

```
In [115]: 1 # how to stack data using pandas concatenate
          2 df1 = pd.DataFrame(np.random.randint(0,100,size=(100, 4)),
          3                        columns=list('ABCD'))
```

```
In [118]: 1 df2 = pd.DataFrame(np.random.randint(0,100,size=(150, 4)),
          2                        columns=list('ABCD'))
          3 df3 = pd.DataFrame(np.random.randint(0,100,size=(110, 4)),
          4                        columns=list('ABCD'))
```

```
In [120]: 1 frames = [df1, df2, df3]
          2
          3 df123 = pd.concat(frames)
          4 print(df123.shape)
```

```
(360, 4)
```

```
In [62]: 1
```

```
In [121]: 1 #titanic data
          2 #http://web.stanford.edu/class/archive/cs/cs109/cs109.1166/problem12.html
          3 # put the data in the directory of workshop, which is in the working
          4 directory
```

```
In [56]: 1 import numpy as np
          2 import pandas as pd
          3 import os
```

```
In [122]: 1 titanic_train = pd.read_csv("workshop/titanic.csv")
          2 char_cabin = titanic_train["Pclass"].astype(str)
          3 # Convert cabin to str -- categorical data
```

```
In [ ]: 1
```

```
In [125]: 1 my_tab = pd.crosstab(titanic_train["Survived"], # Make a crosstab
          2                        columns="count")           # Name the count
          3 column
          4 my_tab
```

```
Out[125]:
```

col_0	count
<b>Survived</b>	
0	545
1	342

```
In [126]: 1 my_tab = pd.crosstab(titanic_train["Survived"], # Make a crosstab
2                               columns="count", margins = True) #
3       Name the count column
4       my_tab
```

```
Out[126]:
```

col_0	count	All
<b>Survived</b>		
0	545	545
1	342	342
<b>All</b>	887	887

```
In [70]: 1 my_tab = pd.crosstab(titanic_train["Sex"], # Make a crosstab
2                               columns="count") # Name the count
3       column
4       my_tab
```

```
Out[70]:
```

col_0	count
<b>Sex</b>	
<b>female</b>	314
<b>male</b>	573

```
In [ ]: 1
```

```
In [127]: 1 # Table of survival vs. sex
2 survived_sex = pd.crosstab(titanic_train["Survived"],
3                             columns=titanic_train["Sex"])
4
5 survived_sex.index= ["died","survived"]
6
7 survived_sex
```

```
Out[127]:
```

Sex	female	male
<b>died</b>	81	464
<b>survived</b>	233	109

```
In [ ]: 1
```

```
In [128]: 1 # Table of survival vs passenger class
2 survived_class = pd.crosstab(titanic_train["Survived"],
3                               columns=titanic_train["Pclass"])
4
5 survived_class.columns = ["class1","class2","class3"]
6 survived_class.index= ["died","survived"]
7
8 survived_class
```

```
Out[128]:
```

	class1	class2	class3
died	80	97	368
survived	136	87	119

```
In [129]: 1 # Table of survival vs passenger class
2 survived_class = pd.crosstab(titanic_train["Survived"],
3                               columns=titanic_train["Pclass"], margins =
4                               True)
5 survived_class.columns = ["class1","class2","class3", "rows_total"]
6 survived_class.index= ["died","survived", "cols_total"]
7
8 survived_class
```

```
Out[129]:
```

	class1	class2	class3	rows_total
died	80	97	368	545
survived	136	87	119	342
cols_total	216	184	487	887

```
In [131]: 1 # simple analysis for survival by class
2 print( survived_class.iloc[0, :] / survived_class.iloc[2, :] )
```

```
class1      0.370370
class2      0.527174
class3      0.755647
rows_total  0.614431
dtype: float64
```

```
In [85]: 1 # observed values --
2 observed = survived_class.iloc[0:2,0:3] # Get table without totals
3 observed
```

```
Out[85]:
```

	class1	class2	class3
died	80	97	368
survived	136	87	119

In [ ]:

1

In [107]:

```
1 #expected
2 expected2 = np.outer(survived_class["rows_total"][0:2],
3                       survived_class.iloc[2, 0:3] )/ 487
4 expected2 = pd.DataFrame(expected2)
5
6 expected2.columns = ["class1","class2","class3"]
7 expected2.index = ["died","survived"]
8
9 expected2
10
```

Out[107]:

	class1	class2	class3
died	241.724846	205.913758	545.0
survived	151.687885	129.215606	342.0

In [108]:

```
1 chi_squared_stat = (((observed-expected2)**2)/expected2).sum().sum()
2
3 print(chi_squared_stat)
```

384.11430821033935

In [111]:

```
1 crit = stats.chi2.ppf(q = 0.95, # Find the critical value for 95%
2                          confidence*
3                          df = 5)    # *
4 print("Critical value")
5 print(crit)
6
7 p_value = 1 - stats.chi2.cdf(x=chi_squared_stat, # Find the p-value
8                              df=8)
9 print("P value")
10 print(p_value)
```

Critical value  
11.0704976935  
P value  
0.0

In [ ]:

1

```

In [132]: 1 #independence test
          2 import numpy as np
          3 import pandas as pd
          4 import scipy.stats as stats
          5 np.random.seed(10)
          6
          7 # Sample data randomly at fixed probabilities
          8 voter_race = np.random.choice(a=
          9     ["asian","black","hispanic","other","white"],
         10     p = [0.05, 0.15 ,0.25, 0.05, 0.5],
         11     size=1000)
         12
         13 # Sample data randomly at fixed probabilities
         14 voter_party = np.random.choice(a=
         15     ["democrat","independent","republican"],
         16     p = [0.4, 0.2, 0.4],
         17     size=1000)
         18
         19 voters = pd.DataFrame({"race":voter_race,
         20                     "party":voter_party})
         21
         22 voter_tab = pd.crosstab(voters.race, voters.party, margins = True)
         23
         24 voter_tab.columns =
         25     ["democrat","independent","republican","row_totals"]
         26
         27 voter_tab.index =
         28     ["asian","black","hispanic","other","white","col_totals"]
         29
         30 observed = voter_tab.iloc[0:5,0:3] # Get table without totals for
         31     later use
         32 voter_tab

```

Out[132]:

	democrat	independent	republican	row_totals
asian	21	7	32	60
black	65	25	64	154
hispanic	107	50	94	251
other	15	8	15	38
white	189	96	212	497
col_totals	397	186	417	1000

```

In [102]: 1 voter_tab.loc["col_totals"][0:3]

```

```

Out[102]: democrat      397
independent    186
republican     417
Name: col_totals, dtype: int64

```

```
In [89]: 1 expected = np.outer(voter_tab["row_totals"][0:5],
2                       voter_tab.loc["col_totals"][0:3]) / 1000
3
4 expected = pd.DataFrame(expected)
5
6 expected.columns = ["democrat", "independent", "republican"]
7 expected.index = ["asian", "black", "hispanic", "other", "white"]
8
9 expected
```

Out[89]:

	democrat	independent	republican
asian	23.820	11.160	25.020
black	61.138	28.644	64.218
hispanic	99.647	46.686	104.667
other	15.086	7.068	15.846
white	197.309	92.442	207.249

```
In [52]: 1 chi_squared_stat = (((observed-expected)**2)/expected).sum().sum()
2
3 print(chi_squared_stat)
```

7.169321280162059

```
In [53]: 1 crit = stats.chi2.ppf(q = 0.95, # Find the critical value for 95%
2                               confidence*
3                               df = 8)  # *
4 print("Critical value")
5 print(crit)
6
7 p_value = 1 - stats.chi2.cdf(x=chi_squared_stat, # Find the p-value
8                               df=8)
9 print("P value")
10 print(p_value)
```

Critical value  
15.5073130559  
P value  
0.518479392949

```
In [ ]: 1 #Given the high p-value, the test result does not detect a significant
relationship between the variables.
```

```
In [ ]: 1
```

```
In [ ]: 1
```

###iris dataset

<https://gist.githubusercontent.com/curran/a08a1080b88344b0c8a7/raw/d546eae765268bf2f487608c537c05e22e4b221/iris.csv>  
(<https://gist.githubusercontent.com/curran/a08a1080b88344b0c8a7/raw/d546eae765268bf2f487608c537c05e22e4b221/iris.csv>)

```
In [133]: 1 import pandas as pd
          2 filepath =
            'https://gist.githubusercontent.com/curran/a08a1080b88344b0c8a7/raw/d5
            46eae765268bf2f487608c537c05e22e4b221/iris.csv'
          3 iris = pd.read_csv(filepath)
```

```
In [15]: 1 print(iris.head())
```

	sepal_length	sepal_width	petal_length	petal_width	species
0	5.1	3.5	1.4	0.2	setosa
1	4.9	3.0	1.4	0.2	setosa
2	4.7	3.2	1.3	0.2	setosa
3	4.6	3.1	1.5	0.2	setosa
4	5.0	3.6	1.4	0.2	setosa

```
In [16]: 1 iris.columns
```

```
Out[16]: Index(['sepal_length', 'sepal_width', 'petal_length', 'petal_width',
               'species'],
              dtype='object')
```

```
In [17]: 1 iris_groups = iris.groupby("species")
```

```
In [18]: 1 iris.describe()
```

```
Out[18]:
```

	sepal_length	sepal_width	petal_length	petal_width
<b>count</b>	150.000000	150.000000	150.000000	150.000000
<b>mean</b>	5.843333	3.054000	3.758667	1.198667
<b>std</b>	0.828066	0.433594	1.764420	0.763161
<b>min</b>	4.300000	2.000000	1.000000	0.100000
<b>25%</b>	5.100000	2.800000	1.600000	0.300000
<b>50%</b>	5.800000	3.000000	4.350000	1.300000
<b>75%</b>	6.400000	3.300000	5.100000	1.800000
<b>max</b>	7.900000	4.400000	6.900000	2.500000

```
In [19]: 1 iris_groups.describe()
```

Out[19]:

	petal_length								petal_width			...	sepal_length		se
	count	mean	std	min	25%	50%	75%	max	count	mean	...	75%	max	co	
species															
setosa	50.0	1.464	0.173511	1.0	1.4	1.50	1.575	1.9	50.0	0.244	...	5.2	5.8		
versicolor	50.0	4.260	0.469911	3.0	4.0	4.35	4.600	5.1	50.0	1.326	...	6.3	7.0		
virginica	50.0	5.552	0.551895	4.5	5.1	5.55	5.875	6.9	50.0	2.026	...	6.9	7.9		

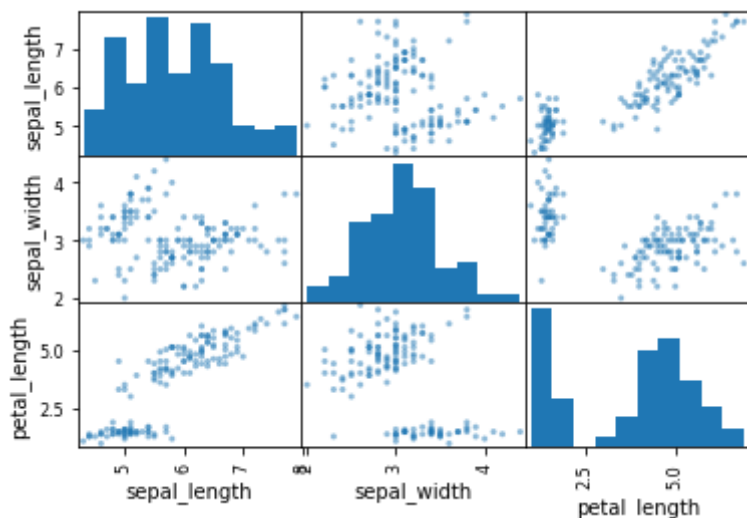
3 rows × 32 columns

```
In [20]: 1 import os
2 os.getcwd()
```

Out[20]: '/Users/byunglee'

```
In [135]: 1 %matplotlib inline
2 from pandas.plotting import scatter_matrix
3 scatter_matrix(irisData[['sepal_length', 'sepal_width',
    'petal_length']])
```

Out[135]: array([[<matplotlib.axes.\_subplots.AxesSubplot object at 0x1a14517e10>,  
<matplotlib.axes.\_subplots.AxesSubplot object at 0x1a146232e8>,  
<matplotlib.axes.\_subplots.AxesSubplot object at 0x1a145a23c8>],  
[<matplotlib.axes.\_subplots.AxesSubplot object at 0x1a146bb4a8>,  
<matplotlib.axes.\_subplots.AxesSubplot object at 0x1a146f54a8>,  
<matplotlib.axes.\_subplots.AxesSubplot object at 0x1a146f54e0>],  
[<matplotlib.axes.\_subplots.AxesSubplot object at 0x1a14754f98>,  
<matplotlib.axes.\_subplots.AxesSubplot object at 0x1a1478ef98>,  
<matplotlib.axes.\_subplots.AxesSubplot object at 0x1a147d2048>]],  
dtype=object)



```
In [ ]: 1
```



```
In [138]: 1 irisData.drop(['species'], axis=1).corr(method='spearman')
```

Out[138]:

	sepal_length	sepal_width	petal_length	petal_width
sepal_length	1.000000	-0.159457	0.881386	0.834421
sepal_width	-0.159457	1.000000	-0.303421	-0.277511
petal_length	0.881386	-0.303421	1.000000	0.936003
petal_width	0.834421	-0.277511	0.936003	1.000000

```
In [22]: 1 #t-test
2 # 1-sample t-test: testing the value of a population mean
3 # scipy.stats.ttest_1samp() tests if the population mean of data is
4 # likely to be equal to a given value
5
6 from scipy import stats
7 stats.ttest_1samp(irisData['sepal_length'], 3)
```

Out[22]: Ttest\_1sampResult(statistic=42.054104134903668, pvalue=1.4781832542930679e-84)

```
In [48]: 1 #2-sample t-test: testing for difference across populations
2 setosaSL = irisData[irisData['species'] == 'setosa']['sepal_length']
3
4 versicolorSL = irisData[irisData['species'] == 'versicolor']
5 ['sepal_length']
6
7 #print(setosaSL)
8 #print(versicolorSL)
9 stats.ttest_ind(setosaSL, versicolorSL)
```

Out[48]: Ttest\_indResult(statistic=-10.520986267549111, pvalue=8.9852350374870789e-18)

```
In [ ]: 1
```

```
In [49]: 1 # Regression: including multiple factors
2 import pandas as pd
3 filepath =
4 'https://gist.githubusercontent.com/curran/a08a1080b88344b0c8a7/raw/d5
5 46eae765268bf2f487608c537c05e22e4b221/iris.csv'
6 irisData = pd.read_csv(filepath)
7 print(irisData.columns)
```

Index(['sepal\_length', 'sepal\_width', 'petal\_length', 'petal\_width',  
      'species'],  
      dtype='object')

```
In [12]: 1 #simple regression
2 from statsmodels.formula.api import ols
3 model = ols("petal_width ~ sepal_width", irisData).fit()
4 print(model.summary())
5
```

# OLS Regression Results

```
=====
=====
Dep. Variable:          petal_width    R-squared:
0.127
Model:                  OLS          Adj. R-squared:
0.121
Method:                Least Squares    F-statistic:
21.55
Date:                  Wed, 20 Jun 2018    Prob (F-statistic):          7.5
2e-06
Time:                  16:42:53    Log-Likelihood:          -1
61.60
No. Observations:          150    AIC:
327.2
Df Residuals:              148    BIC:
333.2
Df Model:                  1
```

Covariance Type: nonrobust

```
=====
=====
              coef      std err          t      P>|t|      [0.025
0.975]
-----
-----
Intercept          3.1152         0.417        7.472      0.000         2.291
3.939
sepal_width       -0.6275         0.135       -4.643      0.000        -0.895
-0.360
=====
=====
```

```
Omnibus:              14.660    Durbin-Watson:
0.523
Prob(Omnibus):          0.001    Jarque-Bera (JB):
6.402
Skew:                   0.266    Prob(JB):
0.0407
Kurtosis:               2.140    Cond. No.
24.3
=====
=====
```

## Warnings:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

```
In [13]: 1 #multiple regression
2 from statsmodels.formula.api import ols
3 model = ols("petal_width ~ sepal_width + petal_length",
4 irisData).fit()
5 print(model.summary())
```

#### OLS Regression Results

```
=====
=====
Dep. Variable:          petal_width    R-squared:
0.930
Model:                  OLS          Adj. R-squared:
0.929
Method:                Least Squares    F-statistic:
972.6
Date:                  Wed, 20 Jun 2018    Prob (F-statistic):          1.7
2e-85
Time:                  16:43:43    Log-Likelihood:          2
7.367
No. Observations:      150    AIC:          -
48.73
Df Residuals:          147    BIC:          -
39.70
Df Model:              2
```

Covariance Type: nonrobust

```
=====
=====
              coef      std err          t      P>|t|      [0.025
0.975]
-----
-----
Intercept      -0.7221      0.151      -4.776      0.000      -1.021
-0.423
sepal_width      0.1033      0.042       2.436      0.016       0.019
0.187
petal_length      0.4271      0.010     40.978      0.000       0.406
0.448
=====
=====
```

```
Omnibus:          3.971    Durbin-Watson:
1.538
Prob(Omnibus):    0.137    Jarque-Bera (JB):
3.764
Skew:            0.244    Prob(JB):
0.152
Kurtosis:        3.604    Cond. No.
48.2
=====
=====
```

#### Warnings:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

In [ ]:

1

In [8]:

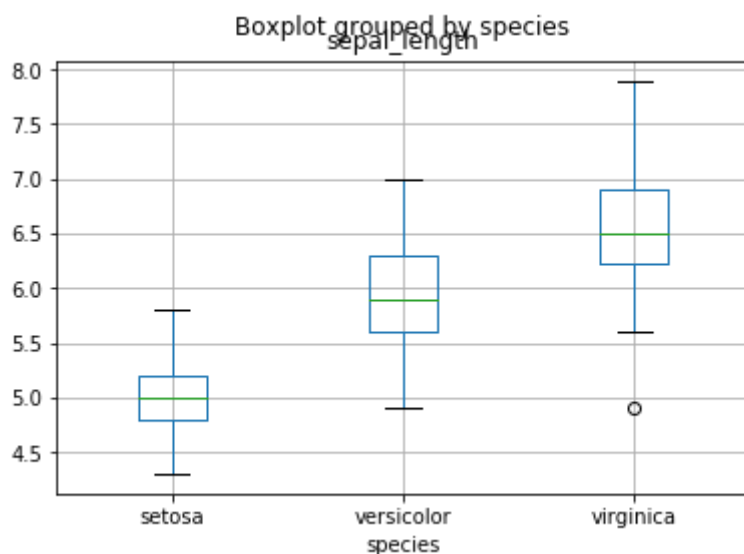
```
1 # Anova
2 irisData.columns
```

Out[8]: Index(['sepal\_length', 'sepal\_width', 'petal\_length', 'petal\_width',  
'species'],  
dtype='object')

In [11]:

```
1 %matplotlib inline
2 irisData.boxplot("sepal_length", by = "species")
3 ## what it means top, middle, bottom
```

Out[11]: <matplotlib.axes.\_subplots.AxesSubplot at 0x1alc8417b8>



In [12]:

```
1 #https://www.r-bloggers.com/anova-%E2%80%93-type-iiiiii-ss-explained/
2 #type 1, 2, 3 sums of squares
3
4
5 # compare R and Python
6 # https://dpaniukov.github.io/2016/10/25/You-and-Your-R-Doing-
  Statistics-in-Python.html
```

In [15]:

```
1 import statsmodels.api as sm
2 from statsmodels.formula.api import ols
3
4 mod = ols('sepal_length ~ species',
5           data=irisData).fit()
6
7 aov_table = sm.stats.anova_lm(mod, typ=2)
8 print (aov_table)
9
```

	sum_sq	df	F	PR(>F)
species	63.212133	2.0	119.264502	1.669669e-31
Residual	38.956200	147.0	NaN	NaN

```
In [18]: 1 esq_sm = aov_table['sum_sq'][0]/(aov_table['sum_sq']
          2 [0]+aov_table['sum_sq'][1])
          3 print(esq_sm)
```

0.618705730738

```
In [ ]: 1
```

```
In [19]: 1 # four different ways of ANOVA
          2 #https://www.marsja.se/four-ways-to-conduct-one-way-anovas-using-
          3 python/
```

```
In [ ]: 1
```

```
In [ ]: 1 #=====
```

```
In [3]: 1 #from sklearn.datasets import load_iris
          2 from sklearn.decomposition import FactorAnalysis
          3 #iris = load_iris()
          4 #X, y = iris.data, iris.target
          5 X = irisData.iloc[:, 0:4]
          6 factor = FactorAnalysis(n_components=4, random_state=101).fit(X)
          7
```

```
In [7]: 1 import pandas as pd
          2 print(pd.DataFrame(factor.components_, columns=X.columns))
```

	sepal_length	sepal_width	petal_length	petal_width
0	0.707227	-0.153147	1.653151	0.701569
1	0.114676	0.159763	-0.045604	-0.014052
2	-0.000000	0.000000	0.000000	0.000000
3	-0.000000	0.000000	0.000000	-0.000000

```
In [11]: 1 from sklearn.decomposition import PCA
          2 import pandas as pd
          3 pca = PCA().fit(X)
          4 print( "Explained variance by component: %s" %
          5       pca.explained_variance_ratio_)
          6 print( " ")
          7 print (pd.DataFrame(pca.components_, columns=X.columns))
```

Explained variance by component: [ 0.92461621 0.05301557 0.01718514 0.00518309]

	sepal_length	sepal_width	petal_length	petal_width
0	0.361590	-0.082269	0.856572	0.358844
1	0.656540	0.729712	-0.175767	-0.074706
2	-0.580997	0.596418	0.072524	0.549061
3	0.317255	-0.324094	-0.479719	0.751121

#more advanced one including varimax rotation <https://pypi.org/project/factor-analyzer/>  
[\(https://pypi.org/project/factor-analyzer/\)](https://pypi.org/project/factor-analyzer/)

**researchgate.net**

Varimax Rotation and Thereafter: Tutorial on PCA Using Linear Algebra, Visualization, and Python Programming for R and Q analysis

In [ ]:

1	
---	--