



Master Data Science & Pandas EDA Cheatsheet

1. Data Loading

Task	Function / Notes
Read CSV	<code>df = pd.read_csv('filename.csv')</code>
Read Excel	<code>df = pd.read_excel('filename.xlsx')</code>
Read SQL	<code>df = pd.read_sql(query, connection)</code>
Read JSON	<code>df = pd.read_json('filename.json')</code>
Read XML	<code>df = pd.read_xml('filename.xml')</code>
Read HTML Table	<code>dfs = pd.read_html('url')</code>
Chunking Large Files	<code>pd.read_csv('file.csv', chunksize=1000)</code>

2. Basic Data Inspection

Task	Function / Notes
Top / Bottom rows	<code>df.head()</code> <code>df.tail()</code>
Data Types	<code>df.dtypes</code>
Summary Statistics	<code>df.describe()</code>
Info / Columns / Index	<code>df.info()</code> <code>df.columns</code> <code>df.index</code>
Check Missing	<code>df.isna().sum()</code> <code>df.isnull().sum()</code>
Unique / Counts	<code>df['col'].unique()</code> <code>df['col'].nunique()</code> <code>df['col'].value_counts()</code>

Memory Usage	<code>df.memory_usage(deep=True)</code>
--------------	---

3. Data Cleaning & Missing Data

Task	Function / Notes
Drop Rows/Columns	<code>df.drop(columns=['col'])</code> <code>df.dropna()</code>
Fill Missing	<code>df.fillna(value)</code> <code>df['col'].ffill()</code> <code>df['col'].bfill()</code> <code>df['col'].interpolate()</code>
Fill by Group	<code>df['col'] = df.groupby('group')['col'].transform(lambda x: x.fillna(x.median()))</code>
Conditional Fill	<code>df['col'] = df['col'].apply(lambda x: 0 if pd.isna(x) else x)</code>
Rename Columns	<code>df.rename(columns={'old':'new'})</code>
Remove Duplicates	<code>df.drop_duplicates(subset=['col1','col2'], keep='first')</code>
Assert / Validation	<code>assert df.notnull().all().all(), "Missing values!"</code>

4. Data Transformation

Task	Function / Notes
Apply Function / Custom	<code>df['col'].apply(lambda x: func(x))</code>
Apply Across Rows	<code>df.apply(lambda row: func(row['col1'], row['col2'])), axis=1)</code>
GroupBy & Aggregate	<code>df.groupby('col').agg({'col2':'sum'})</code>
Multiple Aggregate	<code>df.groupby('col').agg(['mean', 'sum'])</code>

Transform Function	<code>df.groupby('col').transform(lambda x: x-x.mean())</code>
Pivot Table	<code>df.pivot_table(index='col1', columns='col2', values='col3', aggfunc='sum', fill_value=0)</code>
Melt / Stack / Unstack	<code>df.melt(id_vars=['col1'])</code> <code>df.stack()</code> <code>df.unstack()</code>
Cross-tab	<code>pd.crosstab(df['col1'], df['col2'])</code>
Merge / Join	<code>pd.merge(df1, df2, on='col', how='left')</code>
Concatenate	<code>pd.concat([df1, df2], axis=0)</code>

5. Indexing & Selection

Task	Function / Notes
Select Columns	<code>df['col']</code> <code>df[['col1', 'col2']]</code>
Select Rows by Position	<code>df.iloc[0:5]</code>
Select Rows by Label	<code>df.loc[0:5]</code>
Conditional Selection	<code>df[df['col']>value]</code>
Query Function	<code>df.query('col>value')</code>
Filtering with isin	<code>df[df['col'].isin([v1, v2])]</code>
Reset / Set Index	<code>df.reset_index(drop=True)</code> <code>df.set_index('col')</code>
MultiIndex	<code>df.set_index(['col1', 'col2'])</code> <code>df.xs(key, level='level_name')</code> <code>df.loc[(slice(start1, end1), slice(start2, end2)), :]</code>

6. Data Formatting & Conversion

Task	Function / Notes
Convert Type	<code>df['col'].astype('type')</code>
Datetime Conversion	<code>pd.to_datetime(df['col'], format='%Y-%m-%d', errors='coerce')</code>
String Operations	<code>df['col'].str.lower()</code> <code>df['col'].str.upper()</code> <code>df['col'].str.split()</code> <code>df['col'].str.contains()</code> <code>df['col'].str.extract(r'(regex)')</code>
Categorical	<code>df['col'].astype('category')</code> <code>df['col'].cat.set_categories(['cat1', 'cat2'], ordered=True)</code>
Numeric Conversion	<code>df['col'] = pd.to_numeric(df['col'], errors='coerce')</code>
Mixed-Type Cleanup	<code>df['col'] = df['col'].astype(str).str.replace('\$', '').astype(float)</code>

7. Handling Time Series

Task	Function / Notes
Set Datetime Index	<code>df.set_index(pd.to_datetime(df['date_col']))</code>
Resample	<code>df.resample('frequency').agg_func()</code> 'D' → Daily 'W' → Weekly 'M' → Monthly 'Q' → Quarterly 'Y' → Yearly <code>agg_func() => sum(), mean(), max()</code>
Rolling Operations	<code>df.rolling(window=5).mean()</code>
Lag / Shift	<code>df['col'].shift(1)</code>
Time-Based Groupby	<code>df.groupby(pd.Grouper(key='date_col', freq='M')).sum()</code>

Missing Timestamps	<code>df = df.asfreq('D')</code>
--------------------	----------------------------------

8. Numeric Transformations / Normalization

Task	Function / Notes
Min-Max Normalization	<code>(df['col']-df['col'].min())/(df['col'].max()-df['col'].min())</code>
Z-score Standardization	<code>(df['col']-df['col'].mean())/df['col'].std()</code>
Round Values	<code>df['col'].round(2)</code>
Aggregate Functions	<code>df.sum() df.mean() df.median() df.std()</code>

9. Visualization

Task	Function / Notes
Histogram	<code>df['col'].hist(bins=20) sns.histplot(df['col'])</code>
Boxplot	<code>df.boxplot(column=['col1', 'col2']) sns.boxplot(x='cat', y='num', data=df)</code>
Scatter	<code>df.plot.scatter(x='col1', y='col2') sns.scatterplot(x='col1', y='col2')</code>
Line	<code>df.plot.line() sns.lineplot(x='x', y='y')</code>
Bar	<code>df['col'].value_counts().plot.bar() sns.barplot(data=df, x='col1', y='col2', hue='col3')</code>
Pairplot / Correlation	<code>sns.pairplot(df) sns.heatmap(df.corr(), annot=True)</code>

Matplotlib Customization	<code>plt.style.use('ggplot')</code> <code>plt.show()</code>
-----------------------------	---

10. File Export

Task	Function / Notes
CSV	<code>df.to_csv('filename.csv', index=False)</code>
Excel	<code>df.to_excel('filename.xlsx', index=False)</code>
JSON	<code>df.to_json('filename.json')</code>
SQL	<code>df.to_sql('table_name', connection)</code>

11. Performance & Large Data Handling

Task	Function / Notes
Reduce Memory	<code>df['col'] = pd.to_numeric(df['col'], downcast='float')</code>
Categorical Optimization	<code>df['cat_col'] = df['cat_col'].astype('category')</code>
Swifter	<code>df['col'].swifter.apply(lambda x: func(x))</code>
Dask for Parallel	<code>df = dd.from_pandas(df, npartitions=10)</code>
Sampling	<code>df.sample(n=1000)</code>

12. Statistical Analysis & SciPy / StatsModels

Task	Function / Notes
Correlation / Covariance	<code>df.corr()</code> <code>df.cov()</code>
Skewness / Kurtosis	<code>from scipy.stats import skew, kurtosis;</code> <code>skew(df['col']), kurtosis(df['col'])</code>
Linear Regression	<code>import statsmodels.api as sm;</code>

	<code>sm.OLS(y,X).fit()</code>
--	--------------------------------

13. Text / String Data

Task	Function / Notes
String contains	<code>df[df['col'].str.contains('text')]</code>
String split	<code>df['col'].str.split(' ', expand=True)</code>
Regex extract	<code>df['col'].str.extract(r'(regex)')</code>

14. Outlier Detection & Handling

Task	Function / Notes
IQR Method	<code>Q1, Q3 = df['col'].quantile([0.25,0.75]); IQR=Q3-Q1; df[(df['col']>=Q1-1.5*IQR) & (df['col']<=Q3+1.5*IQR)]</code>
Z-score	<code>from scipy.stats import zscore; df[(np.abs(zscore(df['col']))<3)]</code>