

▼ Introduction :

This notebook performs data cleaning, integration, feature engineering, and preparation on a large historical stock price dataset combined with company metadata. The goal is to create a clean, structured, feature-rich dataset ready for time-series modeling.

```
1 # --- 1. IMPORT LIBRARIES ---
2 import pandas as pd
3 import numpy as np
4 from sklearn.preprocessing import StandardScaler
5 import gc
6
```

```
1 # Load metadata
2 meta_df = pd.read_csv('/content/historical_stocks.csv') # the metadata raw file
3
4 print("Missing values BEFORE cleaning:\n")
5 print(meta_df.isna().sum(), "\n")
6
7 # 1. Drop rows with no ticker (ticker is primary key)
8 meta_df = meta_df.dropna(subset=['ticker'])
9
10 # 2. Normalize sector & industry
11 meta_df['sector'] = meta_df['sector'].str.strip().str.upper()
12 meta_df['industry'] = meta_df['industry'].str.strip().str.upper()
13
14 # 3. Fill missing text values
15 fill_values = {
16     'exchange': 'UNKNOWN',
17     'sector': 'UNKNOWN',
18     'industry': 'UNKNOWN'
19 }
20 meta_df = meta_df.fillna(fill_values)
21
22 # 4. Deduplicate tickers
23 meta_df = meta_df.drop_duplicates(subset='ticker', keep='first')
24
25 # 5. Keep only essential cols
26 meta_df = meta_df[['ticker', 'exchange', 'sector', 'industry']]
27
28 print("Missing values AFTER cleaning:\n")
29 print(meta_df.isna().sum())
30 print("\nMetadata cleaned successfully!")
31
32 gc.collect() # garbage collector - very useful for collecting garbage from colab.
33
```

Missing values BEFORE cleaning:

```
ticker      0
exchange    0
name        0
sector     1440
industry    1440
dtype: int64
```

Missing values AFTER cleaning:

```
ticker      0
exchange    0
sector      0
industry    0
dtype: int64
```

Metadata cleaned successfully!
0

```
1 hist_df = pd.read_csv('/content/historical_stock_prices.csv') # the stock prices raw dataset
2
3 # Convert date to datetime
4 hist_df['date'] = pd.to_datetime(hist_df['date'])
5
6 # 1. Sort by ticker + date
7 hist_df = hist_df.sort_values(['ticker', 'date'])
```

```

8
9 # 2. Interpolate OHLCV per ticker
10 def fast_interpolation(group):
11     group = group.sort_values('date')
12     numeric_cols = ['open', 'close', 'adj_close', 'low', 'high', 'volume']
13     group[numeric_cols] = (
14         group[numeric_cols]
15         .interpolate(method='linear')
16         .bfill()
17         .ffill()
18     )
19     return group
20
21 hist_df_cleaned = hist_df.groupby('ticker', group_keys=False).apply(fast_interpolation)
22
23 print("Missing values AFTER interpolation:\n")
24 print(hist_df_cleaned.isna().sum())
25 gc.collect()
26

```

/tmp/ipython-input-1863111219.py:21: DeprecationWarning: DataFrameGroupBy.apply operated on the grouping columns. This behavior
hist_df_cleaned = hist_df.groupby('ticker', group_keys=False).apply(fast_interpolation)
Missing values AFTER interpolation:

```

ticker      0
open        1
close       1
adj_close   1
low         1
high        1
volume      1
date        1
dtype: int64
0

```

```

1 # Merge Dataset to AVOID CRASH!!!
2
3 merged_df = hist_df_cleaned.merge(meta_df, on='ticker', how='left') # this syntax merged the both files at this point becau
4
5 print("Merged dataset shape:", merged_df.shape)
6 merged_df.head()
7 gc.collect()

```

```

Merged dataset shape: (5246183, 11)
0

```

```

1 # Save the file to Parquet:
2
3 # Note for me: Parquet is more efficient than CSV. Becasue of so manay crashes that experience, I came across parquet. It pr
4 # category, int, datetime, unlike csv that do not preserve them. It is good for large dataset and reduce memory capacity (RA
5
6 merged_df.to_parquet('/content/merged_df.parquet', compression='snappy')
7 print("Saved merged_df.parquet!")
8 gc.collect()

```

```

Saved merged_df.parquet!
0

```

```

1
2 # Download the merged file and save it to avoid data loss
3
4 from google.colab import files
5 files.download('/content/merged_df.parquet')
6

```

```

1 #import gc - this is use to collect the grabage in the system that takes up memory/space (RAM)
2 gc.collect()
3

```

```

22

```

```

1 # Reload the merged dataset for use in the FE. It holds all of the files in one place ( prices and metadata)
2
3 merged_df = pd.read_parquet('/content/merged_df.parquet')

```

```

4 print("Reloaded merged_df shape:", merged_df.shape)
5
6 gc.collect()
7

```

Reloaded merged_df shape: (5246183, 11)
0

```

1 # Feature Engineering (Smooths price movement by averaging over past days – reduces noise.)
2
3 # Group by ticker so each stock is processed separately
4 grp = merged_df.groupby('ticker')['close']
5
6 # 7-day Moving Average – short-term trend
7 merged_df['MA_7'] = grp.transform(lambda x: x.rolling(7).mean())
8
9 # 14-day Moving Average – medium-term trend
10 merged_df['MA_14'] = grp.transform(lambda x: x.rolling(14).mean())
11
12 # 21-day Moving Average – long-ish trend
13 merged_df['MA_21'] = grp.transform(lambda x: x.rolling(21).mean())
14
15 gc.collect()
16 print("Moving averages completed!")
17

```

✓ Moving averages completed!

```

1 # Rolling Standard Deviations (Volatility) Measures how much the price fluctuates – higher std = higher volatility.
2
3 # 7-day volatility
4 merged_df['std_7'] = grp.transform(lambda x: x.rolling(7).std())
5
6 # 14-day volatility
7 merged_df['std_14'] = grp.transform(lambda x: x.rolling(14).std())
8
9 # 30-day volatility (longer-term uncertainty)
10 merged_df['std_30'] = grp.transform(lambda x: x.rolling(30).std())
11
12 gc.collect()
13 print("STD features completed!")
14

```

STD features completed!

```

1 # Price Momentum : Measures whether the stock is rising or falling compared to 1 week ago.
2
3 # Momentum = today's close - close price 7 days ago
4 merged_df['Price_Momentum'] = grp.transform(lambda x: x - x.shift(7))
5
6 gc.collect()
7 print("Price Momentum completed!")
8

```

✓ Price Momentum completed!

```

1 # RSI (Relative Strength Index: 14-day) RSI tells whether a stock is overbought (>70) or oversold (<30). Calculated using a
2
3 # 1. Price change between days
4 merged_df['delta'] = grp.diff()
5
6 # 2. Positive changes (gains)
7 merged_df['gain'] = merged_df['delta'].clip(lower=0)
8
9 # 3. Negative changes (losses)
10 merged_df['loss'] = -merged_df['delta'].clip(upper=0)
11
12 # 4. 14-period rolling average gain
13 merged_df['avg_gain'] = merged_df.groupby('ticker')['gain'].transform(
14     lambda x: x.rolling(14, min_periods=1).mean()
15 )
16
17 # 5. 14-period rolling average loss
18 merged_df['avg_loss'] = merged_df.groupby('ticker')['loss'].transform(
19     lambda x: x.rolling(14, min_periods=1).mean()

```

```

20 )
21
22 # 6. RS = avg gain / avg loss
23 merged_df['RS'] = merged_df['avg_gain'] / merged_df['avg_loss']
24
25 # 7. RSI formula
26 merged_df['RSI14'] = 100 - (100 / (1 + merged_df['RS']))
27
28 gc.collect()
29 print("RSI completed!")
30
31 # 8. Drop temporary columns to reduce RAM
32 merged_df.drop(['delta', 'gain', 'loss', 'avg_gain', 'avg_loss', 'RS'],
33                axis=1, inplace=True)
34
35 gc.collect()
36

```

✓ RSI completed!
0

```

1 # MACD and Signal Line: MACD is a widely used trend indicator.
2 #It compares fast (12-day) and slow (26-day) EMAs.
3 #Signal line is the 9-day EMA of MACD.
4
5 # 12-day Exponential Moving Average (fast)
6 EMA12 = grp.transform(lambda x: x.ewm(span=12, adjust=False).mean())
7
8 # 26-day Exponential Moving Average (slow)
9 EMA26 = grp.transform(lambda x: x.ewm(span=26, adjust=False).mean())
10
11 # MACD = fast EMA - slow EMA
12 merged_df['MACD'] = EMA12 - EMA26
13
14 # MACD Signal = 9-day EMA of MACD
15 merged_df['MACD_signal'] = merged_df.groupby('ticker')['MACD'].transform(
16     lambda x: x.ewm(span=9, adjust=False).mean()
17 )
18
19 # Free memory – no longer needed
20 del EMA12, EMA26
21 gc.collect()
22
23 print("MACD completed!")
24

```

MACD completed!

```

1 #Drop metadata columns to reduce size and avoid crashes
2 cols_to_drop = ['exchange', 'sector', 'industry']
3
4 merged_df.drop(columns=cols_to_drop, inplace=True, errors='ignore')
5
6 print("Dropped metadata columns. New shape:", merged_df.shape)
7

```

Dropped metadata columns. New shape: (5246183, 18)

```

1 # Save FE Dataset to Parquet
2
3 # Save this file to Parquet and avoid data loss during this project:
4
5 merged_df.to_parquet('/content/merged_df_FE.parquet', compression='snappy')
6
7
8 # Download the file to avoid data loss
9
10 files.download('/content/merged_df_FE.parquet')
11
12 print("Feature-engineered dataset saved & Downloaded!")
13 gc.collect()

```

Feature-engineered dataset saved & Downloaded!
118

```

1 # Ensure correct chronological order
2 merged_df = merged_df.sort_values(['ticker', 'date']).reset_index(drop=True)
3
4 print("Sorted dataset shape:", merged_df.shape)
5

```

Sorted dataset shape: (5246183, 18)

```

1 # Create Train / Validation / Test Splits
2
3 train_size = int(len(merged_df) * 0.70)
4 val_size = int(len(merged_df) * 0.15)
5
6 train_df = merged_df.iloc[:train_size]
7 val_df = merged_df.iloc[train_size : train_size + val_size]
8 test_df = merged_df.iloc[train_size + val_size :]
9
10 print("Train:", train_df.shape)
11 print("Validation:", val_df.shape)
12 print("Test:", test_df.shape)
13

```

Train: (3672328, 18)
 Validation: (786927, 18)
 Test: (786928, 18)

```

1 # Save the Train/Val/Test Dataset
2
3 # Save the train/val/test datasets to Parquet files
4 train_df.to_parquet('/content/train_df.parquet', compression='snappy')
5 val_df.to_parquet('/content/val_df.parquet', compression='snappy')
6 test_df.to_parquet('/content/test_df.parquet', compression='snappy')
7
8 # Download each file
9 files.download('/content/train_df.parquet')
10 files.download('/content/val_df.parquet')
11 files.download('/content/test_df.parquet')
12
13 print("All splits saved & Downloaded - James !")
14

```

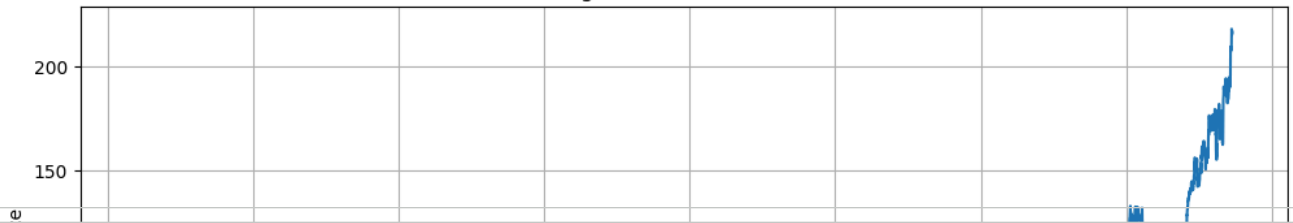
All splits saved & Downloaded - James !

```

1 #Closing Price Over Time (Per Ticker)
2
3 import matplotlib.pyplot as plt
4
5 # Select a sample ticker
6 sample_ticker = merged_df['ticker'].unique()[0]
7
8 df_sample = merged_df[merged_df['ticker'] == sample_ticker]
9
10 plt.figure(figsize=(12,5))
11 plt.plot(df_sample['date'], df_sample['close'])
12 plt.title(f"Closing Price Over Time - {sample_ticker}")
13 plt.xlabel("Date")
14 plt.ylabel("Close Price")
15 plt.grid(True)
16 plt.show()
17

```

Closing Price Over Time - AAPL



```

1 # Moving Averages Overlay Chart (Trend Detection)
2
3 plt.figure(figsize=(12,5))
4 plt.plot(df_sample['date'], df_sample['close'], label='Close Price', alpha=0.6)
5 plt.plot(df_sample['date'], df_sample['MA_7'], label='MA 7', linewidth=2)
6 plt.plot(df_sample['date'], df_sample['MA_14'], label='MA 14', linewidth=2)
7 plt.plot(df_sample['date'], df_sample['MA_21'], label='MA 21', linewidth=2)
8
9 plt.title(f"Moving Averages Trend - {sample_ticker}")
10 plt.xlabel("Date")
11 plt.ylabel("Price")
12 plt.legend()
13 plt.grid(True)
14 plt.show()
15

```

Moving Averages Trend - AAPL



```

1 # Moving Averages Overlay Chart (Trend Detection)
2
3 plt.figure(figsize=(12,5))
4 plt.plot(df_sample['date'], df_sample['close'], label='Close Price', alpha=0.6)
5 plt.plot(df_sample['date'], df_sample['MA_7'], label='MA 7', linewidth=2)
6 plt.plot(df_sample['date'], df_sample['MA_14'], label='MA 14', linewidth=2)
7 plt.plot(df_sample['date'], df_sample['MA_21'], label='MA 21', linewidth=2)
8
9 plt.title(f"Moving Averages Trend - {sample_ticker}")
10 plt.xlabel("Date")
11 plt.ylabel("Price")
12 plt.legend()
13 plt.grid(True)
14 plt.show()
15

```

