

Selecting elements: df[col]	<i>df.col = prob1 numbers as prefix</i>
df.col_filter= df[['col1', 'col5', 'col17',...]]	<i>filter columns</i>
df.row_filter= df[[True,False...]]. copy()	<i>filter rows</i>
#pd.options.mode.chained_assignment = None #default='warn' filterd_rows = ~ (df['manufacturer'] == 'Airbus') & (df['nb_engines'] == 2)	
df.ix['label'] = df. loc ['label'] = ret all 'label' row(s) → ret DF or S	
df.iloc[N] = ret row at position N → ret S	
df[N:N+1] =extract row N <i>df[N] will not work since it does not exist</i>	
df['col']['label'] = ret element → df.iloc[N]['col']=df['col'].iloc[num]	
df.info()	df.head(N)
df.tail(N)	df.sample()
#pd.set_option('display.max_columns', 100/None)	
df.describe ([include='all'])=ret df with statistics, excl NaN; incl strings	
df.dtypes = ret dtypes of the object	
df.astype(dtype, copy=True, raise_on_error=True) <i>cast S/df/row/col</i>	
df.shape =(rows, columns)	<i>size=df.shape[0] * df.shape[1]</i>
df.size = num of elem in df →	<i>len(df) = number of records/rows</i>
df.count ([axis=0/1])=ret S with label & number of not-null values	
df.isnull().sum (axis=0/1)=ret S with label & number of null values	
df.fillna (value=def_vals, method=None, axis=None, inplace=False) <i>method:None/ffill/ bfill propagate forward/backward last/next valid value</i> <i>def_vals = {'c1': 'v1', c2': v2,...'}; df.fillna(def_vals)</i>	
df.dropna (axis=0/1/[0,1], how='any', thresh=N, inplace=False) <i>how: drop label if 'all' or at least one ('any') element is NA</i> <i>df_new= df.dropna(axis=1, thresh=1000)</i>	
df.drop (labels, axis=0, inplace=False, errors='raise') <i>cols_drop = df.columns[df.count() < 1000]; df_new=df.drop(col_drop, axis=1)</i>	
df.drop_duplicates()	
df['column'].map(function) = apply to elements of columns or rows <i>df['by_hour']=map(lambda x:pd.to_datetime(x).hour,df['TimeStr'])</i>	
df.apply(func, axis=0/1) = apply to series over row/col	
df.applymap(func) =apply to all elements	
groups=df.groupby('gr_col')	<i>group elems=DF with all cols & same gr_col val</i>
<i>for group_name, group_elements in groups:</i>	
<i>NaN groups in GroupBy are automatically excluded. This behavior is consistent with R, for example. convert via .astype(str) to string before grouping. That will conserve the NaN's. df.astype(str).groupby('b').size()</i>	
groups.groups = ret dict with group name & row labels	
<i>When grouping, we specify col(s) to group, then col to 'select' then agg/ops on select</i>	
groups['sel_col'].agg (['mean', 'count', 'std', 'nunique'])	<i>column stats</i>
sel_col_groups = groups['sel_col']	<i>ret (Series Name, Series elems)</i>
<i>avg_of_col_per_group = sel_col_groups.mean()</i>	
groups.agg ({'col1':'fun1', 'col2':'fun2'})	<i>apply only on selected col</i>
groups.agg (['mean', 'count', 'std'])	<i>apply alls func to all col except gr_col</i>
df.sort_values (by=['col1', 'col2'], ascending=False)	
<i>delays=df.groupby('TailNum')['DepDelay'].agg(['mean','count']).sort_values(by='mean')</i>	
df.corr() =correlation between the numerical variables	
df.sum (axis=0/1) = 0→def=sum over rows of a column	
df.mean ([axis=0/1,skipna=True,level=None,numeric_only=None])	
df.std ([axis=0/1,skipna=True,level=None,ddof=1]) <i>ddof=degrees of freedom</i>	
df.join (df2/S, on='col', how='left', lsuffix=' ', rsuffix=' ', sort=False) <i>df.join([df2,df3]) no options → left join on indexes</i>	
pd.concat (objs, axis=0, join='outer', join_axes=None,copy= T/F, keys=, ignore_index=T/F, levels=, names=, verify_integrity= T/F,) <i>pd.concat([df1, df2, df3]) no options → outer join on indexes</i>	
df.merge (df2, how='inner', on=, left_on=, right_on=, left_index=T/F, right_index=T/F, sort=T/F, suffixes=('_x', '_y'), indicator=T/F)	
.str = vectorized string functions for Index and Series <i>(always use .str.func)</i>	
.str.split() = split each element to list	
.str.startswith('P').unique() <i>df[col_name.str.startswith('P').fillna(False)].count()</i>	
df.columns=df.columns.str.strip() = remove blank spaces	
filter_cols = df2.columns[df2.columns. str.contains ('Origin')]	
pd.set_option ('display.max_colwidth', -1)	
pd.set_option ('display.max_columns', 100/None)	