

# Concepts Avancés Pandas

28 novembre 2025

## Résumé

Ce document explore les concepts avancés de Pandas : MultiIndex et indexation avancée, méthodes d'agrégation avancées, séries temporelles complexes, et techniques d'optimisation des performances. Chaque concept est expliqué avec sa théorie, ses mathématiques sous-jacentes, sa valeur ajoutée et des exemples pratiques.

## Table des matières

<b>1 MultiIndex &amp; Indexation Avancée</b>	<b>2</b>
1.1 Concept et Définition . . . . .	2
1.2 Implémentation . . . . .	2
1.3 Mathématiques Sous-Jacentes . . . . .	3
1.4 Sélection Avancée . . . . .	3
1.5 Valeur Ajoutée . . . . .	3
<b>2 Méthodes d'Agrégation Avancées</b>	<b>3</b>
2.1 Agrégations avec Dictionnaires . . . . .	3
2.2 Transformation vs Agrégation . . . . .	4
2.3 Mathématiques des Agrégations . . . . .	4
2.4 Valeur Ajoutée . . . . .	5
<b>3 Time Series Avancées</b>	<b>5</b>
3.1 Resampling Complexe . . . . .	5
3.2 Rolling Windows Avancées . . . . .	5
3.3 Mathématiques des Séries Temporelles . . . . .	6
3.3.1 Resampling . . . . .	6
3.3.2 Rolling Statistics . . . . .	6
3.3.3 Régression Linéaire Mobile . . . . .	6
3.4 Valeur Ajoutée . . . . .	7
<b>4 Performance &amp; Optimisation</b>	<b>7</b>
4.1 Opérations Vectorisées vs Apply . . . . .	7
4.2 Optimisation avec Numba . . . . .	8
4.3 Techniques d'Optimisation Avancées . . . . .	8
4.4 Mathématiques de l'Optimisation . . . . .	9
4.5 Valeur Ajoutée . . . . .	9

<b>5 Challenges Pratiques</b>	<b>9</b>
5.1 Challenge 1 : Analyse Financière Multi-Niveaux . . . . .	9
5.2 Challenge 2 : Système de Monitoring Temporel . . . . .	10
5.3 Challenge 3 : Optimisation de Performance . . . . .	10
<b>6 Conclusion</b>	<b>10</b>
6.1 Performances Typiques . . . . .	10
6.2 Recommandations Finales . . . . .	10

# 1 MultiIndex & Indexation Avancée

## 1.1 Concept et Définition

Un **MultiIndex** (index hiérarchique) permet de créer des structures de données multidimensionnelles dans Pandas, similaires aux cubes de données en OLAP.

## 1.2 Implémentation

```

1 import pandas as pd
2 import numpy as np
3
4 # Cr ation de MultiIndex
5 index = pd.MultiIndex.from_product([
6     ['2023', '2024'],
7     ['Q1', 'Q2', 'Q3', 'Q4']
8 ], names=['Ann e', 'Trimestre'])
9
10 df = pd.DataFrame({
11     'ventes': np.random.randint(1000, 5000, 8),
12     'profits': np.random.randint(100, 1000, 8)
13 }, index=index)
14
15 print("DataFrame avec MultiIndex:")
16 print(df)

```

Sortie :

		ventes	profits
	Année	Trimestre	
2023	Q1	3456	789
	Q2	2345	456
	Q3	4123	234
	Q4	2987	567
2024	Q1	3765	890
	Q2	2543	123
	Q3	3987	678
	Q4	3210	345

## 1.3 Mathématiques Sous-Jacentes

Le MultiIndex utilise une structure d'**arbre B** pour l'indexation :

- **Niveaux** : Chaque niveau représente une dimension
  - **Codes** : Références aux labels pour chaque niveau
  - **Performance** : Recherche en  $O(\log n)$  grâce à l'indexation hiérarchique
- La structure peut être représentée comme :

$$\text{MultiIndex} = \{\text{Niveau}_1 \times \text{Niveau}_2 \times \cdots \times \text{Niveau}_k\}$$

## 1.4 Sélection Avancée

```
1 # Selection par tuple
2 ventes_2023_Q1 = df.loc[('2023', 'Q1')]
3 print("Ventes 2023 Q1:")
4 print(ventes_2023_Q1)

5
6 # Selection par niveau avec xs
7 ventes_2023 = df.xs('2023', level='Ann e')
8 print("\nToutes les ventes 2023:")
9 print(ventes_2023)

10
11 # Slicing avanc
12 premier_trimestre = df.loc[(slice(None), 'Q1'), :]
13 print("\nTous les Q1:")
14 print(premier_trimestre)

15
16 # Selection conditionnelle multi-niveaux
17 ventes_elevees = df[df['ventes'] > 3000]
18 print("\nVentes > 3000:")
19 print(ventes_elevees)
```

## 1.5 Valeur Ajoutée

Avantage	Description
Organisation logique	Structure naturelle pour données hiérarchiques
Agrégations flexibles	Groupby sur différents niveaux
Performance	Accès rapide grâce à l'arbre B
Compatibilité	Intègre toutes les opérations Pandas

TABLE 1 – Valeur ajoutée du MultiIndex

## 2 Méthodes d'Agrégation Avancées

### 2.1 Agrégations avec Dictionnaires

```

1 # Agr gations avanc es avec dictionnaires
2 aggregations = {
3     'ventes': [ 'sum', 'mean', 'std'],
4     'profits': [ 'sum', lambda x: x.quantile(0.8)] # 80 me percentile
5 }
6
7 resultat_agg = df.groupby('Ann e').agg(aggregations)
8 print("Agr gations avanc es:")
9 print(resultat_agg)

```

Sortie :

Année	ventes			profits		
	sum	mean		std	sum	<lambda_0>
2023	12911	3227.750000	738.194061	2046	684.6	
2024	13505	3376.250000	663.069391	2036	795.4	

## 2.2 Transformation vs Agrégation

```

1 # Transformation - maintient la forme originale
2 df[ 'ventes_moyenne_groupe' ] = df.groupby('Ann e')[ 'ventes' ].transform(
3     'mean')
4 df[ 'ecart_moyenne' ] = df[ 'ventes' ] - df[ 'ventes_moyenne_groupe' ]
5
6 print("DataFrame avec transformations:")
7 print(df[['ventes', 'ventes_moyenne_groupe', 'ecart_moyenne']])
8
9 # Agr gation avec conditions
10 def taux_croissance(serie):
11     return (serie.iloc[-1] - serie.iloc[0]) / serie.iloc[0] * 100
12
13 croissance_annuelle = df.groupby('Ann e')[ 'ventes' ].agg(taux_croissance)
14 print("\nTaux de croissance annuel:")
15 print(croissance_annuelle)

```

## 2.3 Mathématiques des Agrégations

- **Aggregation** : Fonction de réduction  $f : \mathbb{R}^n \rightarrow \mathbb{R}$
- **Transformation** : Fonction de préservation  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$
- **Window functions** : Agrégations sur fenêtres glissantes

Formules mathématiques communes :

$$\text{Moyenne} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\text{Écart-type} = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\text{Quantile } q = x_{\lceil q \cdot n \rceil}$$

## 2.4 Valeur Ajoutée

Avantage	Description
Flexibilité	Combinaison libre de métriques
Performance	Opérations vectorisées natives
Expressivité	Syntaxe claire et concise
Réutilisabilité	Schémas d'agrégation réutilisables

TABLE 2 – Valeur ajoutée des agrégations avancées

## 3 Time Series Avancées

### 3.1 Resampling Complexe

```

1 # Cr ation de s ries temporelles
2 dates = pd.date_range('2023-01-01', periods=1000, freq='D')
3 ts = pd.Series(np.random.randn(1000).cumsum(), index=dates)
4
5 print("S rie temporelle originale:")
6 print(f" P riode: {ts.index.min()} to {ts.index.max()}")
7 print(f"Shape: {ts.shape}")
8
9 # Resampling mensuel avec agr gations custom
10 resampled = ts.resample('M').agg(['mean', 'std', lambda x: x.max() - x.
11     min()])
12 resampled.columns = ['moyenne', 'ecart_type', 'range']
13 print("\nS rie resampl e (mensuelle):")
14 print(resampled.head())

```

### 3.2 Rolling Windows Avancées

```

1 # Fen tres glissantes avec calculs complexes
2 def calcul_tendance(x):
3     """Calcule la pente de la tendance lin aire"""

```

```

4     if len(x) < 2:
5         return np.nan
6     return np.polyfit(range(len(x)), x, 1)[0] # pente de la regression
7
8 def efficiency_march (x):
9     """Ratio rendement/volatilit (Sharpe simplifi )"""
10    returns = np.diff(x) / x[:-1]
11    return np.mean(returns) / np.std(returns) if len(returns) > 1 else np
12        .nan
13
14 # Application des rolling windows
15 rolling_features = pd.DataFrame({
16     'mean': ts.rolling(window=30).mean(),
17     'volatility': ts.rolling(window=30).std(),
18     'trend': ts.rolling(window=30).apply(calcul_tendance, raw=True),
19     'efficiency': ts.rolling(window=30).apply(efficiency_march , raw=
20         True)
21 })
22
23 print("\nFeatures sur fen tres glissantes (30 jours):")
24 print(rolling_features.tail())

```

### 3.3 Mathématiques des Séries Temporelles

#### 3.3.1 Resampling

Aggrégation temporelle :  $f(\{x_t\}_{t \in T}) \rightarrow y_{T'}$

$$\text{Interpolation} : x_t = \sum w_i x_{t_i}$$

#### 3.3.2 Rolling Statistics

$$\begin{aligned} \text{Moyenne mobile} &= \frac{1}{w} \sum_{i=t-w+1}^t x_i \\ \text{Volatilité} &= \sqrt{\frac{1}{w} \sum_{i=t-w+1}^t (x_i - \bar{x}_w)^2} \end{aligned}$$

#### 3.3.3 Régression Linéaire Mobile

Pour une fenêtre de taille  $w$  :

$$\beta = (X^T X)^{-1} X^T y$$

où  $X = [1, t]$  et  $y$  sont les valeurs dans la fenêtre.

Avantage	Description
Analyse temporelle	Détection de tendances et saisonnalités
Features engineering	Préparation pour machine learning
Performance	Optimisé pour séries temporelles
Visualisation	Support natif pour plots temporels

TABLE 3 – Valeur ajoutée des séries temporelles avancées

### 3.4 Valeur Ajoutée

## 4 Performance & Optimisation

### 4.1 Opérations Vectorisées vs Apply

```

1 # Cr ation d'un grand dataset
2 np.random.seed(42)
3 grand_df = pd.DataFrame({
4     'A': np.random.rand(100000),
5     'B': np.random.rand(100000),
6     'C': np.random.rand(100000)
7 })
8
9 # MAUVAISE APPROCHE: apply row-wise
10 def methode_lente(row):
11     return row['A'] * row['B'] + row['C']
12
13 # BONNE APPROCHE: op rations vectoris es
14 grand_df['resultat_vectorise'] = grand_df['A'] * grand_df['B'] + grand_df
15     ['C']
16
17 # APPROCHE INTERM DIAIRE: apply avec axis=1
18 grand_df['resultat_apply'] = grand_df.apply(methode_lente, axis=1)
19
20 print("V rification quivalence :")
21 print(f"R sultats identiques: {np.allclose(grand_df['resultat_vectorise'],
22     'resultat_apply'])}")
23
24 # Test de performance
25 import time
26
27 def tester_performance():
28     # Vectoris
29     start = time.time()
30     resultat_vec = grand_df['A'] * grand_df['B'] + grand_df['C']
31     temps_vec = time.time() - start
32
33     # Apply
34     start = time.time()
35     resultat_app = grand_df.apply(methode_lente, axis=1)

```

```

34     temps_app = time.time() - start
35
36     print(f"\nPerformance:")
37     print(f"Vectoris : {temps_vec:.4f}s")
38     print(f"Apply: {temps_app:.4f}s")
39     print(f"Acc l ration: {temps_app/temps_vec:.1f}x")
40
41 tester_performance()

```

## 4.2 Optimisation avec Numba

```

1 from numba import jit
2 import numba
3
4 # Fonction optimisée avec Numba
5 @jit(nopython=True)
6 def calcul_rapide(a, b, c):
7     return a * b + c
8
9 # Application sur les arrays numpy sous-jacents
10 valeurs_a = grand_df['A'].values
11 valeurs_b = grand_df['B'].values
12 valeurs_c = grand_df['C'].values
13
14 start = time.time()
15 resultat_numba = calcul_rapide(valeurs_a, valeurs_b, valeurs_c)
16 temps_numba = time.time() - start
17
18 print(f"\nAvec Numba: {temps_numba:.4f}s")
19
20 # Utilisation avec pandas
21 grand_df['resultat_numba'] = calcul_rapide(valeurs_a, valeurs_b,
22     valeurs_c)

```

## 4.3 Techniques d'Optimisation Avancées

```

1 # 1. Utilisation de catégories pour les données textuelles
2     # et partitives
3 df_categories = df.copy()
4 df_categories['Année'] = df_categories['Année'].astype('category')
5 print(f"\nconomie mémoire: {df['Année']. nbytes} -> {df_categories['
6     Année']. nbytes} bytes")
7
8 # 2. Query optimization
9 résultat_query = df.query("ventes > 3000 and profits > 500")
10 print("\nRésultat query optimisé:")
11 print(résultat_query)

```

```

10
11 # 3. Mémoire usage optimization
12 def optimiser_memoire(df):
13     """Optimise l'usage mémoire d'un DataFrame"""
14     for col in df.columns:
15         if df[col].dtype == 'float64':
16             df[col] = df[col].astype('float32')
17         elif df[col].dtype == 'int64':
18             df[col] = df[col].astype('int32')
19     return df
20
21 df_opt = optimiser_memoire(df.copy())
22 print(f"\nMémoire originale: {df.memory_usage(deep=True).sum()} bytes")
23 print(f"Mémoire optimisée: {df_opt.memory_usage(deep=True).sum()} bytes"
      ")

```

## 4.4 Mathématiques de l'Optimisation

- **Complexité algorithmique** : Choix d'algorithmes  $O(n)$  vs  $O(n^2)$
- **Vectorisation** : Utilisation d'instructions SIMD
- **Localité mémoire** : Accès contigus vs aléatoires
- **Parallélisation** : Distribution des calculs

## 4.5 Valeur Ajoutée

Avantage	Description
Performance	Réduction temps calcul 10x-1000x
Échelle	Traitement de datasets massifs
Coût	Réduction infrastructure nécessaire
Productivité	Temps développement réduit

TABLE 4 – Valeur ajoutée de l'optimisation

## 5 Challenges Pratiques

### 5.1 Challenge 1 : Analyse Financière Multi-Niveaux

Créez un système d'analyse financière avec MultiIndex :

```

1 def challenge_finance():
2     # TODO: Créer un MultiIndex avec: [Entreprise, Année, Trimestre]
3     # Données: CA, bénéfice, dépenses
4     # Analyses requises:
5     # 1. Croissance trimestrielle par entreprise
6     # 2. Marges par année
7     # 3. Performance relative entre entreprises

```

```
8  
9     # Votre code ici  
10    pass
```

## 5.2 Challenge 2 : Système de Monitoring Temporel

Développez un système de monitoring avec séries temporelles :

```
1 def challenge_monitoring():  
2     # TODO: Cr er des s ries temporelles de m triques syst me  
3     # CPU, m moire, r seau, stockage  
4     # Analyses requises:  
5     # 1. D tection d'anomalies avec rolling stats  
6     # 2. Resampling pour rapports horaires/quotidiens  
7     # 3. Features engineering pour ML  
8  
9     # Votre code ici  
10    pass
```

## 5.3 Challenge 3 : Optimisation de Performance

Optimisez un traitement de données massives :

```
1 def challenge_optimisation():  
2     # TODO: Prendre un dataset de 1M+ lignes  
3     # Optimiser:  
4     # 1. Temps de traitement (vectorisation, numba)  
5     # 2. Usage m moire (types, cat gories)  
6     # 3. Requ tes complexes (query, indexing)  
7  
8     # Votre code ici  
9     pass
```

# 6 Conclusion

## Points Clés à Retenir

- **MultiIndex** : Structure hiérarchique pour données complexes
- **Agrégations avancées** : Flexibilité dans le calcul des métriques
- **Séries temporelles** : Analyse sophistiquée des données temporelles
- **Optimisation** : Techniques cruciales pour les grands datasets
- **Vectorisation** : Clé pour les performances avec Pandas

## 6.1 Performances Typiques

## 6.2 Recommandations Finales

1. Préférer MultiIndex pour les données hiérarchiques

Technique	Temps d'exécution	Usage mémoire
Boucles Python	100%	100%
Apply pandas	50-80%	100-120%
Opérations vectorisées	1-10%	100%
Numba	0.1-5%	100%

TABLE 5 – Comparaison des performances (valeurs relatives)

2. Utiliser les agrégations dictionnaires pour la flexibilité
3. Exploiter les séries temporelles pour l'analyse temporelle
4. Toujours vectoriser les opérations sur les grands datasets
5. Profiler régulièrement les performances