

Data Analytics:

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

• Data Exploration

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

} General

`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

Name int position

```
df.loc['row1': 'row5']
```

filtering = `df.loc[df['column'] > threshold]`

`df.iloc[0,0]` - First row, First column

`df.iloc[0:3]` - First 3 rows

`df.iloc[:, 1]` - 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 data

`df.groupby('group')`

`df.groupby('group').size()` - Counting total + new
`.count()` & counts all existing values

`agg()` - Summarizing + aggregating.
`df.agg('mean')` ... for each column
`df.groupby('...').agg({ 'columns': 'column' })`

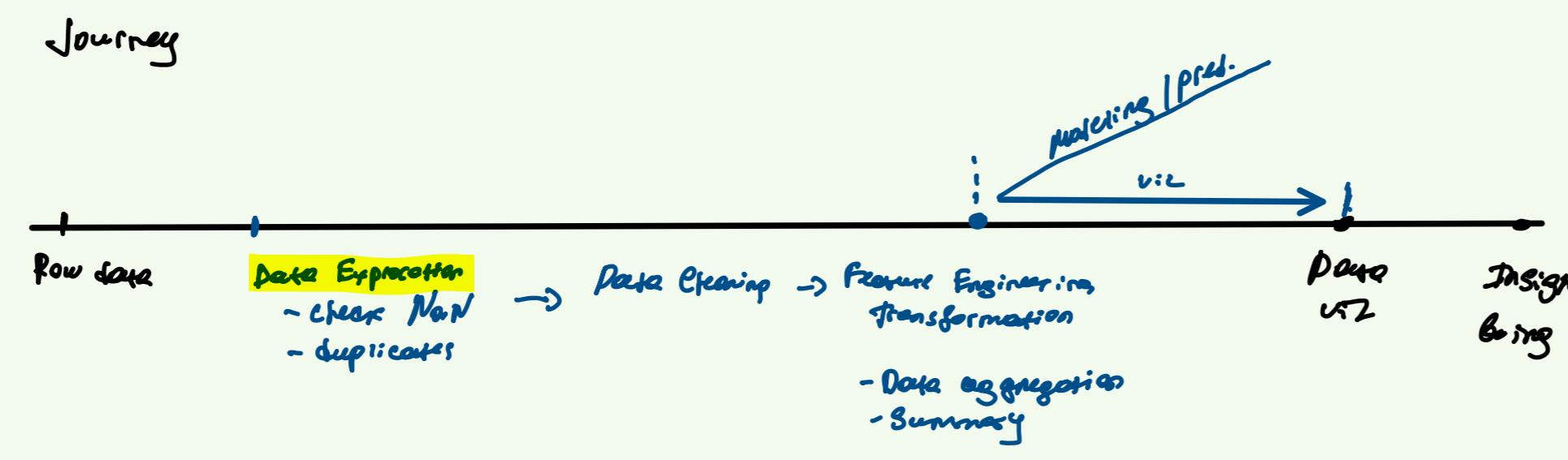
`df.pivot_table(values='columns', index='...', columns='...', aggfunc='mean')` ⇒ better for Deep learning

Browsing

`df.query("column_1 > 2")` - From this column_1 → show only values bigger than 2

`ex_values = 5` ← OR `df = df[df['column_1'] < 5]`

`df.query("column_1 > ex_values")`

**Handling**

`clean = df.dropna() / df.dropna(inplace=True)`
`clean = df.dropna(subset=['column_1', 'column_2'])`
`dropna 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 Null

`df.fillna(0, inplace=True) / df['column'] = df['column'].fillna(0)`
`df['numerical_col'] = df['numerical_col'].fillna(df['numerical_col'].mean())`

`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)`

Concatenating

`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 date is 2, 2, three, 4`
`df['num'] = df['num'].replace('...', errors='coerce')`

`df['category'] = df['category'].astype('category')`

`df['boolean'] = df['boolean'].map({True: True, False: False})`

`df['price'] = df['price'].replace('E\$', '$', regex=True).astype(float)`

Outliers

• Take a look with box plot

`sns.boxplot(x=df['values'].quantile(0.25))`
`Q1 = df['values'].quantile(0.25)`
`Q3 = df['values'].quantile(0.75)`

• Take a look with histogram

`df['values'].hist(bins=20)`
`plt.title(...)`
`plt.show()`

Handling with Stats!

`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 Clipping or`

Feature Engineering

- Creating new columns

`df['new'] = df['old1'] / df['old2']`

- Feature Extraction

`df['year'] = df['date'].dt.year`

- Feature Transformation

Not for now

Merging

`merged_df = pd.merge(df1, df2, on='id')`

structure pretty much same

Plotting - Dynamic**SEABORN**

```

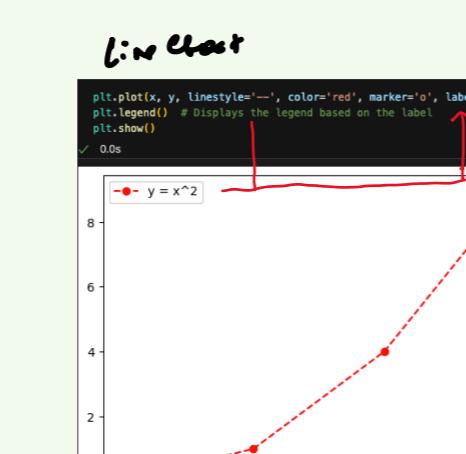
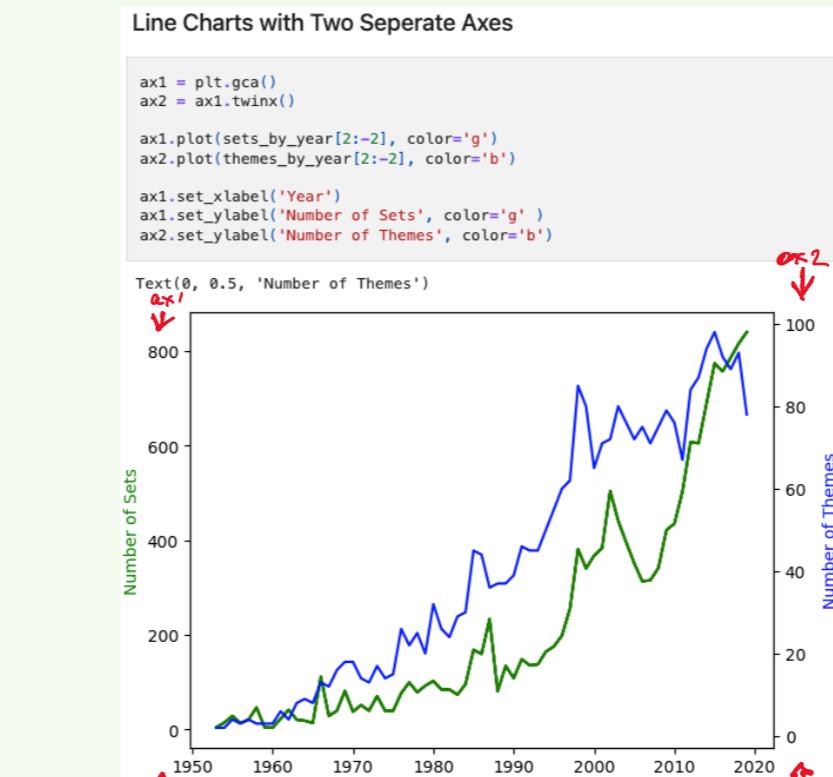
plt.figure(figsize=(10, 6))
ax1 = plt.gca()
ax2 = ax1.twinx()
ax1.plot(sets_by_year[::2], color='g')
ax2.plot(sets_by_year[1::2], color='b')
ax1.set_xlabel('Year')
ax1.set_ylabel('Number of Sets', color='g')
ax2.set_ylabel('Number of Themes', color='b')
ax1.set_xticks([1950, 1960, 1970, 1980, 1990, 2000, 2010, 2020])
ax2.set_xticks([1950, 1960, 1970, 1980, 1990, 2000, 2010, 2020])
ax1.set_yticks([0, 200, 400, 600, 800, 1000])
ax2.set_yticks([0, 0.5, 1.0, 1.5, 2.0, 2.5])
ax1.grid(True)
ax2.grid(True)
ax1.yaxis.set_label_position("left")
ax2.yaxis.set_label_position("right")
ax1.yaxis.set_ticks_position("left")
ax2.yaxis.set_ticks_position("right")
ax1.yaxis.set_tick_params(labelsize=10)
ax2.yaxis.set_tick_params(labelsize=10)
ax1.xaxis.set_tick_params(labelsize=10)
ax2.xaxis.set_tick_params(labelsize=10)
ax1.xaxis.set_label('Year')
ax1.yaxis.set_label('Number of Sets')
ax2.yaxis.set_label('Revenue in $ millions')
ax1.set_title('Free vs Paid Apps by Category')
ax2.set_title('Number of Themes')
    
```

Data UT2**Matplotlib**

Customization:

- `linestyle = '-', '--', ':', '-.'`
- `color = 'blue', 'green'`
- `marker = 'o'`

`size = plt.figure(figsize=(14, 8), dpi=120)`

**Box Plot****Line Charts with Two Separate Axes****Smoothing**