

CorrelationTests

June 16, 2021

1 Testing correlation between employee counts and # facilities

```
[63]: import pandas as pd
from pandas import read_csv
import numpy as np
from sklearn.linear_model import LinearRegression
from numpy import cov
from scipy.stats import pearsonr
from scipy.stats import spearmanr
import matplotlib.pyplot as plt
import seaborn as sn
```

```
pd.set_option('display.max_columns', None)
pd.set_option('display.max_rows', None)
```

```
[292]: # Group A
groupA=pd.read_csv("CompleteSet_GroupA.csv")
groupA=groupA.drop("Unnamed: 0",axis=1)
#Strip all leading whitespace in Area column
groupA['Area'] = groupA['Area'].apply(lambda x: x.strip())

#Filter only for 2006, 2013 and 2018
groupA = groupA.loc[(groupA['Year'] == 2006) | (groupA['Year'] == 2013)|
↳(groupA['Year']==2018)]

#Remove total NZ row
groupA = groupA.loc[(groupA['Area'] != "Total NZ by Regional Council/
↳Statistical Area")]

#Remove total regions
groupA = groupA.loc[(groupA['ParentArea'] != "NewZealand")]

#Only a certain region
groupA = groupA.loc[(groupA['ParentArea'] == "AucklandRegion")] #Only Auckland
groupA = groupA.loc[(groupA['ParentArea'] == "WaikatoRegion")] #Only Waikato
```

```
#groupA = groupA.loc[(groupA['ParentArea'] == "WellingtonRegion")] #Only
↳Wellington
groupA = groupA.loc[(groupA['ParentArea'] == "OtagoRegion")] #Only Otago
#groupA = groupA.loc[(groupA['ParentArea'] == "BayOfPlentyRegion")] #Only
↳BayOfPlenty

#fill in nans caused
groupA=groupA.fillna(0)
```

Useful websites * <https://realpython.com/linear-regression-in-python/#simple-linear-regression-with-scikit-learn> * <https://machinelearningmastery.com/how-to-use-correlation-to-understand-the-relationship-between-variables/>

1.1 Total industry

1.1.1 Scatterplot

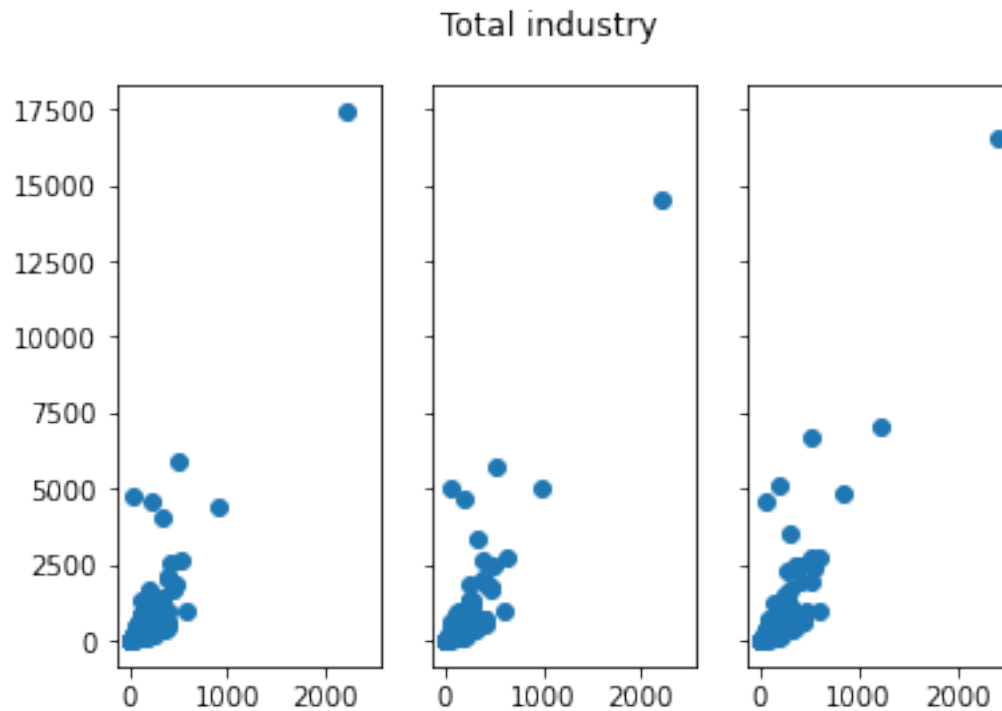
```
[293]: totInd_Fac_2006 = np.array(groupA.loc[(groupA['Year'] == 2006)].
↳TotInd_GeogUnits.tolist())
totInd_Emp_2006 = np.array(groupA.loc[(groupA['Year'] == 2006)].TotInd_EmpCo.
↳tolist())

totInd_Fac_2013 = np.array(groupA.loc[(groupA['Year'] == 2013)].
↳TotInd_GeogUnits.tolist())
totInd_Emp_2013 = np.array(groupA.loc[(groupA['Year'] == 2013)].TotInd_EmpCo.
↳tolist())

totInd_Fac_2018 = np.array(groupA.loc[(groupA['Year'] == 2018)].
↳TotInd_GeogUnits.tolist())
totInd_Emp_2018 = np.array(groupA.loc[(groupA['Year'] == 2018)].TotInd_EmpCo.
↳tolist())
```

```
[294]: fig, (ax1, ax2, ax3) = plt.subplots(1, 3, sharex=True, sharey=True)
fig.suptitle('Total industry')
ax1.scatter(totInd_Fac_2006, totInd_Emp_2006)
ax2.scatter(totInd_Fac_2013, totInd_Emp_2013)
ax3.scatter(totInd_Fac_2018, totInd_Emp_2018)
```

```
[294]: <matplotlib.collections.PathCollection at 0x121b1ab50>
```



1.1.2 Correlation tests

```
[295]: # Covariance

print("2006")
covariance = np.cov([totInd_Fac_2006], [totInd_Emp_2006])
print(covariance)

# Pearson's correlation
corrP, _ = pearsonr(totInd_Fac_2006, totInd_Emp_2006)
print('Pearsons correlation: %.3f' % corrP)

# Spearman's correlation
corrS, _ = spearmanr(totInd_Fac_2006, totInd_Emp_2006)
print('Spearman's correlation: %.3f' % corrS)

print("2013")
covariance = np.cov([totInd_Fac_2013], [totInd_Emp_2013])
print(covariance)

# Pearson's correlation
corrP, _ = pearsonr(totInd_Fac_2013, totInd_Emp_2013)
```

```

print('Pearsons correlation: %.3f' % corrP)

# Spearman's correlation
corrS, _ = spearmanr(totInd_Fac_2013, totInd_Emp_2013)
print('Spearman's correlation: %.3f' % corrS)

print("2018")
covariance = np.cov([totInd_Fac_2018], [totInd_Emp_2018])
print(covariance)

# Pearson's correlation
corrP, _ = pearsonr(totInd_Fac_2018, totInd_Emp_2018)
print('Pearsons correlation: %.3f' % corrP)

# Spearman's correlation
corrS, _ = spearmanr(totInd_Fac_2018, totInd_Emp_2018)
print('Spearman's correlation: %.3f' % corrS)

```

```

2006
[[ 50076.38382353  333027.45392157]
 [ 333027.45392157 2987320.39803922]]
Pearsons correlation: 0.861
Spearman's correlation: 0.774
2013
[[ 53593.46617647  300580.34101307]
 [ 300580.34101307 2368210.04264706]]
Pearsons correlation: 0.844
Spearman's correlation: 0.784
2018
[[ 66749.44117647  396096.83169935]
 [ 396096.83169935 3179827.8417756 ]]
Pearsons correlation: 0.860
Spearman's correlation: 0.785

```

1.1.3 Linear regression

regressor - # Facilities ; predictor - employee count

```

[296]: totInd_Fac_2006 = totInd_Fac_2006.reshape((-1, 1))
      totInd_Fac_2013 = totInd_Fac_2013.reshape((-1, 1))
      totInd_Fac_2018 = totInd_Fac_2018.reshape((-1, 1))

```

1.1.4 Create model

```
[297]: model_2006 = LinearRegression().fit(totInd_Fac_2006, totInd_Emp_2006)
model_2013 = LinearRegression().fit(totInd_Fac_2013, totInd_Emp_2013)
model_2018 = LinearRegression().fit(totInd_Fac_2018, totInd_Emp_2018)
```

```
[298]: print("2006")
r_sq = model_2006.score(totInd_Fac_2006, totInd_Emp_2006)
print('coefficient of determination:', r_sq)
print('intercept:', model_2006.intercept_)
print('slope:', model_2006.coef_)

print("2013")
r_sq = model_2013.score(totInd_Fac_2013, totInd_Emp_2013)
print('coefficient of determination:', r_sq)
print('intercept:', model_2013.intercept_)
print('slope:', model_2013.coef_)

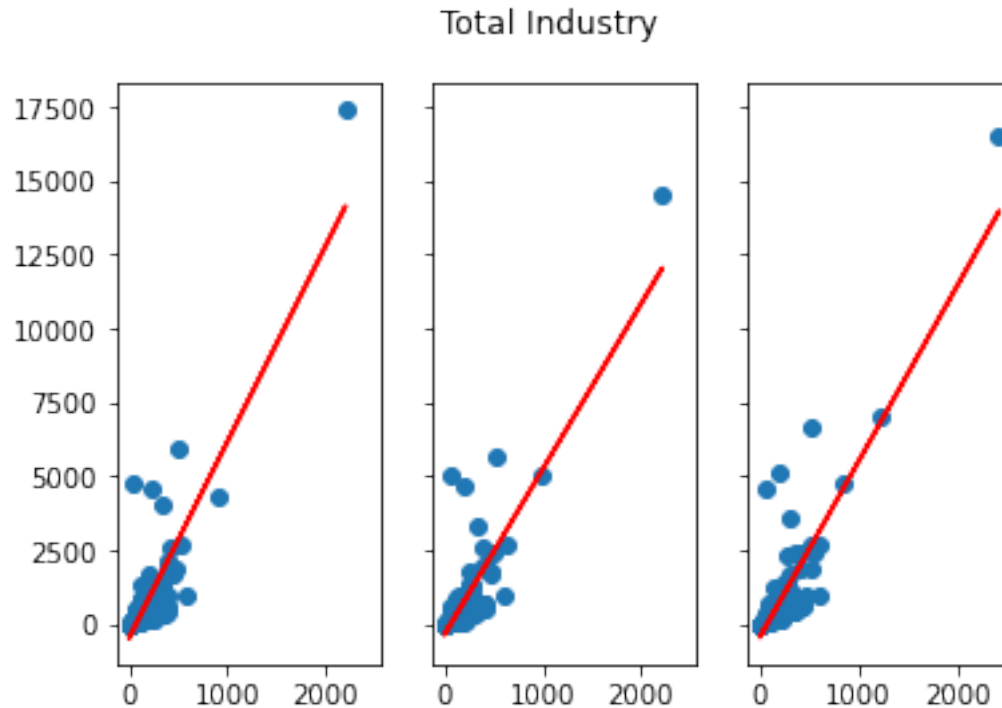
print("2018")
r_sq = model_2018.score(totInd_Fac_2018, totInd_Emp_2018)
print('coefficient of determination:', r_sq)
print('intercept:', model_2018.intercept_)
print('slope:', model_2018.coef_)
```

```
2006
coefficient of determination: 0.7413875868718665
intercept: -478.69491104780195
slope: [6.65038943]
2013
coefficient of determination: 0.7118509698967026
intercept: -359.2368427154354
slope: [5.60852586]
2018
coefficient of determination: 0.7391821917837775
intercept: -438.4199272873399
slope: [5.93408461]
```

```
[299]: fig, (ax1, ax2, ax3) = plt.subplots(1, 3, sharex=True, sharey=True)
fig.suptitle('Total Industry')
ax1.scatter(totInd_Fac_2006, totInd_Emp_2006)
ax1.plot(totInd_Fac_2006, model_2006.coef_*totInd_Fac_2006+model_2006.
    ↳intercept_, 'r')
ax2.scatter(totInd_Fac_2013, totInd_Emp_2013)
ax2.plot(totInd_Fac_2013, model_2013.coef_*totInd_Fac_2013+model_2013.
    ↳intercept_, 'r')
ax3.scatter(totInd_Fac_2018, totInd_Emp_2018)
```

```
ax3.plot(totInd_Fac_2018,model_2018.coef_*totInd_Fac_2018+model_2018.
↪intercept_,'r')
```

[299]: [<matplotlib.lines.Line2D at 0x121d73f40>]



1.2 Wholesale (F)

1.2.1 Scatterplot

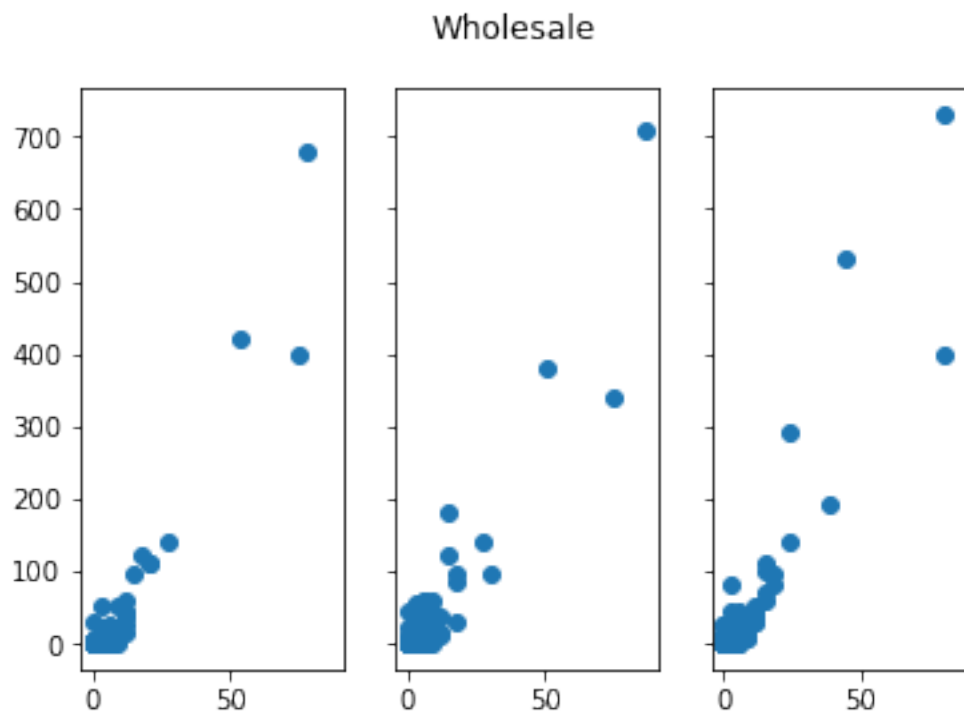
```
[300]: whole_Fac_2006 = np.array(groupA.loc[(groupA['Year'] == 2006)].F_GeogUnits.
↪tolist())
whole_Emp_2006 = np.array(groupA.loc[(groupA['Year'] == 2006)].F_EmpCo.tolist())

whole_Fac_2013 = np.array(groupA.loc[(groupA['Year'] == 2013)].F_GeogUnits.
↪tolist())
whole_Emp_2013 = np.array(groupA.loc[(groupA['Year'] == 2013)].F_EmpCo.tolist())

whole_Fac_2018 = np.array(groupA.loc[(groupA['Year'] == 2018)].F_GeogUnits.
↪tolist())
whole_Emp_2018 = np.array(groupA.loc[(groupA['Year'] == 2018)].F_EmpCo.tolist())
```

```
[301]: fig, (ax1, ax2, ax3) = plt.subplots(1, 3, sharex=True, sharey=True)
fig.suptitle('Wholesale')
ax1.scatter(whole_Fac_2006, whole_Emp_2006)
ax2.scatter(whole_Fac_2013, whole_Emp_2013)
ax3.scatter(whole_Fac_2018, whole_Emp_2018)
```

```
[301]: <matplotlib.collections.PathCollection at 0x121ed2a30>
```



1.2.2 Correlation tests

```
[302]: # Covariance

print("2006")
covariance = np.cov([whole_Fac_2006], [whole_Emp_2006])
print(covariance)

# Pearson's correlation
corrP, _ = pearsonr(whole_Fac_2006, whole_Emp_2006)
print('Pearsons correlation: %.3f' % corrP)

# Spearman's correlation
corrS, _ = spearmanr(whole_Fac_2006, whole_Emp_2006)
```

```

print('Spearman's correlation: %.3f' % corrS)

print("2013")
covariance = np.cov([whole_Fac_2013], [whole_Emp_2013])
print(covariance)

# Pearson's correlation
corrP, _ = pearsonr(whole_Fac_2013, whole_Emp_2013)
print('Pearson's correlation: %.3f' % corrP)

# Spearman's correlation
corrS, _ = spearmanr(whole_Fac_2013, whole_Emp_2013)
print('Spearman's correlation: %.3f' % corrS)

print("2018")
covariance = np.cov([whole_Fac_2018], [whole_Emp_2018])
print(covariance)

# Pearson's correlation
corrP, _ = pearsonr(whole_Fac_2018, whole_Emp_2018)
print('Pearson's correlation: %.3f' % corrP)

# Spearman's correlation
corrS, _ = spearmanr(whole_Fac_2018, whole_Emp_2018)
print('Spearman's correlation: %.3f' % corrS)

```

```

2006
[[ 115.30931373  798.1874183 ]
 [ 798.1874183  6129.3119281 ]]
Pearson's correlation: 0.949
Spearman's correlation: 0.650
2013
[[ 128.1995098  817.57941176]
 [ 817.57941176 6048.39455338]]
Pearson's correlation: 0.928
Spearman's correlation: 0.652
2018
[[ 130.72107843  938.94362745]
 [ 938.94362745 8044.47401961]]
Pearson's correlation: 0.916
Spearman's correlation: 0.713

```

1.2.3 Linear regression

regressor - # Facilities ; predictor - employee count

1.2.4 Create model

```
[303]: whole_Fac_2006 = whole_Fac_2006.reshape((-1, 1))
       whole_Fac_2013 = whole_Fac_2013.reshape((-1, 1))
       whole_Fac_2018 = whole_Fac_2018.reshape((-1, 1))
```

```
[304]: model_2006 = LinearRegression().fit(whole_Fac_2006, whole_Emp_2006)
       model_2013 = LinearRegression().fit(whole_Fac_2013, whole_Emp_2013)
       model_2018 = LinearRegression().fit(whole_Fac_2018, whole_Emp_2018)
```

```
[305]: print("2006")
       r_sq = model_2006.score(whole_Fac_2006, whole_Emp_2006)
       print('coefficient of determination:', r_sq)
       print('intercept:', model_2006.intercept_)
       print('slope:', model_2006.coef_)

       print("2013")
       r_sq = model_2013.score(whole_Fac_2013, whole_Emp_2013)
       print('coefficient of determination:', r_sq)
       print('intercept:', model_2013.intercept_)
       print('slope:', model_2013.coef_)

       print("2018")
       r_sq = model_2018.score(whole_Fac_2018, whole_Emp_2018)
       print('coefficient of determination:', r_sq)
       print('intercept:', model_2018.intercept_)
       print('slope:', model_2018.coef_)
```

```
2006
coefficient of determination: 0.901433397786837
intercept: -16.57922212633539
slope: [6.92214178]
2013
coefficient of determination: 0.8620519029651237
intercept: -10.195551510933843
slope: [6.37739889]
2018
coefficient of determination: 0.8383700632721791
intercept: -11.244578525598953
slope: [7.18280203]
```

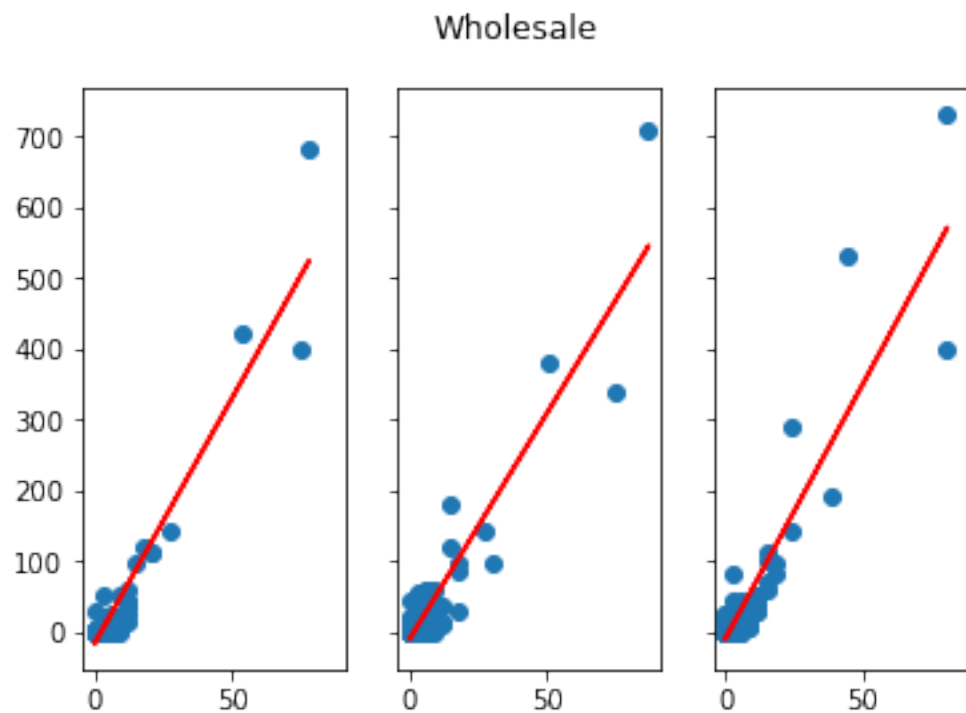
```
[306]: fig, (ax1, ax2, ax3) = plt.subplots(1, 3, sharex=True, sharey=True)
       fig.suptitle('Wholesale')
       ax1.scatter(whole_Fac_2006, whole_Emp_2006)
```

```

ax1.plot(whole_Fac_2006,model_2006.coef_*whole_Fac_2006+model_2006.
↪intercept_,'r')
ax2.scatter(whole_Fac_2013, whole_Emp_2013)
ax2.plot(whole_Fac_2013,model_2013.coef_*whole_Fac_2013+model_2013.
↪intercept_,'r')
ax3.scatter(whole_Fac_2018, whole_Emp_2018)
ax3.plot(whole_Fac_2018,model_2018.coef_*whole_Fac_2018+model_2018.
↪intercept_,'r')

```

[306]: [<matplotlib.lines.Line2D at 0x121c4dfa0>]



1.3 Retail (G)

1.3.1 Scatterplot

```

[307]: retail_Fac_2006 = np.array(groupA.loc[(groupA['Year'] == 2006)].G_GeogUnits.
↪tolist())
retail_Emp_2006 = np.array(groupA.loc[(groupA['Year'] == 2006)].G_EmpCo.
↪tolist())

retail_Fac_2013 = np.array(groupA.loc[(groupA['Year'] == 2013)].G_GeogUnits.
↪tolist())

```

```

retail_Emp_2013 = np.array(groupA.loc[(groupA['Year'] == 2013)].G_EmpCo.
    ↳tolist())

retail_Fac_2018 = np.array(groupA.loc[(groupA['Year'] == 2018)].G_GeogUnits.
    ↳tolist())
retail_Emp_2018 = np.array(groupA.loc[(groupA['Year'] == 2018)].G_EmpCo.
    ↳tolist())

```

```

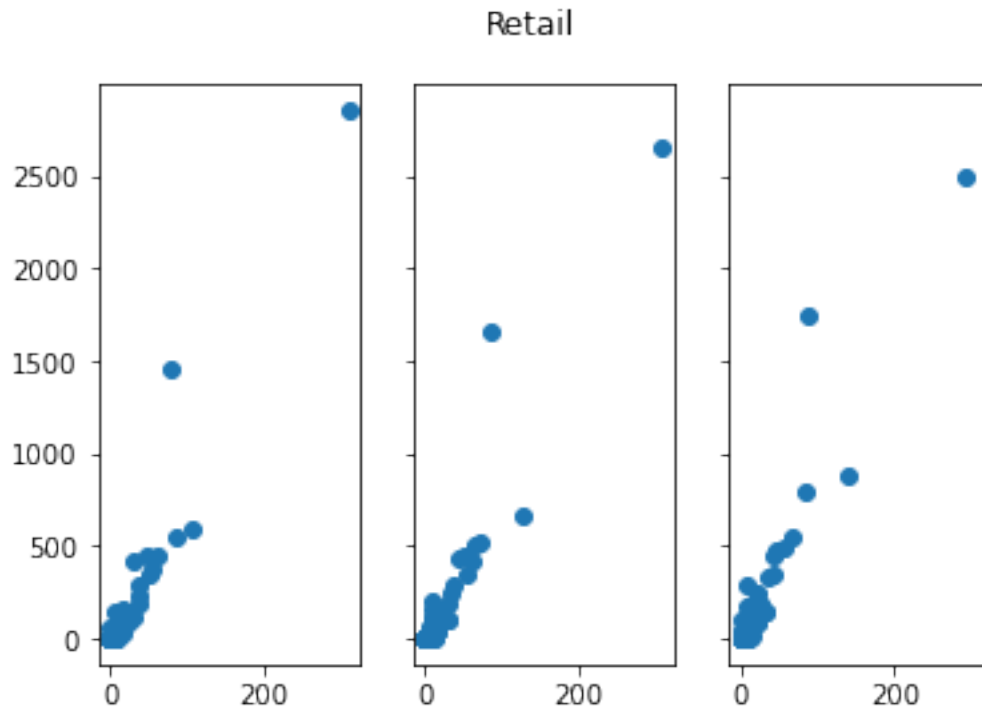
[308]: fig, (ax1, ax2, ax3) = plt.subplots(1, 3, sharex=True, sharey=True)
fig.suptitle('Retail')
ax1.scatter(retail_Fac_2006, retail_Emp_2006)
ax2.scatter(retail_Fac_2013, retail_Emp_2013)
ax3.scatter(retail_Fac_2018, retail_Emp_2018)

```

```

[308]: <matplotlib.collections.PathCollection at 0x1220b8d90>

```



1.3.2 Correlation tests

```

[309]: # Covariance

print("2006")
covariance = np.cov([retail_Fac_2006], [retail_Emp_2006])

```

```

print(covariance)

# Pearson's correlation
corrP, _ = pearsonr(retail_Fac_2006, retail_Emp_2006)
print('Pearsons correlation: %.3f' % corrP)

# Spearman's correlation
corrS, _ = spearmanr(retail_Fac_2006, retail_Emp_2006)
print('Spearman's correlation: %.3f' % corrS)

print("2013")
covariance = np.cov([retail_Fac_2013], [retail_Emp_2013])
print(covariance)

# Pearson's correlation
corrP, _ = pearsonr(retail_Fac_2013, retail_Emp_2013)
print('Pearsons correlation: %.3f' % corrP)

# Spearman's correlation
corrS, _ = spearmanr(retail_Fac_2013, retail_Emp_2013)
print('Spearman's correlation: %.3f' % corrS)

print("2018")
covariance = np.cov([retail_Fac_2018], [retail_Emp_2018])
print(covariance)

# Pearson's correlation
corrP, _ = pearsonr(retail_Fac_2018, retail_Emp_2018)
print('Pearsons correlation: %.3f' % corrP)

# Spearman's correlation
corrS, _ = spearmanr(retail_Fac_2018, retail_Emp_2018)
print('Spearman's correlation: %.3f' % corrS)

```

```

2006
[[ 922.26470588 8372.51895425]
 [ 8372.51895425 83032.84270152]]
Pearsons correlation: 0.957
Spearman's correlation: 0.836
2013
[[ 942.38431373 8198.2         ]
 [ 8198.2         80748.37685185]]
Pearsons correlation: 0.940
Spearman's correlation: 0.843
2018
[[ 907.66862745 8150.69934641]
 [ 8150.69934641 83251.55206972]]

```

Pearsons correlation: 0.938
Spearman's correlation: 0.812

1.3.3 Linear regression

regressor - # Facilities ; predictor - employee count

1.3.4 Create model

```
[310]: retail_Fac_2006 = retail_Fac_2006.reshape((-1, 1))
retail_Fac_2013 = retail_Fac_2013.reshape((-1, 1))
retail_Fac_2018 = retail_Fac_2018.reshape((-1, 1))
```

```
[311]: model_2006 = LinearRegression().fit(retail_Fac_2006, retail_Emp_2006)
model_2013 = LinearRegression().fit(retail_Fac_2013, retail_Emp_2013)
model_2018 = LinearRegression().fit(retail_Fac_2018, retail_Emp_2018)
```

```
[312]: print("2006")
r_sq = model_2006.score(retail_Fac_2006, retail_Emp_2006)
print('coefficient of determination:', r_sq)
print('intercept:', model_2006.intercept_)
print('slope:', model_2006.coef_)

print("2013")
r_sq = model_2013.score(retail_Fac_2013, retail_Emp_2013)
print('coefficient of determination:', r_sq)
print('intercept:', model_2013.intercept_)
print('slope:', model_2013.coef_)

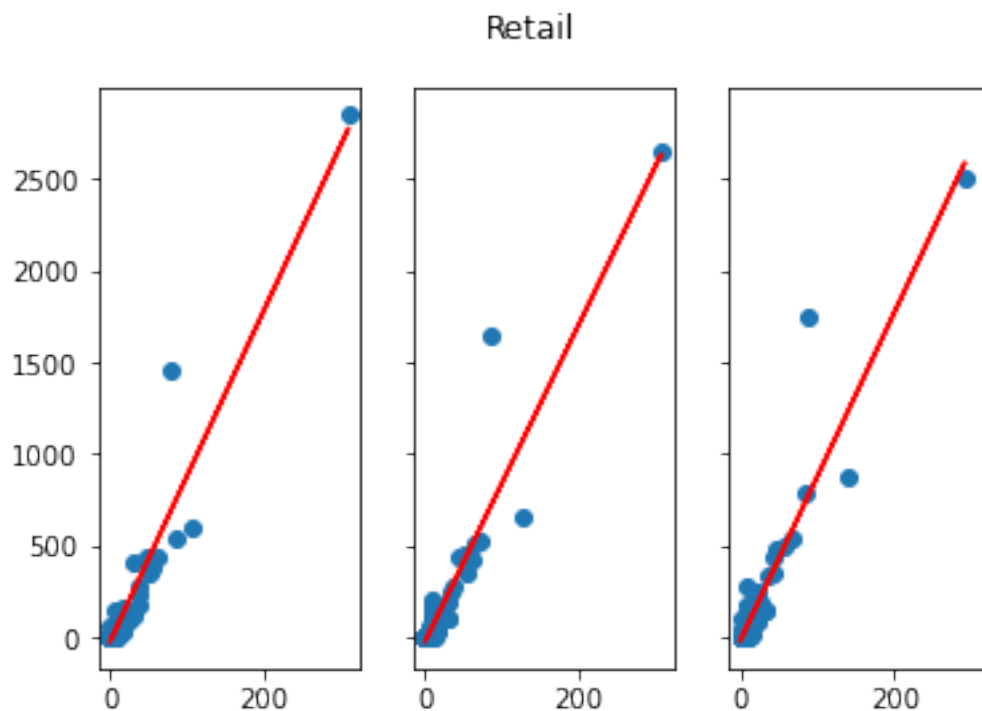
print("2018")
r_sq = model_2018.score(retail_Fac_2018, retail_Emp_2018)
print('coefficient of determination:', r_sq)
print('intercept:', model_2018.intercept_)
print('slope:', model_2018.coef_)
```

```
2006
coefficient of determination: 0.9153913056777588
intercept: -28.861621541176348
slope: [9.07821681]
2013
coefficient of determination: 0.8832327582070519
intercept: -27.873655891605807
slope: [8.69942324]
2018
```

```
coefficient of determination: 0.8791643793621488
intercept: -23.34918591262688
slope: [8.97981829]
```

```
[313]: fig, (ax1, ax2, ax3) = plt.subplots(1, 3, sharex=True, sharey=True)
fig.suptitle('Retail')
ax1.scatter(retail_Fac_2006, retail_Emp_2006)
ax1.plot(retail_Fac_2006, model_2006.coef_*retail_Fac_2006+model_2006.
↪intercept_, 'r')
ax2.scatter(retail_Fac_2013, retail_Emp_2013)
ax2.plot(retail_Fac_2013, model_2013.coef_*retail_Fac_2013+model_2013.
↪intercept_, 'r')
ax3.scatter(retail_Fac_2018, retail_Emp_2018)
ax3.plot(retail_Fac_2018, model_2018.coef_*retail_Fac_2018+model_2018.
↪intercept_, 'r')
```

```
[313]: [<matplotlib.lines.Line2D at 0x120e1c760>]
```



1.4 TransPostWare

```
[314]: groupA['TransPostWare_GeogUnits']=groupA['I461_GeogUnits']+groupA['I471_GeogUnits']+groupA['I481_GeogUnits']
groupA['TransPostWare_EmpCo']=groupA['I461_EmpCo']+groupA['I471_EmpCo']+groupA['I481_EmpCo']
```

1.4.1 Scatterplot

```
[315]: tpw_Fac_2006 = np.array(groupA.loc[(groupA['Year'] == 2006)].
      ↪TransPostWare_GeogUnits.tolist())
tpw_Emp_2006 = np.array(groupA.loc[(groupA['Year'] == 2006)].
      ↪TransPostWare_EmpCo.tolist())

tpw_Fac_2013 = np.array(groupA.loc[(groupA['Year'] == 2013)].
      ↪TransPostWare_GeogUnits.tolist())
tpw_Emp_2013 = np.array(groupA.loc[(groupA['Year'] == 2013)].
      ↪TransPostWare_EmpCo.tolist())

tpw_Fac_2018 = np.array(groupA.loc[(groupA['Year'] == 2018)].
      ↪TransPostWare_GeogUnits.tolist())
tpw_Emp_2018 = np.array(groupA.loc[(groupA['Year'] == 2018)].
      ↪TransPostWare_EmpCo.tolist())
```

```
[316]: fig, (ax1, ax2, ax3) = plt.subplots(1, 3, sharex=True, sharey=True)
fig.suptitle('Transport Post Warehouse')
ax1.scatter(tpw_Fac_2006, tpw_Emp_2006)
ax2.scatter(tpw_Fac_2013, tpw_Emp_2013)
ax3.scatter(tpw_Fac_2018, tpw_Emp_2018)
```

```
[316]: <matplotlib.collections.PathCollection at 0x12238dca0>
```



1.4.2 Correlation tests

```
[317]: # Covariance

print("2006")
covariance = np.cov([tpw_Fac_2006], [tpw_Emp_2006])
print(covariance)

# Pearson's correlation
corrP, _ = pearsonr(tpw_Fac_2006, tpw_Emp_2006)
print('Pearsons correlation: %.3f' % corrP)

# Spearman's correlation
corrS, _ = spearmanr(tpw_Fac_2006, tpw_Emp_2006)
print('Spearman's correlation: %.3f' % corrS)

print("2013")
covariance = np.cov([tpw_Fac_2013], [tpw_Emp_2013])
print(covariance)

# Pearson's correlation
corrP, _ = pearsonr(tpw_Fac_2013, tpw_Emp_2013)
print('Pearsons correlation: %.3f' % corrP)

# Spearman's correlation
corrS, _ = spearmanr(tpw_Fac_2013, tpw_Emp_2013)
print('Spearman's correlation: %.3f' % corrS)

print("2018")
covariance = np.cov([tpw_Fac_2018], [tpw_Emp_2018])
print(covariance)

# Pearson's correlation
corrP, _ = pearsonr(tpw_Fac_2018, tpw_Emp_2018)
print('Pearsons correlation: %.3f' % corrP)

# Spearman's correlation
corrS, _ = spearmanr(tpw_Fac_2018, tpw_Emp_2018)
print('Spearman's correlation: %.3f' % corrS)
```

2006

```
[[ 18.10196078 157.56078431]
 [157.56078431 2802.82783224]]
```

Pearsons correlation: 0.699

Spearman's correlation: 0.664


```

2013
[[ 17.71715686 107.12303922]
 [ 107.12303922 1429.84395425]]
Pearsons correlation: 0.673
Spearman's correlation: 0.626
2018
[[ 19.52941176 115.68823529]
 [ 115.68823529 1530.63044662]]
Pearsons correlation: 0.669
Spearman's correlation: 0.582

```

1.4.3 Linear regression

regressor - # Facilities ; predictor - employee count

1.4.4 Create model

```

[318]: tpw_Fac_2006 = tpw_Fac_2006.reshape((-1, 1))
tpw_Fac_2013 = tpw_Fac_2013.reshape((-1, 1))
tpw_Fac_2018 = tpw_Fac_2018.reshape((-1, 1))

```

```

[319]: model_2006 = LinearRegression().fit(tpw_Fac_2006, tpw_Emp_2006)
model_2013 = LinearRegression().fit(tpw_Fac_2013, tpw_Emp_2013)
model_2018 = LinearRegression().fit(tpw_Fac_2018, tpw_Emp_2018)

```

```

[320]: print("2006")
r_sq = model_2006.score(tpw_Fac_2006, tpw_Emp_2006)
print('coefficient of determination:', r_sq)
print('intercept:', model_2006.intercept_)
print('slope:', model_2006.coef_)

print("2013")
r_sq = model_2013.score(tpw_Fac_2013, tpw_Emp_2013)
print('coefficient of determination:', r_sq)
print('intercept:', model_2013.intercept_)
print('slope:', model_2013.coef_)

print("2018")
r_sq = model_2018.score(tpw_Fac_2018, tpw_Emp_2018)
print('coefficient of determination:', r_sq)
print('intercept:', model_2018.intercept_)
print('slope:', model_2018.coef_)

```

2006
 coefficient of determination: 0.48929888585578496
 intercept: -11.552642980935879
 slope: [8.70407279]
 2013
 coefficient of determination: 0.45298424530439696
 intercept: -6.9176604045043195
 slope: [6.04628835]
 2018
 coefficient of determination: 0.44773277064596595
 intercept: -4.920331325301211
 slope: [5.92379518]

```
[321]: fig, (ax1, ax2, ax3) = plt.subplots(1, 3, sharex=True, sharey=True)
fig.suptitle('Transport Post Warehouse')
ax1.scatter(tpw_Fac_2006, tpw_Emp_2006)
ax1.plot(tpw_Fac_2006, model_2006.coef_*tpw_Fac_2006+model_2006.intercept_, 'r')
ax2.scatter(tpw_Fac_2013, tpw_Emp_2013)
ax2.plot(tpw_Fac_2013, model_2013.coef_*tpw_Fac_2013+model_2013.intercept_, 'r')
ax3.scatter(tpw_Fac_2018, tpw_Emp_2018)
ax3.plot(tpw_Fac_2018, model_2018.coef_*tpw_Fac_2018+model_2018.intercept_, 'r')
```

[321]: [<matplotlib.lines.Line2D at 0x12252b130>]



1.5 Transport

```
[322]: groupA['Transport_GeogUnits']=groupA['I461_GeogUnits']+groupA['I471_GeogUnits']+groupA['I481_GeogUnits']
groupA['Transport_EmpCo']=groupA['I461_EmpCo']+groupA['I471_EmpCo']+groupA['I481_EmpCo']
```

1.5.1 Scatterplot

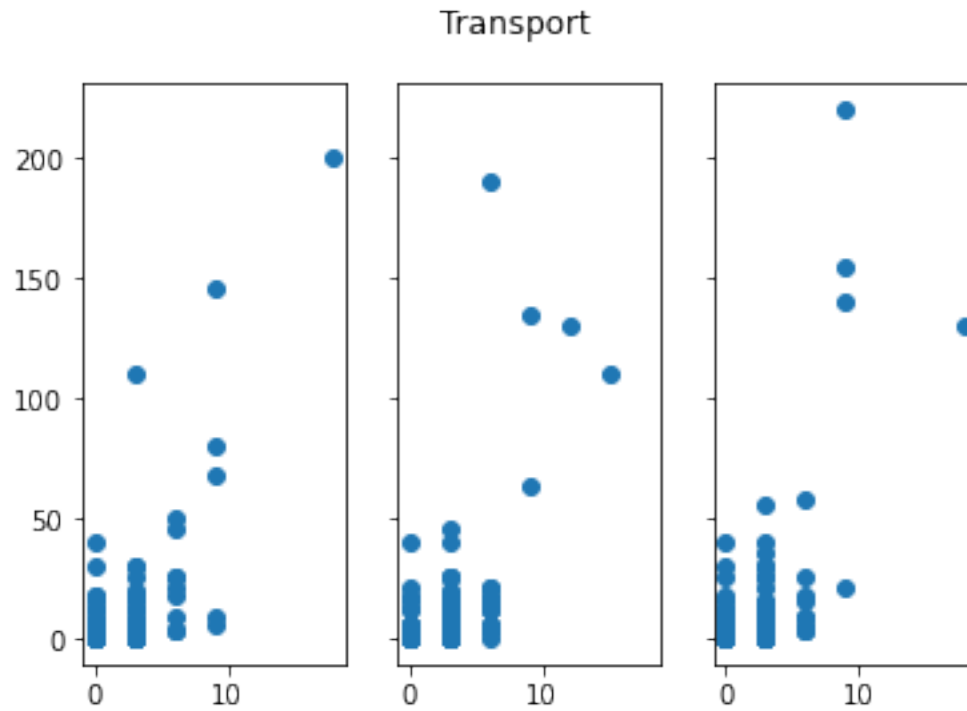
```
[323]: trans_Fac_2006 = np.array(groupA.loc[(groupA['Year'] == 2006)].
↳Transport_GeogUnits.tolist())
trans_Emp_2006 = np.array(groupA.loc[(groupA['Year'] == 2006)].Transport_EmpCo.
↳tolist())

trans_Fac_2013 = np.array(groupA.loc[(groupA['Year'] == 2013)].
↳Transport_GeogUnits.tolist())
trans_Emp_2013 = np.array(groupA.loc[(groupA['Year'] == 2013)].Transport_EmpCo.
↳tolist())

trans_Fac_2018 = np.array(groupA.loc[(groupA['Year'] == 2018)].
↳Transport_GeogUnits.tolist())
trans_Emp_2018 = np.array(groupA.loc[(groupA['Year'] == 2018)].Transport_EmpCo.
↳tolist())
```

```
[324]: fig, (ax1, ax2, ax3) = plt.subplots(1, 3, sharex=True, sharey=True)
fig.suptitle('Transport')
ax1.scatter(trans_Fac_2006, trans_Emp_2006)
ax2.scatter(trans_Fac_2013, trans_Emp_2013)
ax3.scatter(trans_Fac_2018, trans_Emp_2018)
```

```
[324]: <matplotlib.collections.PathCollection at 0x12269abe0>
```



1.5.2 Correlation tests

```
[325]: # Covariance

print("2006")
covariance = np.cov([trans_Fac_2006], [trans_Emp_2006])
print(covariance)

# Pearson's correlation
corrP, _ = pearsonr(trans_Fac_2006, trans_Emp_2006)
print('Pearsons correlation: %.3f' % corrP)

# Spearman's correlation
corrS, _ = spearmanr(trans_Fac_2006, trans_Emp_2006)
print('Spearman's correlation: %.3f' % corrS)

print("2013")
covariance = np.cov([trans_Fac_2013], [trans_Emp_2013])
print(covariance)

# Pearson's correlation
corrP, _ = pearsonr(trans_Fac_2013, trans_Emp_2013)
```

```

print('Pearsons correlation: %.3f' % corrP)

# Spearman's correlation
corrS, _ = spearmanr(trans_Fac_2013, trans_Emp_2013)
print('Spearman's correlation: %.3f' % corrS)

print("2018")
covariance = np.cov([trans_Fac_2018], [trans_Emp_2018])
print(covariance)

# Pearson's correlation
corrP, _ = pearsonr(trans_Fac_2018, trans_Emp_2018)
print('Pearsons correlation: %.3f' % corrP)

# Spearman's correlation
corrS, _ = spearmanr(trans_Fac_2018, trans_Emp_2018)
print('Spearman's correlation: %.3f' % corrS)

```

```

2006
[[ 7.51715686  47.53447712]
 [ 47.53447712 648.12674292]]
Pearsons correlation: 0.681
Spearman's correlation: 0.613
2013
[[ 7.20784314  42.33006536]
 [ 42.33006536 652.54765795]]
Pearsons correlation: 0.617
Spearman's correlation: 0.665
2018
[[ 7.16862745  51.25294118]
 [ 51.25294118 848.49885621]]
Pearsons correlation: 0.657
Spearman's correlation: 0.647

```

1.5.3 Linear regression

regressor - # Facilities ; predictor - employee count

1.5.4 Create model

```

[326]: trans_Fac_2006 = trans_Fac_2006.reshape((-1, 1))
trans_Fac_2013 = trans_Fac_2013.reshape((-1, 1))
trans_Fac_2018 = trans_Fac_2018.reshape((-1, 1))

```

```
[327]: model_2006 = LinearRegression().fit(trans_Fac_2006, trans_Emp_2006)
model_2013 = LinearRegression().fit(trans_Fac_2013, trans_Emp_2013)
model_2018 = LinearRegression().fit(trans_Fac_2018, trans_Emp_2018)
```

```
[328]: print("2006")
r_sq = model_2006.score(trans_Fac_2006, trans_Emp_2006)
print('coefficient of determination:', r_sq)
print('intercept:', model_2006.intercept_)
print('slope:', model_2006.coef_)

print("2013")
r_sq = model_2013.score(trans_Fac_2013, trans_Emp_2013)
print('coefficient of determination:', r_sq)
print('intercept:', model_2013.intercept_)
print('slope:', model_2013.coef_)

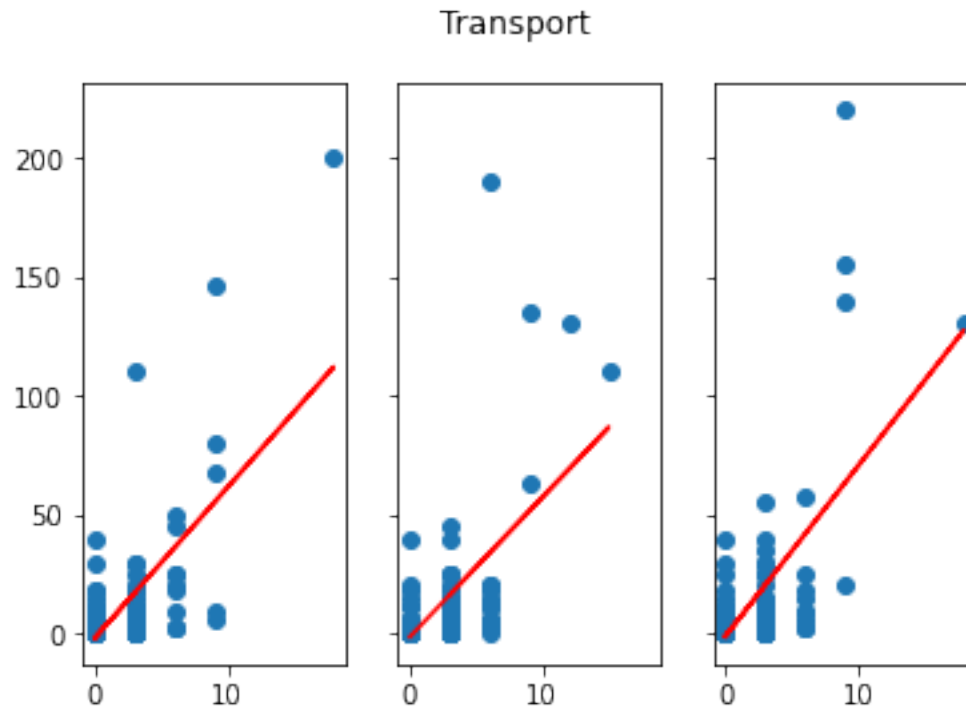
print("2018")
r_sq = model_2018.score(trans_Fac_2018, trans_Emp_2018)
print('coefficient of determination:', r_sq)
print('intercept:', model_2018.intercept_)
print('slope:', model_2018.coef_)
```

```
2006
coefficient of determination: 0.4637713195811307
intercept: -1.8794261493315965
slope: [6.32346484]
2013
coefficient of determination: 0.38096082291467825
intercept: -1.4039717083786751
slope: [5.87277838]
2018
coefficient of determination: 0.43186729166508797
intercept: -1.3464168490153163
slope: [7.14961707]
```

```
[329]: fig, (ax1, ax2, ax3) = plt.subplots(1, 3, sharex=True, sharey=True)
fig.suptitle('Transport')
ax1.scatter(trans_Fac_2006, trans_Emp_2006)
ax1.plot(trans_Fac_2006, model_2006.coef_*trans_Fac_2006+model_2006.
    ↳intercept_, 'r')
ax2.scatter(trans_Fac_2013, trans_Emp_2013)
ax2.plot(trans_Fac_2013, model_2013.coef_*trans_Fac_2013+model_2013.
    ↳intercept_, 'r')
ax3.scatter(trans_Fac_2018, trans_Emp_2018)
```

```
ax3.plot(trans_Fac_2018,model_2018.coef_*trans_Fac_2018+model_2018.
        ↪intercept_,'r')
```

[329]: [<matplotlib.lines.Line2D at 0x122820160>]



1.6 Post and storage

```
[330]: groupA['Storage_GeogUnits']=groupA['I51_GeogUnits']+groupA['I53_GeogUnits']
        groupA['Storage_EmpCo']=groupA['I51_EmpCo']+groupA['I53_EmpCo']
```

1.6.1 Scatterplot

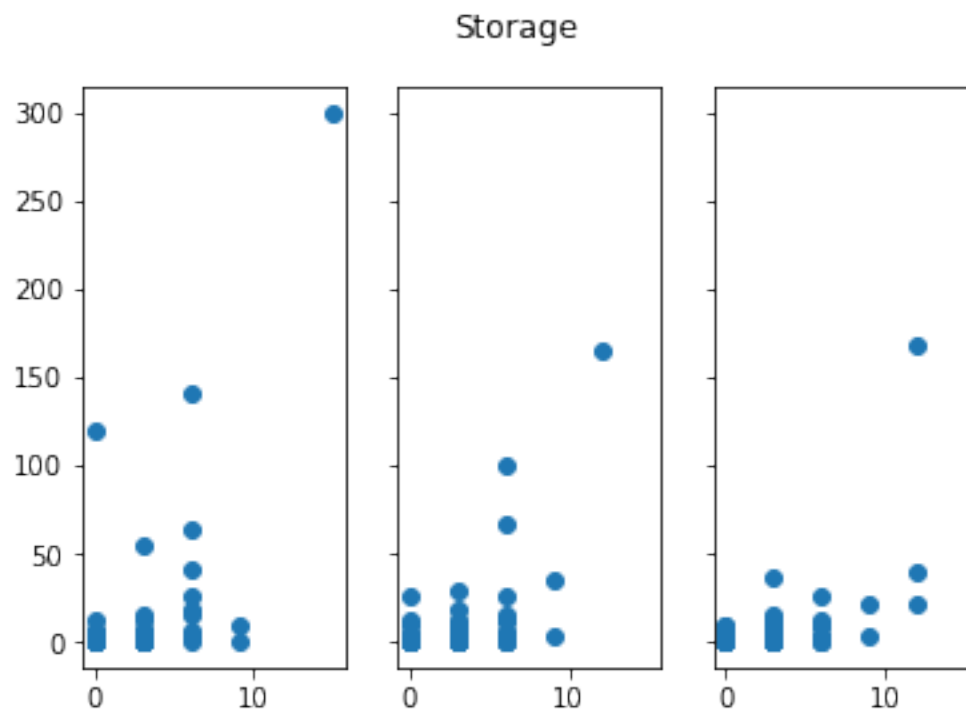
```
[331]: stor_Fac_2006 = np.array(groupA.loc[(groupA['Year'] == 2006)].Storage_GeogUnits.
        ↪tolist())
        stor_Emp_2006 = np.array(groupA.loc[(groupA['Year'] == 2006)].Storage_EmpCo.
        ↪tolist())

        stor_Fac_2013 = np.array(groupA.loc[(groupA['Year'] == 2013)].Storage_GeogUnits.
        ↪tolist())
        stor_Emp_2013 = np.array(groupA.loc[(groupA['Year'] == 2013)].Storage_EmpCo.
        ↪tolist())
```

```
stor_Fac_2018 = np.array(groupA.loc[(groupA['Year'] == 2018)].Storage_GeogUnits.
    ↳tolist())
stor_Emp_2018 = np.array(groupA.loc[(groupA['Year'] == 2018)].Storage_EmpCo.
    ↳tolist())
```

```
[332]: fig, (ax1, ax2, ax3) = plt.subplots(1, 3, sharex=True, sharey=True)
fig.suptitle('Storage')
ax1.scatter(stor_Fac_2006, stor_Emp_2006)
ax2.scatter(stor_Fac_2013, stor_Emp_2013)
ax3.scatter(stor_Fac_2018, stor_Emp_2018)
```

```
[332]: <matplotlib.collections.PathCollection at 0x122986b80>
```



1.6.2 Correlation tests

```
[333]: # Covariance
print("2006")
covariance = np.cov([stor_Fac_2006], [stor_Emp_2006])
print(covariance)
```



```

# Pearson's correlation
corrP, _ = pearsonr(stor_Fac_2006, stor_Emp_2006)
print('Pearsons correlation: %.3f' % corrP)

# Spearman's correlation
corrS, _ = spearmanr(stor_Fac_2006, stor_Emp_2006)
print('Spearman's correlation: %.3f' % corrS)

print("2013")
covariance = np.cov([stor_Fac_2013], [stor_Emp_2013])
print(covariance)

# Pearson's correlation
corrP, _ = pearsonr(stor_Fac_2013, stor_Emp_2013)
print('Pearsons correlation: %.3f' % corrP)

# Spearman's correlation
corrS, _ = spearmanr(stor_Fac_2013, stor_Emp_2013)
print('Spearman's correlation: %.3f' % corrS)

print("2018")
covariance = np.cov([stor_Fac_2018], [stor_Emp_2018])
print(covariance)

# Pearson's correlation
corrP, _ = pearsonr(stor_Fac_2018, stor_Emp_2018)
print('Pearsons correlation: %.3f' % corrP)

# Spearman's correlation
corrS, _ = spearmanr(stor_Fac_2018, stor_Emp_2018)
print('Spearman's correlation: %.3f' % corrS)

```

```

2006
[[ 5.5877451  41.05294118]
 [ 41.05294118 949.53899782]]
Pearsons correlation: 0.564
Spearman's correlation: 0.496
2013
[[ 5.9877451  23.13954248]
 [ 23.13954248 323.34618736]]
Pearsons correlation: 0.526
Spearman's correlation: 0.520
2018
[[ 6.65882353  21.59738562]
 [ 21.59738562 238.70653595]]
Pearsons correlation: 0.542
Spearman's correlation: 0.485

```

1.6.3 Linear regression

regressor - # Facilities ; predictor - employee count

1.6.4 Create model

```
[334]: stor_Fac_2006 = stor_Fac_2006.reshape((-1, 1))
stor_Fac_2013 = stor_Fac_2013.reshape((-1, 1))
stor_Fac_2018 = stor_Fac_2018.reshape((-1, 1))
```

```
[335]: model_2006 = LinearRegression().fit(stor_Fac_2006, stor_Emp_2006)
model_2013 = LinearRegression().fit(stor_Fac_2013, stor_Emp_2013)
model_2018 = LinearRegression().fit(stor_Fac_2018, stor_Emp_2018)
```

```
[336]: print("2006")
r_sq = model_2006.score(stor_Fac_2006, stor_Emp_2006)
print('coefficient of determination:', r_sq)
print('intercept:', model_2006.intercept_)
print('slope:', model_2006.coef_)

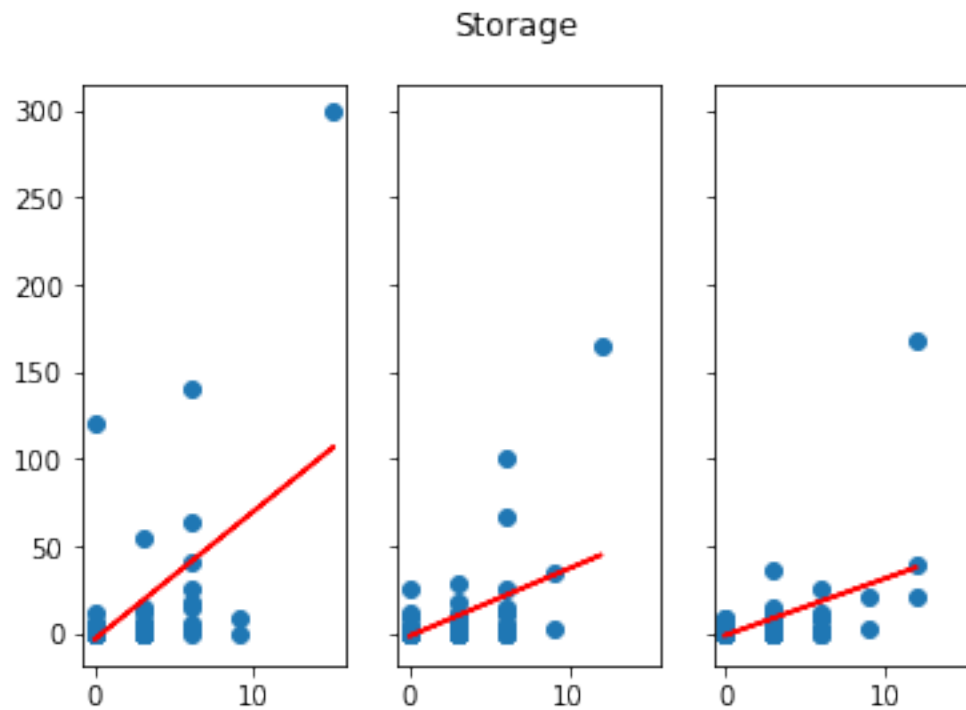
print("2013")
r_sq = model_2013.score(stor_Fac_2013, stor_Emp_2013)
print('coefficient of determination:', r_sq)
print('intercept:', model_2013.intercept_)
print('slope:', model_2013.coef_)

print("2018")
r_sq = model_2018.score(stor_Fac_2018, stor_Emp_2018)
print('coefficient of determination:', r_sq)
print('intercept:', model_2018.intercept_)
print('slope:', model_2018.coef_)
```

```
2006
coefficient of determination: 0.31764290677708873
intercept: -3.53364330204404
slope: [7.34696026]
2013
coefficient of determination: 0.2765530721347862
intercept: -1.6788374948833429
slope: [3.86448356]
2018
coefficient of determination: 0.2934543485079212
intercept: -1.3719081272084797
slope: [3.24342364]
```

```
[337]: fig, (ax1, ax2, ax3) = plt.subplots(1, 3, sharex=True, sharey=True)
fig.suptitle('Storage')
ax1.scatter(stor_Fac_2006, stor_Emp_2006)
ax1.plot(stor_Fac_2006, model_2006.coef_*stor_Fac_2006+model_2006.intercept_, 'r')
ax2.scatter(stor_Fac_2013, stor_Emp_2013)
ax2.plot(stor_Fac_2013, model_2013.coef_*stor_Fac_2013+model_2013.intercept_, 'r')
ax3.scatter(stor_Fac_2018, stor_Emp_2018)
ax3.plot(stor_Fac_2018, model_2018.coef_*stor_Fac_2018+model_2018.intercept_, 'r')
```

[337]: [<matplotlib.lines.Line2D at 0x122b24af0>]



[]:

[]: