

Python for Data Analysis: Pandas

Andreas Weller, PhD

WTCHG - NHS

23.4.2014

DataFrame

- New object introduced by Pandas
- Inspired by R data.frame

DF creation

- `pd.read_csv(filename, sep=",")`
- `pd.DataFrame(data)`

Basic methods I

DF summary

Method	Effect
<code>df.describe()</code>	summary stats
<code>df.head()</code>	head
<code>df.tail(10)</code>	select rows+cols by index

Selection

<code>df[["A","B"]]</code>	select cols
<code>df.loc[:5, "A"]</code>	select rows+cols by name
<code>df.iloc[:5, :5]</code>	select rows+cols by index

Basic methods II

Boolean indexing

Method	Effect
<code>df.value > 12</code>	Query all rows
<code>df[df.value > 12]</code>	Select True rows
<code>df[cond1 & cond2]</code>	Select rows
<code>df.shape</code>	Count rows and cols

Column creation

<code>df["new"] = 1</code>	New col, same values
<code>df["new"] = df.old * 2</code>	New col from old col
<code>df["new"] = df.old.str.upper()</code>	New col from string method

Groupby

```
df.groupby(df.value2).func()
```

```
df.value1.groupby([df.value2, df.value3]).func()
```

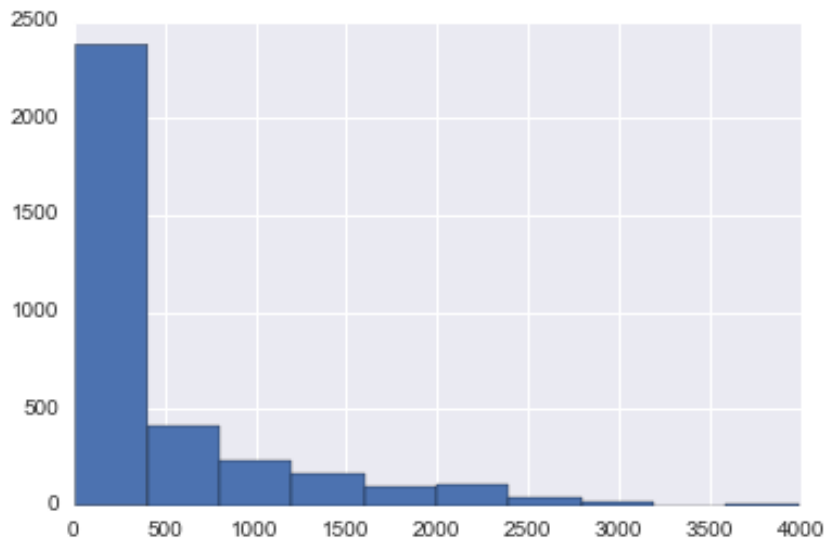
Groupby functions

Method	Effect
.sum(), .mean(), median()	Sum, Mean, Median
.count()	No. of entries
.nunique()	No. of unique entries

Histograms

```
# plot a histogram from a pandas DataFrame  
  
data = df.DP  
  
plt.hist(data, bins=100)
```

Histograms



Histograms

plot a histogram from a pandas DataFrame

```
data = df.DP
```

```
plt.hist(data, bins=100)
```

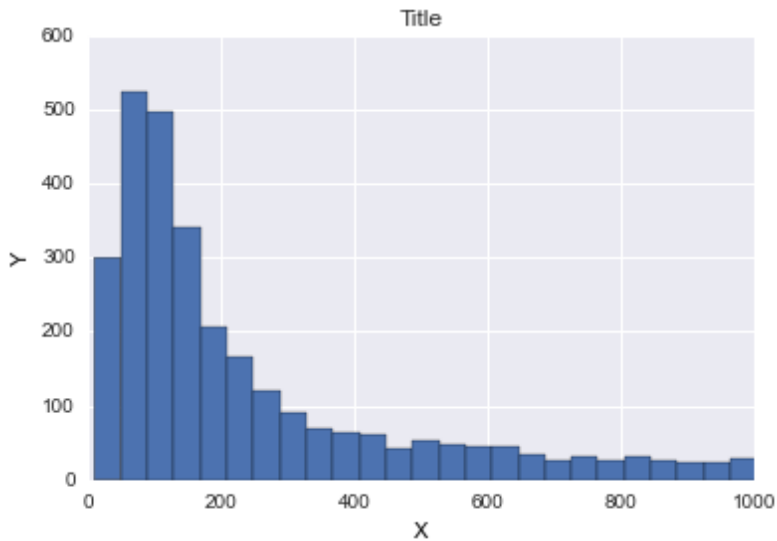
```
plt.xlabel("X")
```

```
plt.ylabel("Y")
```

```
plt.xlim(0,1000)
```

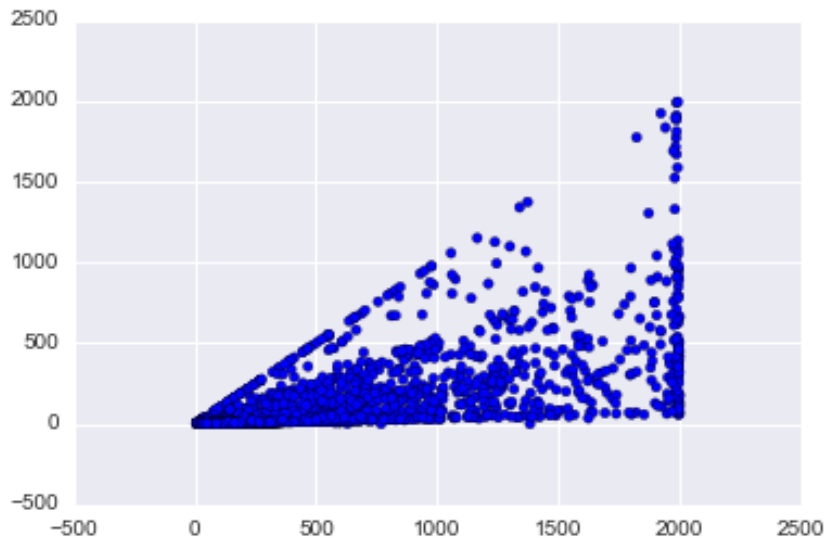
```
plt.title("Title")
```


Histograms



Scatterplots

```
plt.scatter(df.FDP, df.FAO)
```



Scatterplots

```
sns.jointplot("FDP", "FAO", df)
```

