

# Cheatography

## Data Science Cheat Sheet by elhamsh via cheatography.com/31327/cs/13764/

### Pandas

```
import pandas as pd
df.iloc[:5,:]
return slice of data:all columns
first 5 rows
type(df)
DataFrame
df.shape
(len, #ofcols)
df.columns
name of cols
df.index
return index column
df.head(3)
return first 3 rows
df.iloc[-5:,:]
return last 5 rows
df.tail()
return last 5 rows
df.info()
return index, column types, #
of row, # of not null cols
type(df['low'])
Series
type(df['low'].v
alue)
numpy.ndarray
np.log10(df['lo
w'])
return data frame
np.log10(df['lo
w'].values)
return list of list
```

Each column in pandas is a Series.  
You can run numpy on df or a col of df

### Statistical Data Analysis

```
df.describe()
count, mean,std,max, quartiles for
() each col of non-null rows
df['low'].co
return # of not null rows
unt()
df[cols].co
return a series
unt()
df['low'].mean()
return mean ignoring nulls
ean()
df.std()
df.median()
df.quantile(
q=.5:median q=[.25,.75]:IQRrange
q)
df['low'].mi
alphabetic order for non-numerics
n()
df['low'].m
alphabetic order for non-numerics
ax()
```

### Statistical Data Analysis (cont)

```
df.mean(axis='colu
mns')
mean of all columns for
each row
df.low
df['low']
```

### Time series

```
index_col='Date', parse_date=True
df.loc['2015-
2']
df.loc['2015-
2-20']
df.loc['2015-
2-20':
'2015-3']
newD =
y-m-d h:m:s
pd.to_datedate
me('Date' )
df.reindex(n
ewD)
reindexing with matching dates.
if doesn't match,new rows w.
null value
df.reindex(n
ewD,metho
d='ffill')
fill empty values forward
fill:value of previous rows
method='bfil
l'
backward fill: value of later rows
l'
```

```
df.resample
('D').mean()
'H', 'min',
'2W'
'Y', 'Q', 'M',
'B'
df.resample
('W').sum().
max()
df.resample
('4h').ffill()
df1+df2
df['Tempera
ture']
['2010-augu
st']
```

```
df['Tempera
ture']
['2010-2']
```

### Time series (cont)

```
unsmooth.rolling(
window=24).mea
n()
df['type'].str.upper
()
df['product'].str.co
ntains('ware')
True+True
2
False + False
0
df['product'].str.co
ntains('ware').su
m()
df['date'].dt.hour
return hour of each row 0-
23
df['date'].dt.tz_loc
alize('US/Central'
)
df['date'].dt.tz_convert('US/Eastern')
df['date'].resampl
e('A').first()
df['date'].resampl
e('A').first().interpol
ate('linear')
df.columns.str stri
p()
df.set_index('Date', inplace=True)
newD = pd.to_datetime('Date_list',
format='%Y-%M-%D %H:%M')
pd.Series(Colum
ns_list,
index=newD)
ts2_interp =
ts2.reindex(ts1.in
dex).interpolate(h
ow='linear')
timezone.dt.tz_loc
alize('US/Central'-'
)
```

# Cheatography

## Data Science Cheat Sheet by elhamsh via cheatography.com/31327/cs/13764/

### Build DF

```
df=pd.read_csv("filepath", add index column 0-  
index_col=0) len(inp)  
  
index_col='nameofacolumn'  
  
df.index=['A', 'B', ...] assign index to df.  
len(index)==len(df)  
  
pd.DataFrame({'id': [1,2,3], 'gen':'M'}) key: columns,  
values: row  
  
pd.DataFrame(dict_of_lists)  
  
zipped=list(zip(list_labels, list_values))  
  
pd.DataFrame(dict(zippe d)) list_labels,  
list_values = list of  
list  
  
pd.read_csv("filepath", header=None) no header  
  
pd.read_csv("filepath", col_n:list of column  
options) names  
  
header=0, rename the header  
names=col_n  
  
header=None, no header in file &  
names=col_n header is col_n  
  
na_values='-1' convert specific  
value (-1) to a nan  
  
na_values={'colname':[-1, "]} define a dic for  
each col  
  
parse_dates=[[0,1,1]] convert 3 columns  
of date to one col  
  
parse_dates=True convert column with  
date to dateformat  
  
delimiter=' ' delimiter  
  
header=3 header is in index 3  
  
comment='#' ignore all lines start  
with '#' in the input  
  
index_col = 'dates' set a column as  
index  
  
df[cols] take specific  
columns  
  
df.to_csv('outputpath')  
  
df.to_excel('outputpath')  
  
pd.DataFrame({'smoothed': smoothed, 'unsmoothed': unsmoothed}) create df.if they  
have index, will  
merge based on  
index
```

### categorical

```
df['type'].decri count not null,# of unique,top  
be() item,freq. of top  
  
df['type'].uniq #of unique items  
ue()  
  
df.loc[df['type' ]==x,:] df[df['type']==x]  
  
del delete a column  
def['type']
```

### Numpy+Df

```
df.values Create array of DataFrame  
values  
  
df[colname] create a columns with zero  
=0 elements in df
```

### Cleanning

```
df_dropped Remove the appropriate columns  
= list_to_drop  
  
df.drop(list_ to_drop,  
axis='colu mns')  
  
df.set_index Set colname as the index  
(colname)
```

```
pd.to_num ric() It converts a Series of values to  
floating-point values.  
Furthermore, by specifying the  
keyword argument  
errors='coerce', you can force  
strings like 'M' to be interpreted  
as NaN.
```

```
df.reset_in dex() Extract the colname column from  
df using .reset_index()  
[colname]
```

```
df.loc[df[col name]=='st h'] choose the rows in df for  
df[colname]='sth'
```

```
df.loc[df[col name].str.c ontain('sth')] choose the rows in df where the  
column df[colname] contain 'sth'
```

Not published yet.

Last updated 20th December, 2017.

Page 2 of 3.

### Plot

```
import matplotlib.pyplot as plt  
  
plt.plot(df['low'].values) x axis= index of value  
  
plt.show() show the image  
  
plt.plot(df['low']) x axis is index of df  
(eg date)  
  
df['low'].plot() plot series directly.  
has also x label  
  
df.plot() show all columns in df  
with legend  
  
plt.yscale('log') log scale on vertical  
axis  
  
df['low'].plot(color='b',style='.-', legend=True)  
  
plt.axis((minx, maxx,miny,maxy)) zoom  
  
plt.title('title')  
  
plt.ylabel('label')  
  
plt.xlabel('xlabel')  
  
plt.savefig('a.pdf')  
  
plt.savefig('a.jpg')  
  
df.plot(subplots=True) Draw each column in  
one subplot.  
  
df.plot(x='colname',y=' colname',kind='scatter  
') plot 2 columns  
  
kind = 'box' box plot  
  
kind = 'hist' histogram  
  
kind='area'  
  
bins=30 integer:#of bins  
  
range=(4,8) tuple (min,max)  
  
normed=True boolean. normalize to  
one for hist  
  
cumulative=True boolean for hist  
  
alpha=0.3 visibility of several  
histograms  
  
s=sizes sizes= array of size of  
each circle in scatter  
plot  
  
fig, axes=subplots(nrows=1,ncols=1)  
  
df['low'].plot(ax=axes[0 ..., kind, bins,  
, ...]) normed,cumulative
```

Sponsored by [CrosswordCheats.com](#)

Learn to solve cryptic crosswords!

<http://crosswordcheats.com>



By [elhamsh](#)

[cheatography.com/elhamsh/](http://cheatography.com/elhamsh/)

## Plot (cont)

```
df.plot(y='colname',kind='box')
```

```
style=' color,marker,line type  
k.-'
```

plt.clf() clears the entire current figure with all its axes, but leaves the window opened, such that it may be reused for other plots

## Indexing

```
df['colname'][rowname]      rowname is index_col
```

```
df.colname['rowname']
```

```
df.loc['rowname','colname']
```

```
df.loc['rownstart','rownend',:]    row names are inclusive.
```

```
df[['low']]                  returns a single column data frame
```

```
df['low']                   returns a series with index of df
```



By **elhamsh**

[cheatography.com/elhamsh/](http://cheatography.com/elhamsh/)

Not published yet.

Last updated 20th December, 2017.

Page 3 of 3.

Sponsored by **CrosswordCheats.com**

Learn to solve cryptic crosswords!

<http://crosswordcheats.com>