| Steps | SQL | | SEABORN | | PANDAS | |
|---|---|---|---|---|---|---|
| | **code** | **args** | **code** | **args** | **code** | **args/attributes** |
| **1 Import & export** | | | *set seaborn style*<br>　sns.set_style()<br>　sns.set_context()<br>*palette*<br>　sns.set_palette()<br>*close & clear*<br>　plt.cla()<br>　plt.clf()<br>*save*<br>　plt.savefig() | <br>whitegrid<br>notebook<br><br><br><br>clear axis<br>clear figure<br><br>foo.png | *load from file*<br>　pd.read_csv-excel-html<br>　pd.read_json()<br>　pd.read_pickle()<br><br>*load from db*<br>　pd.read_sql(query, conn)<br><br>*load from web/API* | <br><br>orient='split' |
| **2 Describe & Subset** | *subset*<br>　SELECT<br><br><br>　SELECT DISTINCT | | *correlation*<br>　sns.heatmap()<br>　sns.pairplot()<br><br>*relation plot*<br>　sns.relplot()<br><br>*usual plots -> 1 plot*<br>　sns.scatterplot()<br>　sns.histplot()<br>　sns.countplot() | relplot<br>(relational)<br><br>scatterplot<br>lineplot | *basic descriptive stats*<br>　df.describe()/info/shape<br>　df.value_counts()<br>　df.nunique()<br>　df[[c1,c2,c3]].duplicated()<br><br>*slice*<br>　df.iloc()<br>　df.loc()<br>*correlate*<br>　df1.corr(df2)<br>*plotting*<br>　df.plot<br>　.hist()<br>　.scatter() | <br><br><br><br><br><br><br><br><br><br><br>kind, ax<br><br>x, y |
| **3 Filter & Clean** | *filter*<br>　WHERE<br><br>*within a group statement*<br>　HAVING | <br>a = b<br>a IN …<br>a IS (NOT) NULL | *set facetgrid*<br>　g = sns.FacetGrid()<br><br>*map a plot to facetgrid*<br>　g = g.map(plt.hist, 'age') | | *using logical criteria*<br>　df[df.col1 >0]<br>*using query*<br>　df.query()<br>*using regex*<br>　df.filter()<br>*rename/drop columns*<br>　df.rename()<br>　df.drop()<br>*drop missing or duplicated data*<br>　df.dropna()<br>　df.drop_duplicates()<br>*change data type & operate on it*<br>　s.astype('int')<br>　s.astype('category').cat.rename_categories() | <br><br><br>query<br><br>regex<br><br>columns={}<br>columns=[]<br><br><br>[c1, c2…], keep='first' |
| **4 Impute (num & categ)** | | | *distr plots onto FacetGrid*<br>　sns.displot()<br><br>*reg plot onto FacetGrid*<br>　sns.lmplot() | kind={'hist','kde','ecdf'}<br>x=<br><br>displot<br>(distributions)<br><br>histplot<br>kdeplot<br>ecdfplot<br>rugplot | *impute missing data*<br>　df.fillna()<br><br>*stats for imputation*<br>　s.median()/mean()/mode()<br>　s.var()/std()<br><br>*stats for distribution*<br>　s.quantile() | <br>method = {'ffill', 'bfill'}<br><br><br># check for average<br># check for dispersion<br><br><br>[0,25, 0,75] |

| Steps | | SQL | | SEABORN | | PANDAS | |
|---|---|---|---|---|---|---|---|
| | | code | args | code | args | code | args/attributes |
| 4bis | **Impute (dates & text)** | *parse dates*<br>  DATE(string)<br>  TO_DATE(string, format)<br>  TO_TIMESTAMP(string, format)<br>  EXTRACT('year' FROM 'date_ts' :: timestamp)<br>  DATE(string) + INTERVAL '3 days'<br><br><br>*clean & extract text data*<br><br>  TRIM/LTRIM(string)<br>  LOWER/UPPER(string) | | | | *dates : parse (pd.Timestamp / pd.Period)*<br>  datetime = pd.to_datetime('2017-01')<br>  timestamp = pd.Timestamp('date_ts')<br>  timestamp.strftime()<br>  period = pd.Period('2017-01')<br>  period/DatetimeIndex.asfreq('D')<br><br>*text data : clean & extract*<br>  s.str.strip()/lstrip()/split()<br>  s.str.lower()/upper()<br>  s.str.title()/capitalize()<br>*text data : match pattern*<br>  s.str.contains() | <br><br>attr : {year, month_name()...}<br># -> string<br><br># define frequency<br><br><br>pat='/', expand=True<br><br><br><br># False True True |
| 5 | **Summarize & Scale** | *group*<br>  GROUP BY<br><br><br><br>*window functions*<br>  LAG(value, offset) OVER()<br>  LEAD(value, offset) OVER()<br><br>*window functions (rolling summary)*<br>  SUM() OVER(<br>    ROWS BETWEEN CURRENT AND 1 FOLLOWING) | <br>COUNT()<br>MEDIAN()<br>AVG()<br><br><br># call to last value<br># call to following value | *cat plot -> FacetGrid (1 cat + 1 discr/contin)*<br>  sns.catplot()<br>*cat plot -> FacetGrid (1 cat + 1 discr-var dist)*<br>  sns.catplot()<br>*cat plot -> FacetGrid (1 cat + 1 contin-var dist)*<br>  sns.catplot() | kind=<br>{'point','bar','count'}<br><br>kind=<br>{'box','violin','boxen'}<br><br>kind=<br>{'strip','swarm'}<br><br>catplot (categorical)<br>stripplot<br>swarmplot<br>boxplot<br>violinplot<br>pointplot<br>barplot | *group statistics*<br>  df.groupby().size()<br>  df.groupby()[col].count()<br>  df.groupby()[col].rank()<br><br>*multiple aggregate*<br>  df.groupby().agg(func)<br><br>*window functions*<br>  df.expending()<br>  df.rolling() | <br># size of each group<br># non-NA count of each obj<br>method, pct<br><br><br>np.mean<br>np.cumsum<br><br><br># -> expending object<br># -> rolling object for summary funcs<br># to be applied to windows of length n |
| 6 | **Reshape & Combine** | *reshape*<br>  ORDER BY<br><br><br>*combine by rows*<br>  UNION / UNION ALL<br>  INTERSECT / EXCEPT<br>*combine by columns*<br>  INNER/LEFT JOIN … ON … (ids)<br>  FULL OUTER JOIN … ON … (ids) | | *multiple plots*<br>  f, ax1, ax2 = plt.subplots<br><br>*add several plots on a same fig*<br>  ax1 = sns.boxplot()<br><br>*handle grid facets*<br>  g = sns.catplot()<br><br>  g.set_[xy]ticklabels()<br>  g.set_axis_labels() | <br>nrows, ncols<br>figsize=(5,6)<br><br><br><br><br><br>col = 'catcol1'<br>row= 'catcol2'<br>[men, women]<br>xlabel, ylabel | *reshape*<br>  df.sort_values()<br>  pd.melt()<br>  pd.pivot()<br><br>*combine by rows*<br>  pd.concat()<br>  s.str.cat(sep=' ')<br><br>*combine by columns*<br>  pd.merge()<br>  df1.merge(df2) | <br>ascending=<br>id_vars, value_vars<br>columns, values<br><br><br>[df1, df2], axis=0/1<br># series string concat<br><br><br>how={'left', 'outer'}<br>how={'left', 'outer'} |
| 7 | **Customize & Transform** | | | *customize labels*<br>  ax.set_[xy]label()<br><br>*customize title*<br>  ax.set_title()<br><br>*other*<br>  ax.set_[xy]lim()<br>  ax.set_[xy]ticks()<br>  ax.set_[xy]ticklabels()<br>  ax.bar_label() | | *add/transform columns*<br>  df.assign()<br>  df.insert()<br>  pd.qcut(df.col, n) | <br>new=lambda df: …<br># insert col at specific location<br># bin col into n buckets |