

Exploratory Data Analysis (EDA)

CHEAT SHEET

EVERY DATA ANALYSTS MUST KNOW

UNDERSTAND • EXPLORE • ANALYZE DATA



RAMA GOPALA KRISHNA MASANI

What is EDA?

Exploratory Data Analysis (EDA) is the process of:

- Understanding the structure of data
- Identifying patterns & trends
- Detecting missing values & outliers
- Preparing data for modeling & dashboards

👉 **EVERY DATA ANALYST MUST MASTER EDA.**

Import Required Libraries



```
1 import pandas as pd  
2 import numpy as np
```

Pandas is the core library for EDA in Python

Loading Data (All Common Sources)



```
1 pd.read_csv('file.csv')
2 pd.read_excel('file.xlsx')
3 pd.read_sql(query, connection)
4 pd.read_json('file.json')
5 pd.read_html(url)
```

First step in any EDA process.

Basic Data Overview



- 1 `df.head()`
- 2 `df.tail()`
- 3 `df.sample()`
- 4 `df.shape`
- 5 `df.info()`
- 6 `df.describe()`

Understand rows, columns, datatypes & summary statistics.

Understanding Columns & Data Types



- 1 `df.columns`
- 2 `df.dtypes`
- 3 `df.index`

Helps identify numerical, categorical & datetime columns.

Selecting Data

```
1 df['column']
2 df[['col1','col2']]
3 df.iloc[0]
4 df.loc[0]
5 df.iloc[0,0]
6 df.loc[0,'column']
```

Used for column & row-level analysis

Filtering Data



```
1 df[df['col'] > 5]
2 df.query('col > 5')
3 df.filter(like='sales')
```

Focus only on relevant data.

Handling Missing Values



- 1 df.isnull().sum()
- 2 df.notnull()
- 3 df.dropna()
- 4 df.fillna(value)

Missing data must be handled before analysis.

Data Cleaning Operations



```
1 df.drop_duplicates()  
2 df.rename(columns={'old': 'new'})  
3 df.astype('int')  
4 df.replace(1, 'One')  
5 df.reset_index()
```

Clean data = accurate insights.

Sorting & Sampling



```
1 df.sort_values('col')
2 df.sort_values('col', ascending=False)
3 df.sample(5)
4 df.nlargest(3,'col')
5 df.nsmallest(3,'col')
```

Useful for ranking & trend analysis.

Aggregations & Statistics



```
1 df.mean()  
2 df.median()  
3 df.std()  
4 df.min()  
5 df.max()  
6 df.sum()  
7 df.corr()
```

Core statistical understanding of data.

GroupBy (Most Important for Analysts)



```
1 df.groupby('col').mean()  
2 df.groupby('col').sum()  
3 df.groupby('col').count()  
4 df.groupby('col')['sales'].max()
```

Used in almost every business analysis.

Pivot Tables

```
1 pd.pivot_table(  
2     df,  
3     values='sales',  
4     index='region',  
5     aggfunc='mean'  
6 )
```

Excel-like analysis using Python.

Merging & Combining Data



```
1 pd.concat([df1, df2])
2 df1.merge(df2, on='key', how='inner')
3 df1.merge(df2, on='key', how='left')
```

Required for multi-table analysis.

Apply & Transform



```
1 df.apply(np.mean)  
2 df.transform(lambda x: x + 10)
```

Custom column-wise operations.

Advanced Missing Value Analysis



```
1 df.isnull().sum()  
2 df.isnull().mean()*100
```

Identify columns with high missing percentage.

Outlier Detection



```
1 df.describe(percentiles=[.01,.05,.95,.99])
```

Helps detect extreme values before modeling.

Distribution Analysis



```
1 df['col'].value_counts()  
2 df['col'].nunique()
```

Understand frequency & cardinality.

Date & Time EDA



```
1 df['date'] = pd.to_datetime(df['date'])
2 df['year'] = df['date'].dt.year
3 df['month'] = df['date'].dt.month
```

Time-based trend analysis.

Categorical Data Analysis



```
1 df['category'].value_counts(normalize=True)
```

Essential for business insights.

Numerical vs Categorical Analysis



```
1 df.groupby('category')['sales'].mean()
```

Most common real-world business analysis.

Duplicate Analysis



```
1 df.duplicated().sum()
```

Prevents inflated metrics.

Feature Engineering



```
1 df['profit'] = df['revenue'] - df['cost']
```

Convert raw data into meaningful features.

Scaling Awareness



1 `df.describe()`

**Helps decide normalization or
standardization.**

Visualization (Quick Business EDA)



```
1 df.plot(kind='line')
2 df.plot(kind='bar')
3 df.plot(kind='box')
4 df.plot(kind='hist')
```

Fast insight before dashboards.

Thank You!

**Was this helpful?
Follow me for More**

