

Installation et Imports

```
In [1]: # Requirements d'installation
# pip install pandas pyarrow fastparquet

import pandas as pd
import pyarrow as pa
import pyarrow.parquet as pq
import numpy as np
import time

print("Démonstration Parquet Simplifiée - NYC Taxi Data")

Démonstration Parquet Simplifiée - NYC Taxi Data
```

Fonctions de Téléchargement et Simulation

```
In [2]: def download_taxi_data():
    """Télécharge les données NYC Taxi 2025 ou crée des données simulées"""

    # URLs des données NYC Taxi 2025
    url_2025 = [
        "https://d37id6v9uayc8a.cloudfront.net/trip-data/yellow_tripdata_2025-01.parquet",
        "https://d37id6v9uayc8a.cloudfront.net/trip-data/yellow_tripdata_2024-12.parquet",
        "https://d37id6v9uayc8a.cloudfront.net/trip-data/yellow_tripdata_2024-11.parquet"
    ]

    for url in url_2025:
        try:
            print(f"📡 Tentative de téléchargement: {url.split('/')[-1]}")
            df = pd.read_parquet(url)
            print(f"📊 Données téléchargées: {len(df)}, {len(df.columns)} lignes")
            return df.head(10000) # Limiter pour la démo
        except Exception as e:
            print(f"❌ Erreur: {e}")
            continue

    print("Tous les téléchargements ont échoué, création de données simulées...")
    return create_simulated_taxi_data()

def create_simulated_taxi_data(n_rows=50000):
    """Crée des données de taxi simulées NYC style"""
    np.random.seed(42)

    data = {
        'tpep_pickup_datetime': pd.date_range('2025-01-01', periods=n_rows, freq='min'),
        'tpep_dropoff_datetime': pd.date_range('2025-01-01 00:10:00', periods=n_rows, freq='min'),
        'passenger_count': np.random.choice(1, 2, 3, 4, 5, n_rows, p=[0.5, 0.3, 0.1, 0.05, 0.05]),
        'trip_distance': np.random.exponential(1.5, n_rows),
        'fare_amount': np.random.gamma(2, 5, n_rows),
        'extra': np.random.choice([0, 0.5, 1], n_rows, p=[0.7, 0.2, 0.1]),
        'trip_amount': np.random.gamma(1, 2, n_rows),
        'ratecodeID': np.random.gamma(0.5, 1, n_rows),
        'pickup_location_id': np.random.randint(1, 265, n_rows),
        'dropoff_location_id': np.random.randint(1, 265, n_rows),
        'payment_type': np.random.choice(1, 2, 3, 4, n_rows, p=[0.6, 0.3, 0.05, 0.05])
    }

    df = pd.DataFrame(data)
    df['total_amount'] = df['fare_amount'] + df['tip_amount'] + df['extra']

    return df
```

Chargement des Données

```
In [3]: # Chargement des données
taxi_df = download_taxi_data()
print(f"📊 Dataset: {taxi_df.shape[0]}, {len(taxi_df.columns)} colonnes")

# Aperçu des données
print(f"📄 {len(taxi_df)} lignes des données")
print(taxi_df.head())
print(f"📄 Info sur le dataset:")
print(taxi_df.info())

📡 Tentative de téléchargement: yellow_tripdata_2025-01.parquet
Données téléchargées: 3,475,028 lignes
📊 Dataset: 50,000 lignes, 20 colonnes

📄 Aperçu des données:
VendorID tpep_pickup_datetime tpep_dropoff_datetime passenger_count \
0 1 2025-01-01 01:10:18 2025-01-01 01:16:19 1.0
1 1 2025-01-01 01:13:49 2025-01-01 01:19:13 1.0
2 1 2025-01-01 01:14:16 2025-01-01 01:16:01 1.0
3 2 2025-01-01 01:14:17 2025-01-01 01:20:01 3.0
4 2 2025-01-01 01:21:14 2025-01-01 01:29:16 1.0

trip_distance ratecodeID store_and_fwd_flag PULocationID DOLocationID \
0 1.60 1.0 N 229 237
1 0.50 1.0 N 236 237
2 0.40 1.0 N 141 141
3 0.52 1.0 N 244 244
4 0.46 1.0 N 244 116

payment_type fare_amount extra mta_tax tip_amount tolls_amount \
0 1 10.0 3.5 0.5 3.00 0.0
1 1 5.1 1.5 0.5 2.02 0.0
2 1 5.1 3.5 0.5 2.00 0.0
3 2 7.2 1.0 0.5 0.00 0.0
4 2 5.8 1.0 0.5 0.00 0.0

improvement_surcharge total_amount congestion_surcharge Airport_fee \
0 1.0 18.00 2.5 0.0
1 1.0 12.12 2.5 0.0
2 1.0 12.10 2.5 0.0
3 1.0 9.70 0.0 0.0
4 1.0 6.30 0.0 0.0

cbb_congestion_fee
0 0.0
1 0.0
2 0.0
3 0.0
4 0.0

📄 Info sur le dataset:
<class 'pandas.core.frame.DataFrame'>
Int64Index: 50000 entries, 0 to 49999
Data columns (total 20 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   VendorID              50000 non-null   int32
 1   tpep_pickup_datetime  50000 non-null   datetime64[us]
 2   tpep_dropoff_datetime 50000 non-null   datetime64[us]
 3   passenger_count        50000 non-null   float64
 4   trip_distance          50000 non-null   float64
 5   ratecodeID            50000 non-null   float64
 6   store_and_fwd_flag     50000 non-null   object
 7   PULocationID          50000 non-null   int32
 8   DOLocationID          50000 non-null   int32
 9   payment_type          50000 non-null   int64
10   fare_amount           50000 non-null   float64
11   extra                 50000 non-null   float64
12   mta_tax               50000 non-null   float64
13   tip_amount            50000 non-null   float64
14   tolls_amount          50000 non-null   float64
15   improvement_surcharge 50000 non-null   float64
16   total_amount          50000 non-null   float64
17   congestion_surcharge  50000 non-null   float64
18   Airport_fee           50000 non-null   float64
19   cbb_congestion_fee    50000 non-null   float64
dtypes: datetime64[us](2), float64(13), int32(3), int64(1), object(1)
memory usage: 7.1+ MB
None
```

Sauvegarde en Différents Formats

```
In [4]: print("📄 * * * * 60")
print("📄 SAUVEGARDE EN DIFFÉRENTS FORMATS")
print("📄 * * * * 60")

# CSV
start_time = time.time()
taxi_df.to_csv('taxi_data.csv', index=False)
csv_time = time.time() - start_time
csv_size = os.path.getsize('taxi_data.csv') / 1024**2

# Parquet avec différentes compressions
print("📄 Sauvegarde Parquet (sans compression)...")
start_time = time.time()
taxi_df.to_parquet('taxi_data_none.parquet', compression=None)
parquet_none_time = time.time() - start_time
parquet_none_size = os.path.getsize('taxi_data_none.parquet') / 1024**2

print("📄 Sauvegarde Parquet (Snappy)...")
start_time = time.time()
taxi_df.to_parquet('taxi_data_snappy.parquet', compression='snappy')
parquet_snappy_time = time.time() - start_time
parquet_snappy_size = os.path.getsize('taxi_data_snappy.parquet') / 1024**2

print("📄 Sauvegarde Parquet (Gzip)...")
start_time = time.time()
taxi_df.to_parquet('taxi_data_gzip.parquet', compression='gzip')
parquet_gzip_time = time.time() - start_time
parquet_gzip_size = os.path.getsize('taxi_data_gzip.parquet') / 1024**2

print("📄 Comparaison des sauvegardes terminée!")
```

```
📄 SAUVEGARDE EN DIFFÉRENTS FORMATS
📄 Sauvegarde CSV...
📄 Sauvegarde Parquet (sans compression)...
📄 Sauvegarde Parquet (Snappy)...
📄 Sauvegarde Parquet (Gzip)...
📄 Toutes les sauvegardes terminées!
```

Comparaison des Tailles

```
In [5]: print("📄 * * * * 60")
print("📄 COMPARAISON DES TAILLES")
print("📄 * * * * 60")

print(f"Format          | Taille (MB) | Temps (s) | Ratio vs CSV")
print("📄 * * * * 60")
print(f"Format (CSV)          | (csv_size*8.16) | (csv_time*6.26) | 1.0x")
print(f"Format (Parquet (aucun)) | (parquet_none_size*8.16) | (parquet_none_time*6.26) | (csv_size/parquet_none_size:16x)")
print(f"Format (Parquet (Snappy)) | (parquet_snappy_size*8.16) | (parquet_snappy_time*6.26) | (csv_size/parquet_snappy_size:16x)")
print(f"Format (Parquet (Gzip)) | (parquet_gzip_size*8.16) | (parquet_gzip_time*6.26) | (csv_size/parquet_gzip_size:16x)")

print("📄 Meilleure compression: (csv_size/min(parquet_none_size, parquet_snappy_size, parquet_gzip_size):16x plus compact)")

📄 COMPARAISON DES TAILLES
Format          | Taille (MB) | Temps (s) | Ratio vs CSV
CSV              | 5.3 | 1.33 | 1.0x
Parquet (aucun) | 1.2 | 0.09 | 4.3x
Parquet (Snappy) | 1.1 | 0.09 | 5.0x
Parquet (Gzip)   | 0.8 | 0.27 | 6.2x

📄 Meilleure compression: 6.2x plus compact
```

Performance de Lecture Complète

```
In [6]: print("📄 * * * * 60")
print("📄 PERFORMANCE DE LECTURE")
print("📄 * * * * 60")

# Lecture complète CSV
start_time = time.time()
df_csv = pd.read_csv('taxi_data.csv')
csv_read_time = time.time() - start_time
csv_read_size = os.path.getsize('taxi_data.csv')

# Lecture complète Parquet
start_time = time.time()
df_parquet = pd.read_parquet('taxi_data_snappy.parquet')
parquet_read_time = time.time() - start_time

print(f"📄 Format          | Temps | Lignes | Lignes")
print("📄 * * * * 60")
print(f"Format (CSV)          | (csv_read_time*6.26) | (len(df_csv)) |")
print(f"Format (Parquet)      | (parquet_read_time*6.26) | (len(df_parquet)) |")
print("📄 Meilleure Parquet est (csv_read_time/parquet_read_time:16x plus rapide pour la lecture)")

📄 PERFORMANCE DE LECTURE
=====
📄 Test Lecture CSV...
📄 Test Lecture Parquet...

Format          | Temps | Lignes | Lignes
CSV              | 0.324 | 50,000
Parquet          | 0.064 | 50,000

📄 Parquet est 9.1x plus rapide pour la lecture
```

Performance de Filtrage

```
In [7]: print("📄 * * * * 60")
print("📄 PERFORMANCE DE FILTRAGE")
print("📄 * * * * 60")

# Test: Courses avec montant > 20$
print("📄 Test: Courses avec total_amount > 20$")

# Filtrage CSV (lecture complète puis filtrage)
start_time = time.time()
df_csv_full = pd.read_csv('taxi_data.csv')
df_csv_filtered = df_csv_full[df_csv_full['total_amount'] > 20]
csv_filter_time = time.time() - start_time

# Filtrage Parquet (avec predicate pushdown)
start_time = time.time()
df_parquet_filtered = pd.read_parquet('taxi_data_snappy.parquet', filters=[('total_amount', '>', 20)])
parquet_filter_time = time.time() - start_time

print(f"📄 Méthode          | Temps | Lignes résultat")
print("📄 * * * * 45")
print(f"Format (CSV (full)) | (csv_filter_time*6.26) | (len(df_csv_filtered)) |")
print(f"Format (Parquet (filter)) | (parquet_filter_time*6.26) | (len(df_parquet_filtered)) |")
print("📄 Meilleure Parquet filtre est (csv_filter_time/parquet_filter_time:16x plus rapide)")

📄 PERFORMANCE DE FILTRAGE
=====
📄 Test: Courses avec total_amount > 20$

Méthode          | Temps | Lignes résultat
CSV (full) | 0.398 | 25,095
Parquet (filter) | 0.036 | 25,095

📄 Parquet filtre est 12.2x plus rapide
```

Performance de Lecture par Colonnes

```
In [8]: # Test: Lecture colonnes spécifiques
print("📄 Test: Lecture colonnes spécifiques")
columns_to_read = ['tpep_pickup_datetime', 'passenger_count', 'total_amount']

# CSV - colonnes spécifiques
start_time = time.time()
df_csv_cols = pd.read_csv('taxi_data.csv', usecols=columns_to_read)
csv_cols_time = time.time() - start_time

# Parquet - colonnes spécifiques
start_time = time.time()
df_parquet_cols = pd.read_parquet('taxi_data_snappy.parquet', columns=columns_to_read)
parquet_cols_time = time.time() - start_time

print(f"📄 Méthode          | Temps | Colonnes")
print("📄 * * * * 60")
print(f"Format (CSV colonnes) | (csv_cols_time*6.26) | (len(df_csv_cols)) |")
print(f"Format (Parquet colonnes) | (parquet_cols_time*6.26) | (len(df_parquet_cols)) |")
print("📄 Meilleure Parquet colonnes est (csv_cols_time/parquet_cols_time:16x plus rapide)")

# Aperçu des données lues
print(f"📄 {len(df_csv_cols)} lignes des données lues")
print(df_csv_cols.head())

📄 Test: Lecture colonnes spécifiques

Méthode          | Temps | Colonnes
CSV colonnes      | 0.246 | 3
Parquet colonnes | 0.026 | 3

📄 Parquet lecture est 13.0x plus rapide
```

Métadonnées Parquet

```
In [9]: print("📄 * * * * 60")
print("📄 MÉTADONNÉES PARQUET")
print("📄 * * * * 60")

# Ouverture du fichier Parquet
parquet_file = pq.ParquetFile('taxi_data_snappy.parquet')

print("📄 INFORMATIONS GÉNÉRALES:")
print(f"📄 Nombre de lignes: {parquet_file.metadata.num_rows}")
print(f"📄 Nombre de colonnes: {parquet_file.metadata.num_columns}")
print(f"📄 Nombre de row groups: {parquet_file.metadata.num_row_groups}")
print(f"📄 Taille du fichier: {os.path.getsize('taxi_data_snappy.parquet')} / 1024**2 (16 MB)")
print(f"📄 Compression: {parquet_file.metadata.row_group(0).column(0).compression}")

print("📄 SCHEMA DES COLONNES:")
schema = parquet_file.schema_arrow
for i, field in enumerate(schema):
    print(f"📄 {i} {field.name:25} | {field.type}")

📄 MÉTADONNÉES PARQUET
=====
📄 INFORMATIONS GÉNÉRALES:
Nombre de lignes: 50,000
Nombre de colonnes: 20
Nombre de row groups: 1
Taille du fichier: 11 MB
Compression: SNAPPY

📄 SCHEMA DES COLONNES:
VendorID          | int32
tpep_pickup_datetime | timestamp[us]
tpep_dropoff_datetime | timestamp[us]
passenger_count    | double
trip_distance       | double
ratecodeID          | double
store_and_fwd_flag  | string
PULocationID        | int32
DOLocationID        | int32
payment_type        | int64
fare_amount          | double
extra                | double
mta_tax              | double
tip_amount           | double
tolls_amount         | double
improvement_surcharge | double
total_amount         | double
congestion_surcharge | double
Airport_fee          | double
cbb_congestion_fee   | double
```

Statistiques des Colonnes

```
In [10]: print("📄 * * * * 60")
print("📄 STATISTIQUES PAR COLONNE")
print("📄 * * * * 60")

# Statistiques du premier row group
rg_metadata = parquet_file.metadata.row_group(0)
schema_fields = parquet_file.schema_arrow

print("📄 STATISTIQUES AUTOMATIQUES:")
for i in range(min(10, rg_metadata.num_columns)): # Limiter à 5 colonnes
    col_metadata = rg_metadata.column(i)
    stats = col_metadata.statistics
    field_name = schema_fields[i].name
    field_type = schema_fields[i].type
    print(f"📄 {i} {col_metadata.path_in_schema}")
    print(f"📄 Type Arrow: {field_type}")
    print(f"📄 Type physique: {col_metadata.physical_type}")
    print(f"📄 Compression: {col_metadata.compression}")
    print(f"📄 Taille compressée: {col_metadata.total_uncompressed_size:,} bytes")
    print(f"📄 Taille non-compressée: {col_metadata.total_uncompressed_size:,} bytes")

    if stats:
        if stats.has_min_max:
            print(f"📄 Minimum: {stats.min}")
            print(f"📄 Maximum: {stats.max}")
            print(f"📄 Valeurs nulles: {stats.null_count}")
            if stats.distinct_count:
                print(f"📄 Valeurs distinctes: {stats.distinct_count}")

=====
📄 STATISTIQUES PAR COLONNE

📄 VendorID:
Type Arrow: int32
Type physique: INT32
Compression: SNAPPY
Taille compressée: 9,871 bytes
Taille non-compressée: 11,538 bytes
Minimum: 1
Maximum: 7
Valeurs nulles: 0

📄 tpep_pickup_datetime:
Type Arrow: timestamp[us]
Type physique: INT64
Compression: SNAPPY
Taille compressée: 284,866 bytes
Taille non-compressée: 339,850 bytes
Minimum: 2024-12-31 20:47:05
Maximum: 2025-01-01 01:02:10
Valeurs nulles: 0

📄 tpep_dropoff_datetime:
Type Arrow: timestamp[us]
Type physique: INT64
Compression: SNAPPY
Taille compressée: 284,866 bytes
Taille non-compressée: 339,850 bytes
Minimum: 2024-12-31 20:47:05
Maximum: 2025-01-01 01:02:10
Valeurs nulles: 0

📄 passenger_count:
Type Arrow: double
Type physique: DOUBLE
Compression: SNAPPY
Taille compressée: 18,185 bytes
Taille non-compressée: 24,685 bytes
Minimum: -0.0
Maximum: 9.0
Valeurs nulles: 0

📄 trip_distance:
Type Arrow: double
Type physique: DOUBLE
Compression: SNAPPY
Taille compressée: 84,737 bytes
Taille non-compressée: 93,499 bytes
Minimum: -0.0
Maximum: 133.3
Valeurs nulles: 0
```

Résumé

```
In [11]: print("📄 * * * * 60")
print("📄 RÉSUMÉ")
print("📄 * * * * 60")

print("📄 AVANTAGES PARQUET DÉMONSTRÉS:")
print("📄 Compression: (csv_size/parquet_snappy_size:16x plus compact que CSV)")
print("📄 Lecture: (csv_read_time/parquet_read_time:16x plus rapide que CSV)")
print("📄 Filtrage: (csv_filter_time/parquet_filter_time:16x plus rapide avec predicate)")
print("📄 Colonnes: (csv_cols_time/parquet_cols_time:16x plus rapide pour lecture sélective)")
print("📄 Métadonnées: schema automatique, statistiques, compression")

print("📄 Démonstration terminée!")
print("📄 * * * * 60")

=====
📄 RÉSUMÉ

📄 AVANTAGES PARQUET DÉMONSTRÉS:
📄 Compression: 5.0x plus compact que CSV
📄 Lecture: 9.1x plus rapide que CSV
📄 Filtrage: 12.2x plus rapide avec predicate
📄 Colonnes: 13.0x plus rapide pour lecture sélective
📄 Métadonnées: schema automatique, statistiques, compression
```

