

Data Analytics:

```
df = pd.read_csv('df')
```

• Data Exploration

`df.shape / columns`
`df.describe()`
`df.info`

`df.isna().values.any()` > Missing /
`df.duplicated().values.any()` Duplicates
`* *`
`df.isna().sum()`
`df.duplicated().sum()` → duplicates = df[df.duplicated()]
 to see duplicates/missing

Loc & iloc
`df.loc` → Name
`df.iloc` → int position

`df.loc[1:row1 : 'rows']`
`filtering = df.loc[df['column'] > threshold]`
`df.iloc[0,0]` - First row, first column
`df.iloc[0:3]` - First 3 rows
`df.iloc[:, 2]` - Second column
`df.iloc[1:4, 0:2]` - Specific rows, specific columns

Sorting & Grouping - Summarizing

`df.sort_values("column", ascending=True)`
`df.value_counts()` - Grouping and counting only unique values
`df.groupby('group')`
`df.groupby('group').size()` - Counting total + Nan/
`.count()` - Counts all existing values
`agg()`

Querying

`df.query("column_1 > 2")` - From this column_1 to show only values bigger than 2
`ex-values = 5` ←
`df.query("column_1 > @ex-values")` OR `df = df[df['column_1'] < 5]`

Dropping

`clean = df.dropna() / df.dropna(inplace=True)`
`Clean = df.dropna(subset=['column_1', 'column_2'])` ↑ dropping in specific columns
`Clean = df.drop_duplicates() / df.drop_duplicates(inplace=True)`
`Clean = df.drop_duplicates(subset=['column_1', 'column_2'])`
`df = df.drop(column='...', axis=1)`

Filling Nan/ values → Median, Mean, Mode

`df.fillna(0, inplace=True)` / `df['column'] = df['column'].fillna(0)`
`df['numerical'] = df['numerical'].fillna(df['numerical'].mean())` Median
 mode() categorical data

Forward Fill / Backward Fill ← replacing with back/forward value

`df.fillna(method='ffill', inplace=True)`
`'ffill'`

Predicting

`imputer = KNNImputer(n_neighbors=3)`
`df_imputed = pd.DataFrame(imputer.fit_transform(df), columns=df.columns)`

Journey

Data
viz
Insights
Being

Concatenating Data

`df['date'] = pd.to_datetime(df['date'])` → `df['year'] = df['date'].dt.year / month / day`
`df['num'] = pd.to_numeric(df['num'])` → handling error: my data is f, 2, three, 4 I
`df['num'].errors='coerce'`
`df['category'] = df['category'].astype('category')`
`df['boolean'] = df['boolean'].map({True: 'Fals'})`
`df['price'] = df['price'].replace('$', '').replace(',', '').replace('.', '').astype(float)`

Outliness

- Take a look with box plot
`sns.boxplot(x=df['values'])`
`plt.title("Least outliers")`
`plt.show()`
- Take a look with histogram
`df['values'].hist(bins=20)`
`plt.title(...)`
`plt.show()`

Handling with Stars!

`Q1 = df['values'].quantile(0.25)`
`Q3 = df['values'].quantile(0.75)`
`IQR = Q3 - Q1`
`lower_bound = Q1 - 1.5 * IQR`
`upper_bound = Q3 + 1.5 * IQR`
`outliers = df[(df['values'] < lower_bound) | (df['values'] > upper_bound)]`
 ↳ remove or keep & called Cap or Clip or trim

Feature Engineering

- Creating new columns
`df['new'] = df['old1'] / df['old2']`
- Feature Extraction
`df['year'] = df['date'].dt.year`
- Feature transformation
 Not for now