

Lire des fichiers

Modèle	Exemple	Description
pd.read_csv(filepath, sep=..., encoding=...)	pd.read_csv("data.csv", sep=";", encoding="utf-8")	Lire un fichier CSV
pd.read_excel(filepath, sheet_name=..., usecols=..., engine=...)	pd.read_excel("data.xlsx", sheet_name="Feuil1", usecols="A:C", engine='calamine')	Lire un fichier Excel
pd.read_pickle(filepath)	pd.read_pickle("data.pkl")	Lire un fichier pickle (.pkl)

Explorer un DataFrame

Modèle	Exemple	Description
df.shape	df.shape	Nombre de lignes et colonnes
df.index	df.index	Accès à l'index du DataFrame
df.columns	df.columns	Liste des noms de colonnes
df.dtypes	df.dtypes	Type de chaque colonne
df.info()	df.info()	Résumé des colonnes, types, non-nuls et mémoire
df.describe()	df.describe()	Statistiques descriptives sur colonnes numériques
df.isna().sum()	df.isna().sum()	Nombre de valeurs manquantes par colonne
df.head(n)	df.head(10)	Afficher les n premières lignes (défaut n=5)
df.tail(n)	df.tail(3)	Afficher les n dernières lignes (défaut n=5)

Explorer une Series

Modèle	Exemple	Description
series.unique()	df["country"].unique()	Valeurs uniques dans la Series
series.nunique()	df["gender"].nunique()	Nombre de valeurs uniques
series.sort_values()	df["score"].sort_values()	Trier les valeurs (ordre croissant par défaut)
series.value_counts(normalize=...)	df["region"].value_counts(normalize=True)	Compter les occurrences (avec % si normalize=True)

Slicer un DataFrame

Modèle	Exemple	Description
df.loc[ligne, colonne]	df.loc[0, "age"]	Accès par nom de ligne et colonne
df.iloc[ligne, colonne]	df.iloc[0, 2]	Accès par index numérique de ligne et colonne
df.at[ligne, colonne]	df.at[0, "name"]	Accès rapide à une cellule (par nom)
df.iat[ligne, colonne]	df.iat[0, 1]	Accès rapide à une cellule (par index numérique)

Filtrer un DataFrame avec .query()

Nom	Utilité
df.query(condition)	Filtre le DataFrame selon la condition (nom de colonne, valeur)

Modèle	Exemple	Description
col == valeur	df.query("year == 2020")	Égal à
col != valeur	df.query("continent != 'Asia'")	Différent de
col > valeur	df.query("age > 18")	Supérieur à
col >= valeur	df.query("score >= 90")	Supérieur ou égal
col < valeur	df.query("price < 100")	Inférieur à
col <= valeur	df.query("weight <= 70")	Inférieur ou égal
col in @liste	df.query("country in @european_countries")	Appartient à une liste (utilisation de @ pour une variable externe)
col not in @liste	df.query("brand not in @excluded_brands")	N'appartient pas à une liste
col.str.contains('texte')	df.query("name.str.contains('John', case=False, na=False)")	Contient une sous-chaîne
col.isna() ou col.notna()	df.query("score.isna()")	Valeurs manquantes
Conditions multiples avec and	df.query("year > 2010 and year < 2020")	Condition composée (ET logique)
Conditions multiples avec or	df.query("continent == 'Europe' or continent == 'Asia'")	Condition composée (OU logique)
col.between(val1, val2)	df.query("age.between(20, 30)")	Valeur comprise entre deux bornes
Négation avec ~	df.query("~(status == 'inactive')")	Négation logique
Filtrage direct de booléens	df.query("is_active") OU df.query("~is_deleted")	Colonnes booléennes (True/False)

Filtrer un DataFrame avec .loc

Modèle	Exemple	Description
<code>df.loc[df[col] == valeur]</code>	<code>df.loc[df[year] == 2020]</code>	Égal à
<code>df.loc[df[col] != valeur]</code>	<code>df.loc[df[continent] != 'Asia']</code>	Différent de
<code>df.loc[df[col] > valeur]</code>	<code>df.loc[df[age] > 18]</code>	Supérieur à
<code>df.loc[df[col] >= valeur]</code>	<code>df.loc[df[score] >= 90]</code>	Supérieur ou égal
<code>df.loc[df[col] < valeur]</code>	<code>df.loc[df[price] < 100]</code>	Inférieur à
<code>df.loc[df[col] <= valeur]</code>	<code>df.loc[df[weight] <= 70]</code>	Inférieur ou égal
<code>df.loc[df[col].isin(liste)]</code>	<code>df.loc[df[country].isin(european_countries)]</code>	Appartient à une liste
<code>df.loc[~df[col].isin(liste)]</code>	<code>df.loc[~df[brand].isin(excluded_brands)]</code>	N'appartient pas à une liste
<code>df.loc[df[col].str.contains('texte', case=..., na=False)]</code>	<code>df.loc[df[name].str.contains('John', case=False, na=False)]</code>	Contient une sous-chaîne
<code>df.loc[df[col].isna()]</code>	<code>df.loc[df[score].isna()]</code>	Valeurs manquantes
<code>df.loc[(cond1) & (cond2)]</code>	<code>df.loc[(df[year] > 2010) & (df[year] < 2020)]</code>	Condition composée (ET logique)
<code>df.loc[(cond1) (cond2)]</code>	<code>`df.loc[(df[continent] == 'Europe') (df[continent] == 'Asia')]</code>	Condition composée (OU logique)
<code>df.loc[df[col].between(val1, val2)]</code>	<code>df.loc[df[age].between(20, 30)]</code>	Valeur comprise entre deux bornes
<code>df.loc[~(condition)]</code>	<code>df.loc[~(df[status] == 'inactive')]</code>	Négation logique

Ajouter une colonne

Modèle	Exemple	Description
<code>df[nouvelle_col] = valeur</code>	<code>df[discounted] = False</code>	Ajouter une colonne avec une valeur fixe
<code>df[col2] = df[col1] * facteur</code>	<code>df[prix_ttc] = df[prix] * 1.2</code>	Colonne calculée à partir d'une autre
<code>df[col3] = df[col1] + df[col2]</code>	<code>df[total] = df[net] + df[taxe]</code>	Combinaison de colonnes existantes
<code>df[col] = df[col].apply(fonction)</code>	<code>df[upper_name] = df[name].apply(str.upper)</code>	Transformation via une fonction (lent)
<code>df[col] = df[col].apply(function lambda)</code>	<code>df[cat_age_2] = df[age].apply(lambda x : "-25" if age < 25 else "25+")</code>	Transformation via une fonction lambda (lent)
<code>df[col] = pd.cut(df[val], bins=...)</code>	<code>df[age_group] = pd.cut(df[age], bins=[0,18,65,100])</code>	Création de n tranches (rapide)
<code>df[col] = np.where(condition, A, B)</code>	<code>df[statut] = np.where(df[score] >= 10, "valid", "retry")</code>	Conditionnelle avec NumPy
<code>df[col] = df[col].map(dictionnaire)</code>	<code>df[continent_name] = df[continent].map(code_to_label)</code>	Remplacement via un dictionnaire
<code>df.loc[condition, col] = valeur</code>	<code>df.loc[df[age] < 18, categorie] = "minor"</code>	Ajout conditionnel à certaines lignes
<code>df = df.assign(col=...)</code>	<code>df = df.assign(next_year = df[year] + 1)</code>	Ajout sans modifier l'original, méthode chainable

Supprimer une ligne ou une colonne

Modèle	Exemple	Description
<code>df.drop(nom_colonne, axis=1)</code>	<code>df.drop(age, axis=1)</code>	Supprimer une colonne
<code>df.drop([col1, col2], axis=1)</code>	<code>df.drop(["age", "gender"], axis=1)</code>	Supprimer plusieurs colonnes
<code>df.drop(index, axis=0)</code>	<code>df.drop(0, axis=0)</code>	Supprimer une ligne par son index
<code>df.drop([index1, index2])</code>	<code>df.drop([0, 5])</code>	Supprimer plusieurs lignes
<code>df.drop(df[df.colonne == valeur].index)</code>	<code>df.drop(df[df[status] == inactive].index)</code>	Supprimer des lignes selon une condition
<code>df = df[df.colonne != valeur]</code>	<code>df = df[df[status] != inactive]</code>	Filtrage inversé (équivalent à suppression)
<code>df = df[~df.colonne.isin(liste)]</code>	<code>df = df[~df[country].isin(banned_countries)]</code>	Supprimer lignes selon une liste de valeurs
<code>df = df.dropna(axis=0)</code>	<code>df = df.dropna(axis=0, how='all')</code> Supprimée si toutes les valeurs sont na	Supprimer les lignes avec valeurs manquantes
	<code>df = df.dropna(axis=0, how='any')</code> Supprimée si au – 1 valeur est na. Valeur par défaut	
<code>df = df.dropna(axis=1)</code>	<code>df = df.dropna(axis=1)</code>	Supprimer les colonnes avec valeurs manquantes
<code>df = df.loc[df[col].notna()]</code>	<code>df = df.loc[df[score].notna()]</code>	Conserver seulement les lignes non nulles

Concaténer deux DataFrames

Modèle	Exemple	Description
<code>pd.concat([df1, df2])</code>	<code>pd.concat([df1, df2])</code>	Concaténation verticale (lignes) par défaut
<code>pd.concat([df1, df2], axis=1)</code>	<code>pd.concat([df1, df2], axis=1)</code>	Concaténation horizontale (colonnes)
<code>pd.concat([df1, df2], ignore_index=True)</code>	<code>pd.concat([df1, df2], ignore_index=True)</code>	Réindexe les lignes après concat
<code>pd.concat([df1, df2], keys=["A", "B"])</code>	<code>pd.concat([df1, df2], keys=["2023", "2024"])</code>	Ajoute un niveau hiérarchique d'index
<code>pd.concat([df1, df2], join="inner")</code>	<code>pd.concat([df1, df2], join="inner")</code>	Conserve uniquement les colonnes communes
<code>pd.concat([df1, df2], axis=1, join="inner")</code>	<code>pd.concat([df1, df2], axis=1, join="inner")</code>	Jointure uniquement sur les index communs

Merge (= Jointure SQL sur valeurs de colonne) deux DataFrames

Modèle	Exemple	Description
<code>df1.merge(df2, how='inner', on="col")</code>	<code>df1.merge(df2, how='inner', #défaut on="id")</code>	Jointure interne sur une clé commune
<code>df1.merge(df2, how="left", on="col")</code>	<code>df1.merge(df2, how="left", on="id")</code>	Jointure à gauche
<code>df1.merge(df2, how="right", on="col")</code>	<code>df1.merge(df2, how="right", on="id")</code>	Jointure à droite
<code>df1.merge(df2, how="outer", on="col")</code>	<code>df1.merge(df2, how="outer", on="id")</code>	Jointure complète
<code>df1.merge(df2, left_on="col1", right_on="col2")</code>	<code>df1.merge(df2, left_on="user_id", right_on="id")</code>	Noms de clés différents
<code>df1.merge(df2, on="col", suffixes=("_A", "_B"))</code>	<code>df1.merge(df2, on="id", suffixes=("_left", "_right"))</code>	Gérer les colonnes en double
<code>df1.merge(df2, on="col", how="outer", sort=False)</code>	<code>df1.merge(df2, on="id", how="outer", sort=False)</code>	Jointure sans trier, gain de temps
<code>df1.merge(df2, on="col", validate="1:1")</code>	<code>df1.merge(df2, on="id", validate="1:1")</code>	Valider le type de jointure (ex : 1:1, m:1...)

Join (=Jointure sur index ou nom de colonne) deux DataFrames

Modèle	Exemple	Description
<code>df1.join(df2)</code>	<code>df1.join(df2)</code>	Jointure à gauche sur les index
<code>df1.join(df2, how="inner")</code>	<code>df1.join(df2, how="inner")</code>	Jointure interne sur les index communs
<code>df1.join(df2, how="outer")</code>	<code>df1.join(df2, how="outer")</code>	Jointure complète sur les index
<code>df1.join(df2, how="left", sort=False)</code>	<code>df1.join(df2, how="left", sort=False)</code>	Jointure gauche sans trier, gain de temps
<code>df1.join(df2, lsuffix="_left", rsuffix="_right")</code>	<code>df1.join(df2, lsuffix="_A", rsuffix="_B")</code>	Gérer les doublons de colonnes
<code>df1.join([df2, df3])</code>	<code>df1.join([df2, df3])</code>	Joindre plusieurs df en même temps

Agréger un DataFrame : simple opération

Modèle	Exemple	Description
<code>df.agg("fonction")</code>	<code>df.agg("mean")</code>	Agrégation sur tout un dataframe
<code>df["col"].agg("fonction")</code>	<code>df["price"].agg("sum")</code>	Agréger une seule colonne
<code>df.agg(["fonction1", "fonction2"])</code>	<code>df.agg(["min", "max"])</code>	Appliquer plusieurs fonctions
<code>df.agg({"col1": "fonction", "col2": "fonction"})</code>	<code>df.agg({"price": "mean", "qty": "sum"})</code>	Agréger colonne par colonne

Exemples de fonctions : median, min, max, sum, prod, first, last, cumsum, quantile, idxmin, idxmax, fonctions lambda...

Agréger un DataFrame avec .groupby()

Modèle	Exemple	Description
<code>df.groupby(« col »).agg(« fonction »)</code>	<code>df.groupby("category").agg("mean")</code>	Agrégation par groupe
<code>df.groupby(« col »)[[« c1 », « c2 »]].agg(« fonction »)</code>	<code>df.groupby("brand")[["price", "qty"]].agg("sum")</code>	Agrégation multi-colonnes par groupe
<code>df.groupby(« col »).agg({ « c1 » : « fonction1 », « c2 » : « fonction2 » })</code>	<code>df.groupby("type").agg({ "price": "mean", "qty": "sum" })</code>	Fonctions différentes selon les colonnes
<code>df.groupby([« col1 », « col2 »]).agg(« fonction »)</code>	<code>df.groupby(["region", "year"]).agg("sum")</code>	Agrégation sur plusieurs colonnes
<code>df.groupby(« col »).size()</code>	<code>df.groupby(« gender »).size()</code>	Taille de chaque groupe
<code>df.groupby(« col »).count()</code>	<code>df.groupby(« team »).count()</code>	Compte non-nul par colonne dans chaque groupe

Les exemples de fonctions possibles avec .agg() sont généralement possibles avec .groupby()

Agréger un DataFrame avec .pivot_table()

Modèle	Exemple	Description
<code>df.pivot_table(values="val", index="row", columns="col", aggfunc="fonction")</code>	<code>df.pivot_table(values="sales", index="product", columns="month", aggfunc="sum")</code>	Tableau croisé dynamique