# How to use Pandas for SQL-like actions?

Hessam Mohammad Hosseini

January 3, 2023

**Abstract**

This document is part 1 of my cheat sheet on **Pandas** which provide overall review on how to use **Pandas** for those familiar with **SQL**.

# Contents

In this cheetsheet, I try to make readable and easy to use reference for SQL users aiming to explain how to do similar action in Pandas. My assumption is you know how to have your data as dataframe in Pandas. As soon as you have a dataframe, you can query like a table in SQL. Many possibilities are avaiable IO Tools (CSV, EXCEL,DB connection . . . )but most important ones are reading from Excel and CSV.

For example, following line of code will read file located at `file_address`, skip first row, consider `Time` column to be parsed as `datetime` format and consider `,` as thousands identifer while reading file.

```
df = pd.read_csv(file_name,

                 skiprows=10, # ignore first 10 rows of file

                 parse_dates=['Time'], # Parse Time as datetime column

                 thousands=',') # numbers have , to demonstrate 000 sepration
```

1. read_csv Function has lots of useful options to fasilitate work and avoid doing extra data cleaning tasks after data loading. Here are options for `read_csv`.
   ```
   df = pandas.read_csv(filepath_or_buffer, sep=', ', delimiter=None,
   header='infer', names=None, index_col=None, usecols=None,
   squeeze=False, prefix=None, mangle_dupe_cols=True, dtype=None,
   engine=None, converters=None, true_values=None, false_values=None,
   skipinitialspace=False, skiprows=None, skipfooter=0, nrows=None,
   na_values=None, keep_default_na=True, na_filter=True, verbose=False,
   skip_blank_lines=True, parse_dates=False, infer_datetime_format=False,
   keep_date_col=False, date_parser=None, dayfirst=False, iterator=False,
   chunksize=None, compression='infer', thousands=None, decimal=b'.',
   lineterminator=None, quotechar='"', quoting=0, doublequote=True,
   escapechar=None, comment=None, encoding=None, dialect=None,
   tupleize_cols=None, error_bad_lines=True, warn_bad_lines=True,
   delim_whitespace=False, low_memory=True, memory_map=False,
   float_precision=None)
   ```

2. read_excel Functions has lots of useful options to fasilitate work. Here are options for read_excel:
   ```
   df = pandas.read_excel(io, sheet_name=0, header=0, names=None,
   index_col=None, parse_cols=None, usecols=None, squeeze=False, dtype=None,
   engine=None, converters=None, true_values=None, false_values=None,
   skiprows=None, nrows=None, na_values=None, keep_default_na=True,
   verbose=False, parse_dates=False, date_parser=None, thousands=None,
   comment=None, skip_footer=0, skipfooter=0, convert_float=True,
   mangle_dupe_cols=True, **kwds)
   ```

I tried to summerized and add to what was available on Pandas comparison with SQL aiming to simplify understanding. For details on functionality, please check and review Pandas documentation.

# 1  SELECT

| SQL Sample | Pandas Sample |
|---|---|
| select * <br> FROM table | table |
| select distinct c5 <br> FROM table | table['c5'].unique() or <br> table[['c5']].drop_duplicates() |
| select c1, c2,c10 <br> FROM table | df[[c1, c2,c10]] |
| select c10, c2, c1 <br> FROM table | df[[c10, c1, c2]] |
| select c10, c2*12 - c1 <br> + c6 <br> FROM table | df['new c'] = df.c2*12 - <br> df.c1 + df.c6df[[c10,'new <br> c']] or df.assign(c10 = <br> df.c10, new_c = df.c2*12 - <br> df.c1 + df.c6) |
| SELECT total_bill, tip, <br> smoker, time <br> FROM tips <br> LIMIT 5 | tips[['total_bill', 'tip', <br> 'smoker','time']].head(5) |

## 1.1  Update or delete

| | SQL Sample | Pandas Sample |
|---|---|---|
| Delete | DELETE <br> FROM tips <br> WHERE tip > 9 | tips = <br> tips.loc[tips['tip'] <= <br> 9] |
| Update | UPDATE tips <br> SET tip = tip*2 <br> WHERE tip < 2 | tips.loc[tips['tip'] < <br> 2, 'tip'] *= 2 |

# 2  Conditioning at WHERE

| | SQL Sample | Pandas Sample |
|---|---|---|
| - | SQL SELECT * <br> FROM tips <br> WHERE time = 'Dinner' <br> LIMIT 5 | tips[tips['time'] == <br> 'Dinner'].head(5) or <br> is_dinner = <br> tips['time'] == <br> 'Dinner'tips[is_dinner].head(5) |
| AND | SELECT * <br> FROM tips <br> WHERE time = 'Dinner' AND <br> tip > 5.00 | tips[(tips['time'] == <br> 'Dinner') & <br> (tips['tip'] > 5.00)] |

|  | SQL Sample | Pandas Sample |
|---|---|---|
| OR | SELECT *<br>FROM tips<br>WHERE size >= 5 OR<br>total_bill > 45 | tips[(tips['size'] >=<br>5') or<br>(tips['total_bill'] ><br>45)] |
| IS NULL | SELECT *<br>FROM t<br>WHERE col2 IS NULL | t[t['col2'].isna()] |
| IS NOT NULL | SELECT *<br>FROM t<br>WHERE col IS NOT NULL | t[t['col2'].notna()] |

| SQL Sample | Pandas Sample |
|---|---|
| SELECT *<br>FROM tips<br>ORDER BY tip DESC<br>LIMIT 10 OFFSET 5 | tips.nlargest(10 + 5,<br>columns='tip').tail(10) |
| SELECT total_bill, tip,<br>smoker, time<br>FROM tips<br>ORDER BY tip<br>DESC<br>LIMIT 10 OFFSET 5 | tips[['total_bill', 'tip',<br>'smoker','time']]<br>tips.nlargest(10 + 5,<br>columns='tip').tail(10) |

# 3   GROUP BY

| SQL Sample | Pandas Sample |
|---|---|
| SELECT sex, count(*)<br>FROM tips<br>GROUP BY sex | tips.groupby('sex').size()<br>or<br>tips.groupby('sex')['total_bill'].count() |
| select A, sum(C),<br>sum(D)<br>FROM df<br>GROUP BY A | df.groupby('A')['B','C'].sum() |

| SQL Sample | Pandas Sample |
|---|---|
| SELECT day, AVG(tip),<br>COUNT(*)<br>FROM tips<br>GROUP BY day | tips.groupby('day').agg({'tip':<br>np.mean, 'day': np.size}) |

| SQL Sample | Pandas Sample |
| --- | --- |
| SELECT smoker, day, COUNT(*), AVG(tip) FROM tips GROUP BY smoker, day | tips.groupby(['smoker,'day']).agg({'tip': [np.size, np.mean]}) |

| SQL Sample | Pandas Sample |
| --- | --- |
| SELECT c1, COUNT(*) FROM df where country='IR' GROUP BY c1 having count(*)>1000 | df[df.country == 'IR'].groupby('c1').filter(lambda g: len(g) > 1000).groupby('c1').size() |
| SELECT c1, COUNT(*) FROM df WHERE country='IR' GROUP BY c1 HAVING count(*)>1000 ORDER BY count(*) desc | df[df.country == 'IR'] .groupby('c1'). filter(lambda g: len(g) > 1000) .groupby('c1').size() .sort_values(ascending=False) |

# 4   ORDER BY

| SQL Sample | Pandas Sample |
| --- | --- |
| SELECT * FROM df ORDER BY A, B | df.sort_values(['A', 'B']) |
| SELECT * FROM df ORDER BY A desc, C | df.sort_values(['A', 'B'],ascending=[False, True]) |

# 5   UNION, JOIN and other set related operations

I will work to provide more comprehensive explanations on this part.

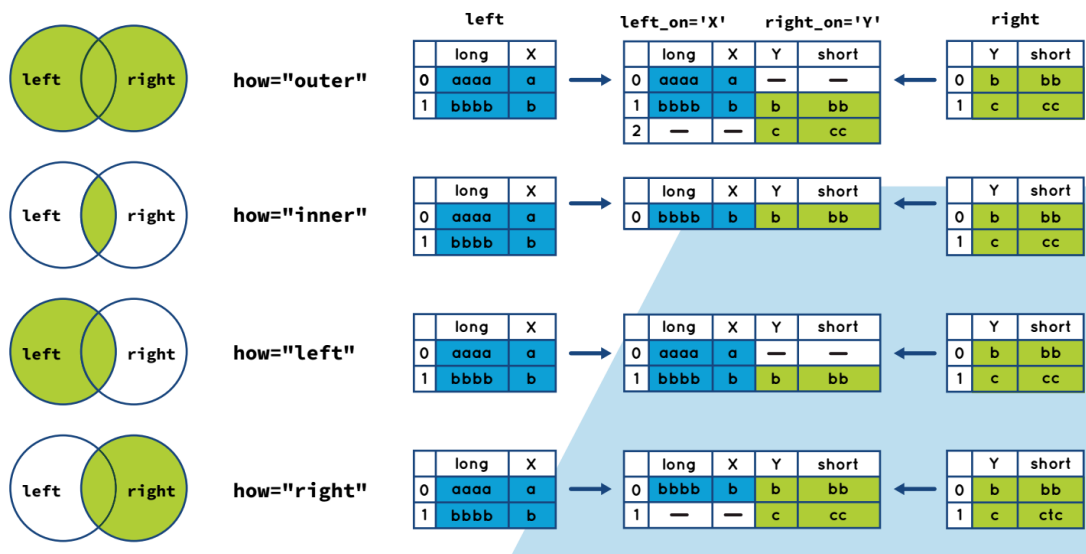## 5.1   Union

| SQL Sample | Pandas Sample |
| --- | --- |
| SELECT c1, c2 FROM df1 UNION ALL SELECT c1, c2 FROM df2 | pd.concat([df1, df2]) |

Difference between **union all** and **union** is that **union** will remove duplicates.

| SQL Sample | Pandas Sample |
|---|---|
| SELECT c1, c2<br>FROM df1<br>UNION<br>SELECT c1, c2<br>FROM df2 | pd.concat([df1, df2])<br>.drop_duplicates() |

## 5.2  Different Join cases

I will add parts to make eplanation on `join` more clear and compehensive. Below image extracted from Enthoughtnamed "Enthought-Python-Pandas-Cheat-Sheets-1-8-v1.0.2" worth more than 100 sentences to explain different types of `join`. You can get whole file via registration on Enthought.



| SQL Sample | Pandas Sample |
|---|---|
| SELECT *<br>FROM df1<br>INNER JOIN df2<br>ON df1.key = df2.key | pd.merge(df1, df2,<br>on='key') |
| SELECT *<br>FROM df1<br>INNER JOIN df2<br>ON df1.c7 = df2.c5 | pd.merge(df1, df2,<br>left_on='c7',right_on='c5') |

# 6  Time functionality

In order to have possibility to use time related functionalities, we need help `Pandas` understand which columns are to be treated as time. Of course, the columns should be in for that converting them to time format is possible. For details, please check and review Pandas documentation. If you manage to let `Pandas` know properly which column(s) to be time related column(s), they will end up having `datetime64[ns]` format. `.dtypes` on dataframe provides you with columns formats. Pass date related column(s) you need to `parse_dates` to `read_csv` or `read_excel` functions. Check Pandas documentation for more details.

Doing so, you can apply `.dt` on column to have date - time selection like

`dt.dayofweek`

`dt.minute`

`dt.hour`

`dt.second`

`dt.quarter`

`dt.month`

`dt.month_name`

`dt.weekday_name`

`dt.weekday`

`dt.weekofyear`

`dt.year`

## 6.1 How to get current date time using pandas?

`pd.datetime.now()`

`pd.datetime.now().date()`

`pd.datetime.now().year`

`pd.datetime.now().month`

`pd.datetime.now().day`

`pd.datetime.now().hour`

`pd.datetime.now().minute`

`pd.datetime.now().second`

`pd.datetime.now().microsecond`

Again, check Pandas documentation for more! Here, we assume `sdate` column to have `datetime64[ns]` format.

| | SQL Sample | Pandas Sample |
|---|---|---|
| sysdate - n | SELECT * <br> FROM df <br> WHERE sdate> sysdate-5 | df[df['sdate'].dt.date()> <br> pd.datetime.now().date()-5] |
| month | SELECT * <br> FROM df <br> WHERE sdate in Q1 | df[(df.sdate.dt.month <br> >= 1) & <br> (df.sdate.dt.month <= <br> 3)] |

| | SQL Sample | Pandas Sample |
|---|---|---|
| between | SELECT * <br> FROM t <br> WHERE to_char(sdate,'yyyy') <br> between 1998 AND 2018 | t[(t.sdate.dt.year >= <br> 1998) & <br> (t.sdate.dt.year <= <br> 2018)] |
| | SELECT * <br> FROM t <br> WHERE to_char(sdate <br> ,'day')= 'Friday' | df[df.sdate.dt.day_name() <br> == 'Friday'] |

# 7   String related functionality like `like`, `Substr`

For columns with **string** content, we could access string related funftionality by applying `.str` on column. Here are few samples: `str.contains` - contains options:

`Series.str.contains(pat, case=True, flags=0, na=nan, regex=True)`

`Here are list of main 'string' functions.`

`str.upper`

`str.lower`

`str.extract`

`str.extractall`

`str.find`

`str.findall`

`str.len`

`str.replace`

`str.slice`

`str.split`

`str.strip`

Check Pandas documentation for more!

| | SQL Sample | Pandas Sample |
|---|---|---|
| regex | SELECT <br> upper(trim(to_char(LAC,'xxxxx')) <br> \|\|'-' \|\| <br> trim(to_char(CI,'xxxxx'))) <br> AS "LAC-CI(HEX)" <br> FROM t | t = t['LAC','CI']\ <br> .apply(lambda x: x\ <br> .astype(str)\ <br> .map(lambda x: int(x, <br> base=16)))t <br> .assig(LAC-CI(HEX) = <br> t['LAC']+'-'+t['CI'] |
| substr | SELECT * <br> FROM tips <br> WHERE substr(time,1,2)  = <br> 'Di' | tips[tips['time']\ <br> .str[:2] == 'Di'] |

| | SQL Sample | Pandas Sample |
|---|---|---|
| regex | `SELECT`<br>`upper(trim(to_char(LAC,'xxxxx'))`<br>`\|\|'-' \|\|`<br>`trim(to_char(CI,'xxxxx')))`<br>`AS "LAC-CI(HEX)"`<br>`FROM t` | `t = t['LAC','CI']\`<br>`.apply(lambda x: x\`<br>`.astype(str)\`<br>`.map(lambda x: int(x,`<br>`base=16)))t`<br>`.assig(LAC-CI(HEX) =`<br>`t['LAC']+'-'+t['CI']` |
| like | `SELECT *`<br>`FROM df`<br>`WHERE Country  like '%IR%'` | `df[df['Country']\`<br>`.str.contains('IR') ==`<br>`True]` |
| like | `SELECT *`<br>`FROM df`<br>`WHERE Country  like 'IR%'` | `df[df['Country']\`<br>`.str.startswith('IR')`<br>`== True]` |
| like | `SELECT *`<br>`FROM df`<br>`WHERE Country  like '%AN'` | `df[df['Country']\`<br>`.str.endswith('AN') ==`<br>`True]` |
| in | `SELECT *`<br>`FROM df`<br>`WHERE City in ('TEHRAN',`<br>`'BERLIN','STOKHOLM')` | `df[df['City']\`<br>`.isin(['TEHRAN',`<br>`'BERLIN','STOKHOLM'])` |
| regex | `SELECT last_name`<br>`FROM contacts`<br>`WHERE REGEXP_LIKE`<br>`(last_name, '^+A(*)')` | `contacts[contacts['last_name']\`<br>`.str.contains('^+A(*)')]` |
| regex | `SELECT c1`<br>`FROM t`<br>`WHERE`<br>`REGEXP_LIKE(c1,'([A-Z][\d]{4})')` | `t[t['c1']\`<br>`.str.contains(([A-Z][\d]{4}))]` |

`\|\|` provide concatenation functonality in PL/SQL. In **Python**, $+$ on `string` values resulted in contanated text.

| SQL Sample | Pandas Sample |
|---|---|
| `SELECT *`<br>`FROM`<br>`(SELECT t.*,`<br>`ROW_NUMBER()`<br>`OVER(PARTITION BY day`<br>`ORDER BY total_bill`<br>`DESC) AS r`<br>`FROM tips t)`<br>`WHERE r <`<br>`ORDER BY day, r;` | `(tips.assign(r=tips\`<br>`.sort_values(['total_bill'],`<br>`ascending=False)\`<br>`.groupby(['day'])\`<br>`.cumcount()+1)\`<br>`.query('r < 3')`<br>`.sort_values(['day', 'r']))` |

**Oracle's ROW_NUMBER() analytic function**

# 8    To check for missing values

- `df.notnull()` Use to Drop Missing Values

- `df.dropna()` Filling Missing Values — Direct Replace

- `df.fillna()`

Besides Pandas comparison with SQL, I also get ideas from following references:

1. pandas-cheatsheet-for-sql-people-part-1

2. how-to-rewrite-your-sql-queries-in-pandas-and-more

3. thinking-like-sql-in-pandas

4. did-you-know-pandas-can-do-so-much

5. 10-python-pandas-tricks-that-make-your-work-more-efficient