389	2018-12-24T00:00:00
481	0.069308	0.063475	2018-12-27T00:00:00
482	-0.006514	0.000512	2018-12-28T00:00:00
483	-0.001417	0.009665	2018-12-31T00:00:00

483 rows x 4 columns

Pour la moyenne et l'écart type, ils sont également fournis par les 2 fonctions prédéfinies sur Python :

100 rows x 8 columns

Entrée [25]: `#Calcul moyenne et ecart type des rendements quotidiens :Returns.mean()
Returns.mean()`

Out[25]: InvestGoog 0.000677
InvestaPPLE 0.000818
dtype: float64

Entrée [26]: `Returns.std()`

Out[26]: InvestGoog 0.014560
InvestaPPLE 0.015458
dtype: float64

Entrée []:

