# # [ Data Cleaning ] {CheatSheet}

# 1. Handling Missing Values

- Identify Missing Values: df.isnull().sum()
- Drop Rows with Missing Values: df.dropna()
- Drop Columns with Missing Values: df.dropna(axis=1)
- Fill Missing Values with a Constant: df.fillna(value)
- Fill Missing Values with Mean/Median/Mode: df.fillna(df.mean())
- Forward Fill Missing Values: df.ffill()
- Backward Fill Missing Values: df.bfill()
- Interpolate Missing Values: df.interpolate()

# 2. Data Type Conversions

- Convert Data Type of a Column: df['col'] = df['col'].astype('type')
- Convert to Numeric: pd.to\_numeric(df['col'], errors='coerce')
- Convert to Datetime: pd.to\_datetime(df['col'], errors='coerce')
- Convert to Categorical: df['col'] = df['col'].astype('category')

# 3. Dealing with Duplicates

- Identify Duplicate Rows: df.duplicated()
- Drop Duplicate Rows: df.drop\_duplicates()
- Drop Duplicates in a Specific Column: df.drop\_duplicates(subset='col')
- Drop Duplicates Keeping the Last Occurrence: df.drop\_duplicates(keep='last')

# 4. Text Data Cleaning

- Trim Whitespace: df['col'] = df['col'].str.strip()
- Convert to Lowercase: df['col'] = df['col'].str.lower()
- Convert to Uppercase: df['col'] = df['col'].str.upper()
- Remove Specific Characters: df['col'] = df['col'].str.replace('[character]', '')

- Replace Text Based on Pattern (Regex): df['col'] = df['col'].str.replace(r'[regex]', 'replacement')
- Split Text into Columns: df[['col1', 'col2']] = df['col'].str.split(',', expand=True)

# 5. Categorical Data Processing

- One-Hot Encoding: pd.get\_dummies(df['col'])
- Label Encoding: from sklearn.preprocessing import LabelEncoder; encoder = LabelEncoder(); df['col'] = encoder.fit\_transform(df['col'])
- Map Categories to Values: df['col'] = df['col'].map({'cat1': 1, 'cat2': 2})
- Convert Category to Ordinal: df['col'] = df['col'].cat.codes

# 6. Normalization and Scaling

- Min-Max Scaling: from sklearn.preprocessing import MinMaxScaler; scaler = MinMaxScaler(); df['col'] = scaler.fit\_transform(df[['col']])
- Standard Scaling (Z-Score): from sklearn.preprocessing import StandardScaler; scaler = StandardScaler(); df['col'] = scaler.fit\_transform(df[['col']])
- Robust Scaling (Median, IQR): from sklearn.preprocessing import RobustScaler; scaler = RobustScaler(); df['col'] = scaler.fit\_transform(df[['col']])

# 7. Handling Outliers

- Remove Outliers with IQR: Q1 = df['col'].quantile(0.25); Q3 =  $df['col'].quantile(0.75); IQR = Q3 - Q1; df = df[\sim((df['col'] < (Q1))]$ -1.5 \* IQR)) | (df['col'] > (Q3 + 1.5 \* IQR)))]
- Remove Outliers with Z-Score: from scipy import stats; df = df[np.abs(stats.zscore(df['col'])) < 3]</pre>
- Capping and Flooring Outliers: df['col'] = df['col'].clip(lower=lower\_bound, upper=upper\_bound)

#### 8. Data Transformation

- Log Transformation: df['col'] = np.log(df['col'])
- Square Root Transformation: df['col'] = np.sqrt(df['col'])
- Power Transformation (Box-Cox, Yeo-Johnson): from sklearn.preprocessing import PowerTransformer; pt = PowerTransformer(method='yeo-johnson'); df['col'] = pt.fit\_transform(df[['col']])
- Binning Data: df['bin\_col'] = pd.cut(df['col'], bins=[range])

### 9. Time Series Data Cleaning

- **Set Datetime Index**: df.set\_index('datetime\_col', inplace=True)
- Resample Time Series Data: df.resample('D').mean()
- Fill Missing Time Series Data: df.asfreq('D', method='ffill')
- Time-Based Filtering: df['year'] = df.index.year; df[df['year'] > 2000]

### 10. Data Frame Operations

- Merge Data Frames: pd.merge(df1, df2, on='key', how='inner')
- Concatenate Data Frames: pd.concat([df1, df2], axis=0)
- Join Data Frames: df1.join(df2, on='key')
- Pivot Table: df.pivot\_table(index='row', columns='col', values='value')

# 11. Column Operations

- Aggregate Functions (sum, mean, etc.): df.groupby('group\_col').agg({'agg\_col': ['sum', 'mean']})
- Rolling Window Calculations: df['col'].rolling(window=5).mean()
- Expanding Window Calculations: df['col'].expanding().sum()

# 12. Handling Complex Data Types

- Explode List to Rows: df.explode('list\_col')
- Work with JSON Columns: df['json\_col'].apply(lambda x: json.loads(x))

• Parse Nested Structures: df['new\_col'] = df['struct\_col'].apply(lambda x: x['nested\_field'])

# 13. Dealing with Geospatial Data

- Handling Latitude and Longitude: df['distance'] = df.apply(lambda x: calculate\_distance(x['lat'], x['long']), axis=1)
- **Geocoding Addresses**: df['coordinates'] = df['address'].apply(geocode\_address)

# 14. Data Quality Checks

- Check for Data Consistency: assert df['col1'].notnull().all()
- Validate Data Ranges: df[(df['col'] >= low\_val) & (df['col'] <= high\_val)]
- Assert Data Types: assert df['col'].dtype == 'expected\_type'

### 15. Efficient Computations

- Use Vectorized Operations: df['col'] = df['col1'] + df['col2']
- Parallel Processing with Dask: import dask.dataframe as dd; ddf = dd.from\_pandas(df, npartitions=10); result = ddf.compute()

# 16. Working with Large Datasets

- Sampling Data for Quick Insights: sampled\_df = df.sample(frac=0.1)
- Chunking Large Files for Processing: for chunk in pd.read\_csv('large\_file.csv', chunksize=10000): process(chunk)

# 17. Feature Engineering

- Creating Polynomial Features: from sklearn.preprocessing import PolynomialFeatures; poly = PolynomialFeatures(degree=2); df\_poly = poly.fit\_transform(df[['col1', 'col2']])
- Encoding Cyclical Features (e.g., hour of day, day of week):  $df['hour_sin'] = np.sin(df['hour'] * (2 * np.pi / 24))$

### 18. Data Imputation

- Impute Missing Values with KNN: from sklearn.impute import KNNImputer; imputer = KNNImputer(n\_neighbors=5); df['col'] = imputer.fit\_transform(df[['col']])
- Iterative Imputation: from sklearn.experimental import enable\_iterative\_imputer; from sklearn.impute import IterativeImputer; imputer = IterativeImputer(); df\_imputed = imputer.fit\_transform(df)

#### 19. Data Validation

- Using Pandera for Schema Validation: import pandera as pa; schema = pa.DataFrameSchema({'col': pa.Column(pa.Int, nullable=False)}); schema.validate(df)
- Validating Range of Values: df['col'].between(low\_value, high\_value)

### 20. Data Anonymization

- Hashing Sensitive Data: df['hashed\_col'] = df['sensitive\_col'].apply(lambda x: hash\_function(x))
- Randomized Noise Addition: df['col'] = df['col'] + np.random.normal(0, 1, df.shape[0])
- Masking Values: df['col'] = df['col'].apply(lambda x: x[:3] + '\*\*\*')

# 21. Data Integration and Alignment

- Aligning Columns from Different DataFrames: df1, df2 = df1.align(df2, join='inner', axis=1)
- Combining Data from Multiple Sources: df\_combined = pd.merge(df1, df2, on='common\_key')

# 22. String Operations and Regular Expressions

Extracting Substrings with Regex: df['extracted'] = df['text\_col'].str.extract(r'(pattern)')

Removing Unwanted Characters: df['clean\_text'] = df['text'].str.replace('[^\w\s]', '', regex=True)

# 23. Handling Time and Date

- Extracting Date Components: df['year'] = df['date\_col'].dt.year
- Calculating Date Differences: df['days\_diff'] = (df['date\_col1'] df['date\_col2']).dt.days
- Date Range Generation for Time Series: pd.date\_range(start='2020-01-01', end='2020-12-31', freq='D')

# 24. Working with Indexes

- Resetting Index: df.reset\_index(drop=True, inplace=True)
- **Setting a Column as Index**: df.set\_index('col', inplace=True)
- Reindexing with α New Index: df.reindex(new\_index)

### 25. Data Compression and Memory Management

- Reducing Memory Usage by Changing Data Types: df['int\_col'] = df['int\_col'].astype('int32')
- Compressing DataFrame using Categories: df['cat\_col'] = df['cat\_col'].astype('category')

# 26. Handling Large and Sparse Data

- Working with Sparse Data Structures: from scipy.sparse import csr\_matrix; sparse\_matrix = csr\_matrix(df)
- Efficiently Storing Large Data with HDF5: df.to\_hdf('data.h5', key='df', mode='w')

#### 27. Data Randomization

- Shuffling Rows Randomly: df = df.sample(frac=1).reset\_index(drop=True)
- Generating Random Samples: df\_sample = df.sample(n=100)

#### 28. Feature Extraction

- Extracting Features from Text: from sklearn.feature\_extraction.text import CountVectorizer: vectorizer = CountVectorizer(); X = vectorizer.fit\_transform(df['text\_col'])
- Dimensionality Reduction (e.g., PCA): from sklearn.decomposition import PCA; pca = PCA(n\_components=2); df\_reduced = pca.fit\_transform(df)

# 29. Combining Data

- Appending Rows of Another DαtαFrame: df = df.append(other\_df)
- Concatenating DataFrames Vertically or Horizontally: pd.concat([df1, df2], axis=0)

### 30. Data Cleaning Automation

- Using Clean Function from CleanPandas: from cleanpandas import clean; df = clean(df)
- Automated Data Cleaning with DataCleaner: from datacleaner import autoclean; df = autoclean(df)

### 31. Handling Numerical Data

- Rounding Numeric Columns: df['col'] = df['col'].round(decimals=2)
- Discretizing Continuous Variables: df['binned\_col'] = pd.qcut(df['col'], q=4)

#### 32. Geospatial Data Processing

- Coordinate Transformation: df['x'], df['y'] = zip(\*df['coordinates'].apply(transform\_coord))
- Distance Calculation Between Coordinates: df['distance'] = df.apply(lambda row: calc\_distance(row['lat1'], row['lon1'], row['lat2'], row['lon2']), axis=1)

### 33. Multilingual and Locale-Specific Operations

- Converting Currencies or Units: df['converted\_col'] = df['amount'].apply(convert\_currency)
- Locale-Specific Sorting: df.sort\_values(by='name', key=lambda col: col.str.normalize('NFKD'))

# 34. Advanced DataFrame Manipulations

- Pivoting and Unpivoting Data: df.pivot(index='date', columns='variable', values='value')
- Stacking and Unstacking Data: df.stack(); df.unstack()

# **35.** Custom Cleaning Functions

- Applying Custom Cleaning Functions: df['clean\_col'] = df['col'].apply(custom\_clean\_function)
- Using Lambda Functions for Quick Cleaning: df['processed\_col'] = df['col'].apply(lambda x: x.strip().lower())