

---

Python

# Pandas/상관관계

---

김선녕(sykim.lecture@gmail.com)

참고문헌 : 파이썬을 이용한 빅데이터 수집, 분석과 시각화 - 비팬북스, 이원화

---

# Pandas Library

상관관계

- Powerful Python data analysis toolkit
- BSD-licensed library providing high-performance, easy-to-use data structures and data analysis tools for the Python programming language.
- The two primary data structures of pandas, Series (1-dimensional) and DataFrame (2-dimensional), handle the vast majority of typical use cases in finance, statistics, social science, and many areas of engineering.
- `pip install pandas`
- Ref : <https://pandas.pydata.org>

## Series(1-dimensional)

- 인덱스(레이블)를 가지는 동일한 데이터형의 1차원 데이터
- 인덱스는 중복 가능
- 레이블 또는 데이터의 위치를 지정한 추출가능.
- 인덱스에 대한 슬라이스가 가능
- 산술 연산이 가능. 통계량을 산출하는 장점을 가지고 있음

Index	Data
1	'A'
2	'B'
3	'C'
4	'D'
5	'E'

```
import numpy as np
import pandas as pd
```

```
s = pd.Series(np.random.randn(5))
s
```

```
0    -0.086872
1     0.260547
2     0.012375
3    -1.185436
4    -0.985884
dtype: float64
```

```
s = pd.Series(np.random.randn(5), index=['A', 'B', 'C', 'D', 'E'])
s
```

```
A     1.172776
B     0.295672
C    -1.450593
D     0.394875
E     1.411907
dtype: float64
```

```
# dictionary  
d = {'a' : 0., 'b' : 1., 'c' : 2.}  
pd.Series(d)
```

```
a    0.0  
b    1.0  
c    2.0  
dtype: float64
```

```
pd.Series(d, index=['a', 'b', 'B', 'c'])
```

```
a    0.0  
b    1.0  
B    NaN  
c    2.0  
dtype: float64
```

```
# 스칼라값  
pd.Series(7, index=['a', 'b', 'c', 'd', 'e'])
```

```
a    7  
b    7  
c    7  
d    7  
e    7  
dtype: int64
```

```
s = pd.Series([1,2,3,4,5], index=['a', 'b', 'c', 'd', 'e'])  
s[0]
```

1

```
s[:3]
```

```
a    1  
b    2  
c    3  
dtype: int64
```

```
s[[4,1]]
```

```
e    5  
b    2  
dtype: int64
```

```
np.power(s, 2)
```

```
a    1  
b    4  
c    9  
d   16  
e   25  
dtype: int64
```

## DataFrame(2-dimensional)

- 행과 열에 레이블을 가진 2차원 데이터
- 열마다 다른 형태를 가질 수 있음
- 테이블형 데이터에 대해 불러오기, 데이터 쓰기가 가능
- DataFrame끼리 여러 가지 조건을 사용한 결합 처리가 가능
- 크로스 집계 가능

The diagram illustrates a DataFrame as a table. The word 'Columns' is positioned above the table with two arrows pointing to the header row. The word 'ROWS' is positioned to the left of the table with three arrows pointing to the first three data rows.

Regd. No	Name	Marks%
1000	Steve	86.29
1001	Mathew	91.63
1002	Jose	72.90
1003	Patty	69.23
1004	Vin	88.30



# Series/Dict 데이터의 활용

```
d = {'one' : pd.Series([1., 2., 3.], index=['a', 'b', 'c']),  
     'two' : pd.Series([1., 2., 3., 4.], index=['a', 'b', 'c', 'd'])}
```

d

```
{'one': a    1.0  
      b    2.0  
      c    3.0  
      dtype: float64, 'two': a    1.0  
      b    2.0  
      c    3.0  
      d    4.0  
      dtype: float64}
```

```
pd.DataFrame(d)
```

	one	two
a	1.0	1.0
b	2.0	2.0
c	3.0	3.0
d	NaN	4.0

# 인덱스 값을 부여하지 않으면 자동으로 0부터 두개의 데이터 중 큰 배열의 길이 - 1 만큼이 부여

```
d = {'one' : pd.Series([1., 2., 3.]), 'two' : pd.Series([1., 2., 3., 4.])}
```

```
pd.DataFrame(d)
```

	one	two
0	1.0	1.0
1	2.0	2.0
2	3.0	3.0
3	NaN	4.0

```
d = {'one' : [1., 2., 3., 4.], 'two' : [4., 3., 2., 1.]}  
pd.DataFrame(d)
```

	one	two
0	1.0	4.0
1	2.0	3.0
2	3.0	2.0
3	4.0	1.0

```
# Dict 리스트 데이터의 활용  
data2 = [{'a': 1, 'b': 2}, {'a': 5, 'b': 10, 'c': 20}]  
pd.DataFrame(data2)
```

	a	b	c
0	1	2	NaN
1	5	10	20.0

```
pd.DataFrame(data2, index=['first', 'second'])
```

	a	b	c
first	1	2	NaN
second	5	10	20.0

```
pd.DataFrame(data2, columns=['a', 'b'])
```

	a	b
0	1	2
1	5	10

```
df = pd.DataFrame(data2, columns=['a', 'b'])  
df.rename(columns={'a': 'COL1'})
```

	COL1	b
0	1	2
1	5	10

```
df.set_index('b')
```

a	
b	
2	1
10	5

```
# 데이터 추가 및 합치기(merge)
data1 = [{ 'name': 'Mark' }, { 'name': 'Eric' }, { 'name': 'Jennifer' }]
df = pd.DataFrame(data1)
df
```

	name
0	Mark
1	Eric
2	Jennifer

```
df['age'] = [10, 11, 12]
pd.DataFrame(data1)
df
```

	name	age
0	Mark	10
1	Eric	11
2	Jennifer	12

```
data2 = [{'sido': '서울'}, {'sido': '경기'}, {'sido': '인천'}]  
df2 = pd.DataFrame(data2)  
df2
```

**sido**

0 서울

1 경기

2 인천

```
pd.merge(df, df2, left_index=True, right_index=True)
```

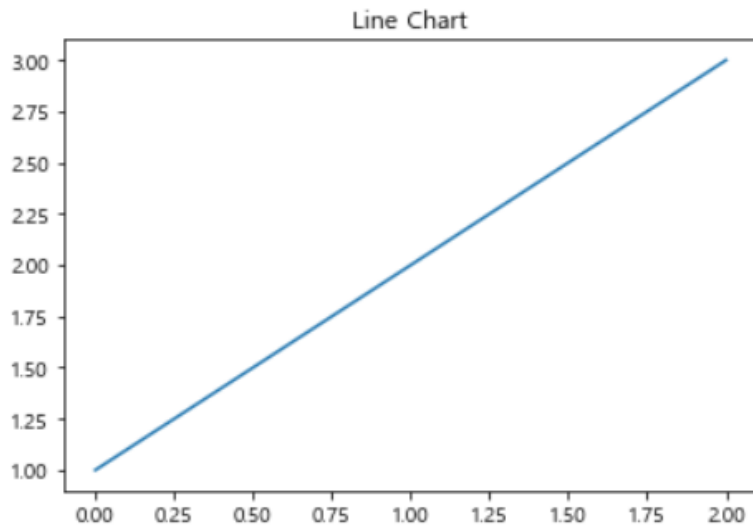
**name age sido**

0 Mark 10 서울

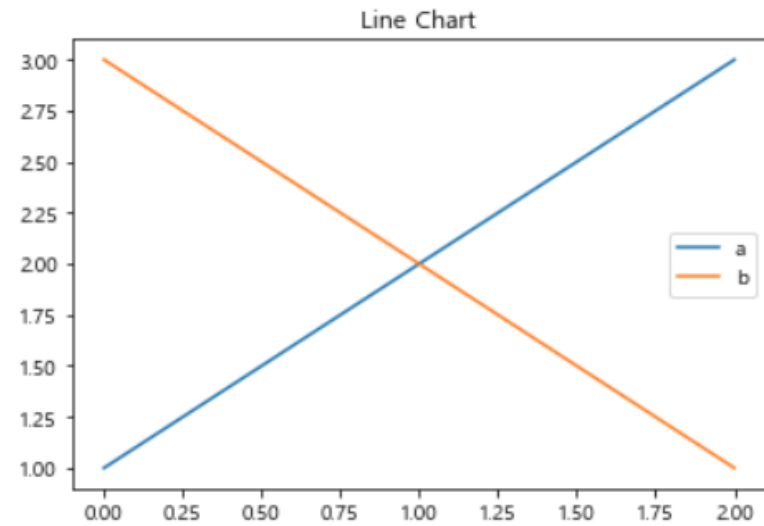
1 Eric 11 경기

2 Jennifer 12 인천

```
# Series에서 시각화  
s = pd.Series([1,2,3])  
ax = s.plot()  
ax.set_title('Line Chart')  
plt.show()
```



```
# DataFrame에서 시각화  
df = pd.DataFrame({'a':[1,2,3], 'b':[3,2,1]})  
ax = df.plot()  
ax.set_title('Line Chart')  
plt.show()
```



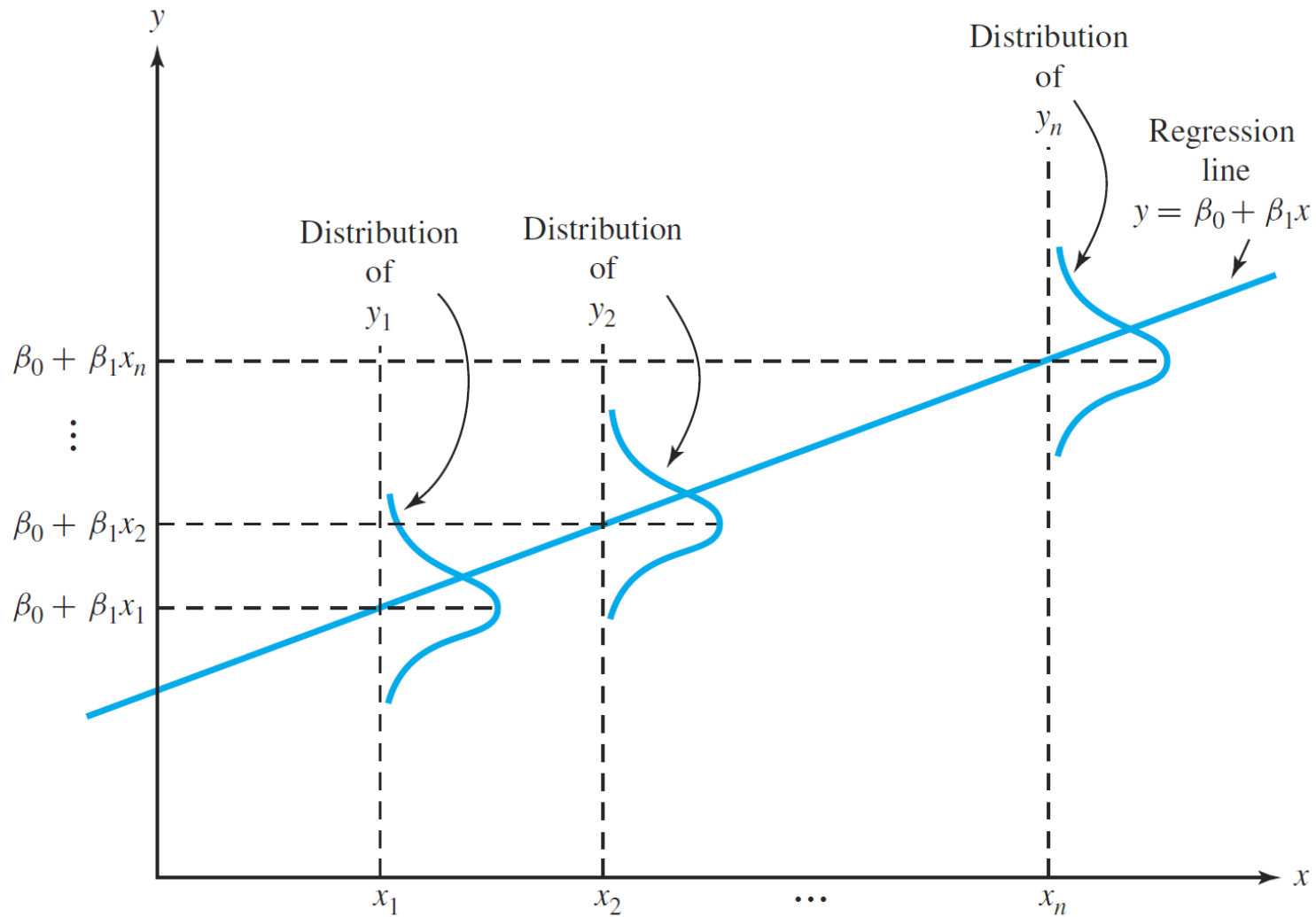
함수	설명
<code>rand(d0, d1, ..., dn)</code>	N차원 배열의 난수 발생
<code>randn(d0, d1, ..., dn)</code>	표준 정규분포에 따른 N차원 난수 발생
<code>randint(low[, high, size])</code>	low 이상 high 미만의 정수형 난수 발생
<code>random_sample([size])</code>	0.0이상 1.0미만의 실수형 난수 발생
<code>random([size])</code>	
<code>ranf([size])</code>	
<code>sample([size])</code>	
<code>choice(a[, size, replace, p])</code>	주어진 1차원 배열을 기반으로 무작위 샘플 추출
<code>bytes(length)</code>	바이트형 난수 발생

Pandas Library

**상관관계**



# 단순선형회귀 (Simple Linear Regression)

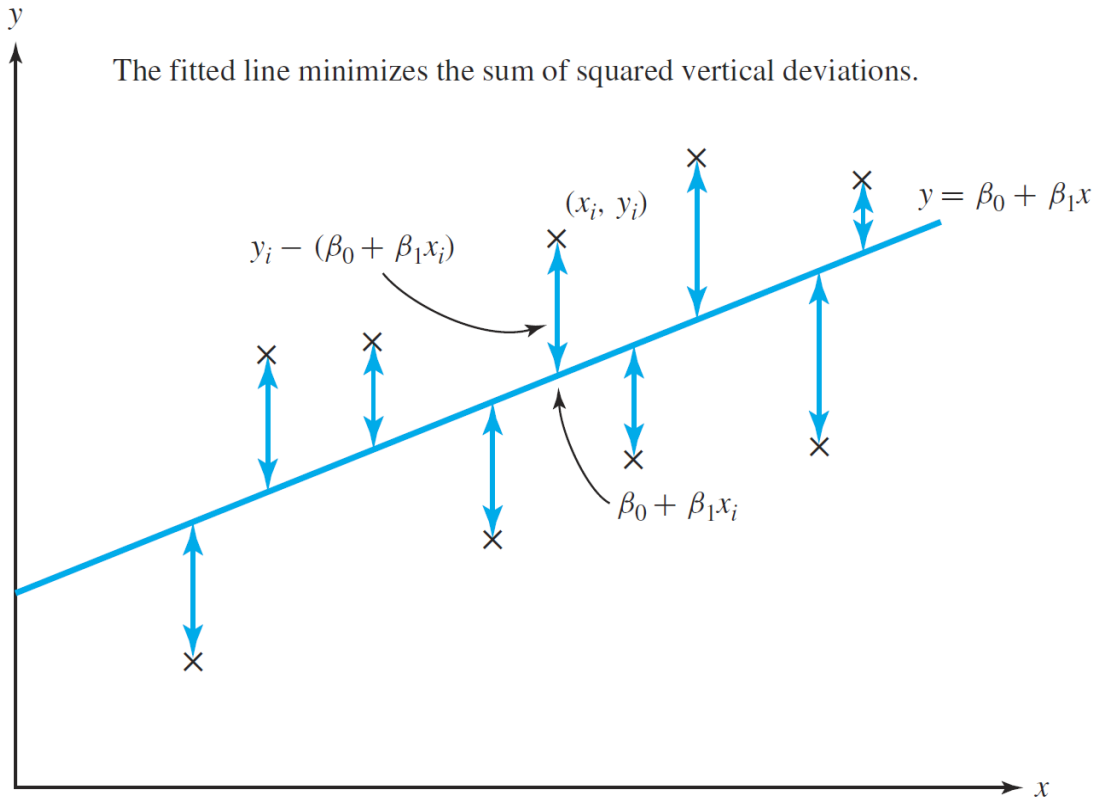


# 최소제곱법(Least Squared Method)

- 회귀직선  $y = \beta_0 + \beta_1 x$  를 데이터  $(x_1, y_1), \dots, (x_n, y_n)$ 에 적합(fit)시키는 과정
  - 데이터들에 가장 근접한 직선을 찾는 과정

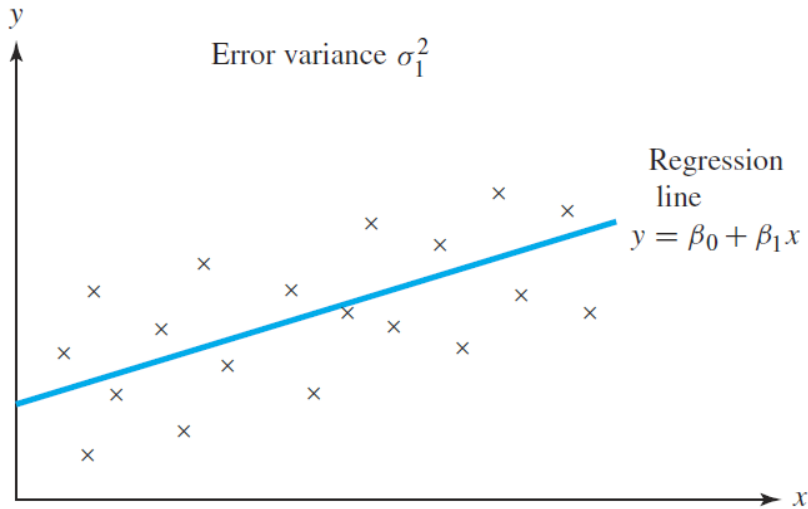
$$Q = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2$$

을 최소화하는 직선을 선택하는 것.

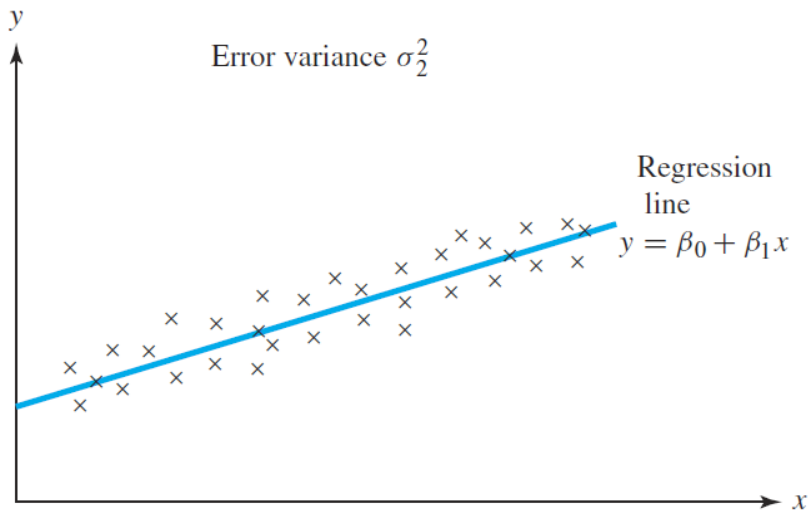


# 오차분산(Error Variance)

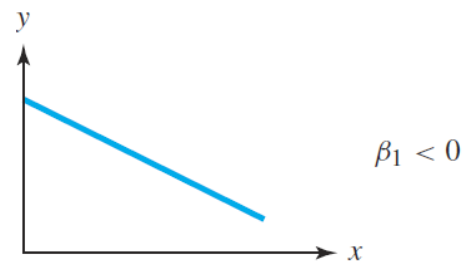
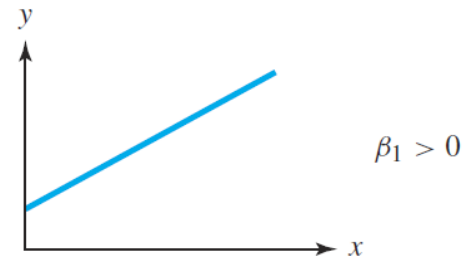
19



$$\sigma_1^2 > \sigma_2^2$$



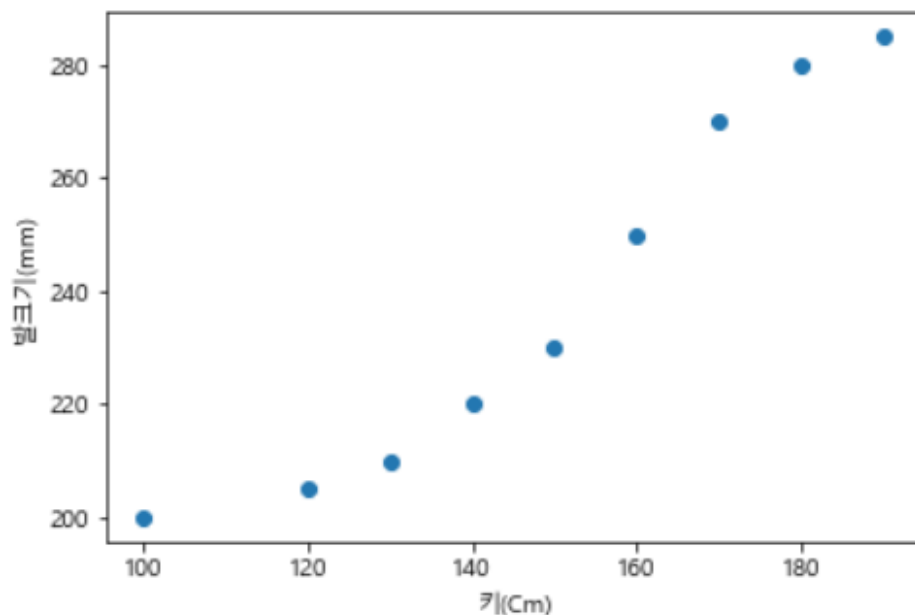
Interpretation of the error variance  $\sigma^2$



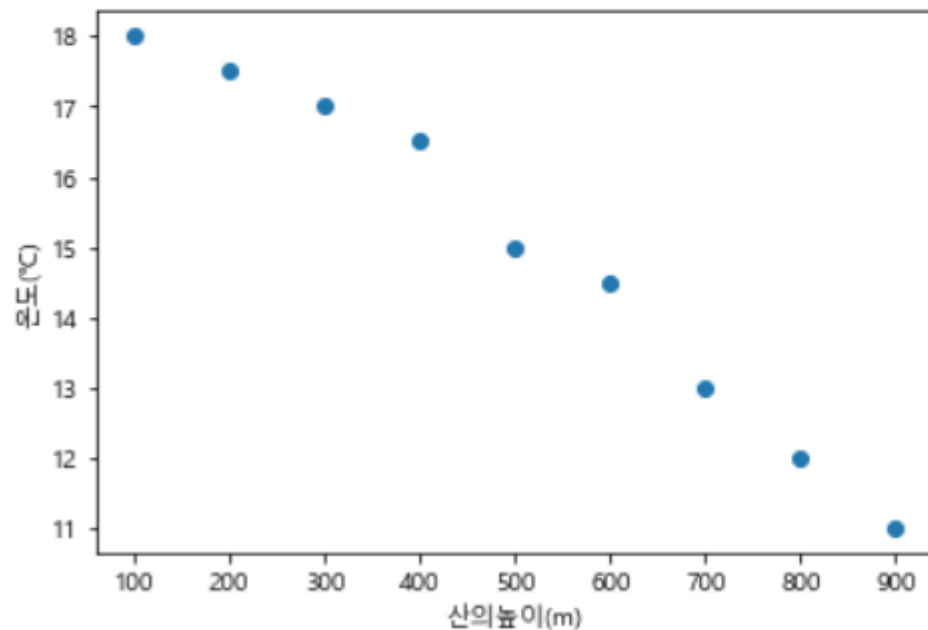
Interpretation of slope parameter  $\beta_1$

```
import matplotlib.pyplot as plt
from matplotlib import font_manager, rc
import matplotlib
font_location = "c:/Windows/fonts/malgun.ttf"
font_name = font_manager.FontProperties(fname=font_location).get_name()

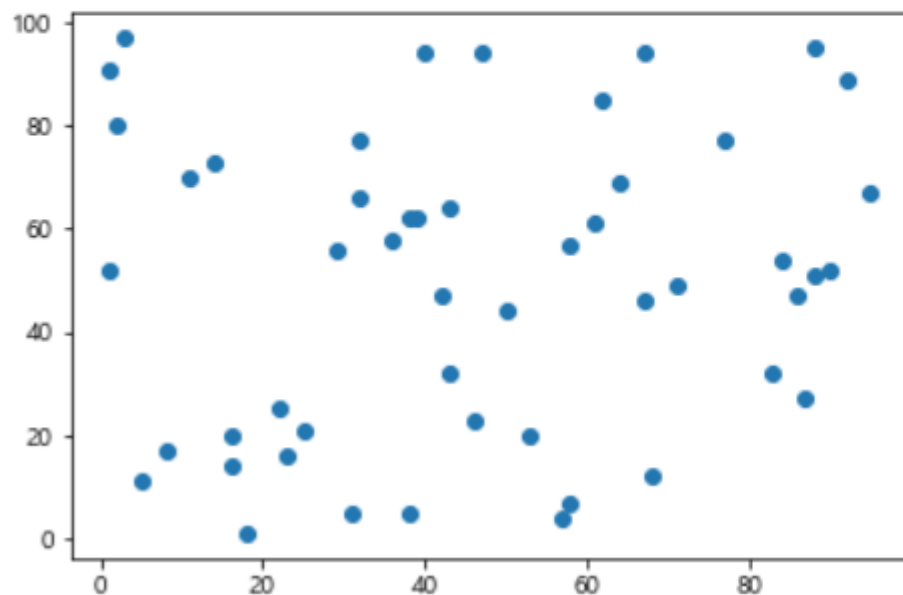
matplotlib.rc('font', family=font_name)
height = [100, 120, 130, 140, 150, 160, 170, 180, 190]
foot_size = [200, 205, 210, 220, 230, 250, 270, 280, 285]
plt.scatter(height, foot_size)
plt.xlabel('키(Cm)')
plt.ylabel('발크기(mm)')
plt.show()
```



```
height = [100, 200, 300, 400, 500, 600, 700, 800, 900]  
temperature = [18.0, 17.5, 17, 16.5, 15, 14.5, 13, 12, 11]  
plt.scatter(height, temperature)  
plt.xlabel('산의높이(m)')  
plt.ylabel('온도(°C)')  
plt.show()
```



```
import numpy as np
random_x = np.random.randint(0, 100, 50) # 1~100 사이의 난수를 50개 생성
random_y = np.random.randint(0, 100, 50) # 1~100 사이의 난수를 50개 생성
plt.scatter(random_x, random_y)
plt.show()
```



0.0 ~ 0.2	상관 관계가 거의 없다
0.2 ~ 0.4	약한 상관 관계
0.4 ~ 0.6	상관 관계가 있다
0.6 ~ 0.8	강한 상관 관계
0.8 ~ 1.0	매우 강한 상관 관계

- $m_x$  is the mean of X
- $m_y$  is the mean of Y
- $r$  is the correlation between X and Y(Pearson's r) :

$$r = \frac{\sum xy}{\sqrt{\sum x^2 \sum y^2}}, \quad (x = X - m_x, y = Y - m_y)$$

```
1 import json
2 import math
3 import numpy as np
4
5 import matplotlib.pyplot as plt
6 import matplotlib
7 from matplotlib import font_manager, rc
8
9 import pandas as pd
10
11 #[CODE 1]
12
13 def correlation(x, y):
14     n = len(x)
15     vals = range(n)
16
17     x_sum = 0.0
18     y_sum = 0.0
19     x_sum_pow = 0.0
20     y_sum_pow = 0.0
21     mul_xy_sum = 0.0
22
23     for i in vals:
24         mul_xy_sum = mul_xy_sum + float(x[i]) * float(y[i])
25         x_sum = x_sum + float(x[i])
26         y_sum = y_sum + float(y[i])
27         x_sum_pow = x_sum_pow + pow(float(x[i]), 2)
28         y_sum_pow = y_sum_pow + pow(float(y[i]), 2)
29
30     try:
31         r = ((n * mul_xy_sum) - (x_sum * y_sum)) / math.sqrt( ((n*x_sum_pow) - pow(x_sum, 2)) * ((n*y_sum_pow) - pow(y_sum, 2)) )
32     except:
33         r = 0.0
34
35     return r
36
```



```

37 #[CODE 2]
38
39 def setScatterGraph(tour_table, visit_table, tourpoint):
40
41     tour = tour_table[tour_table['resNm'] == tourpoint]
42     merge_table = pd.merge(tour, visit_table, left_index=True, right_index=True)
43
44     fig = plt.figure()
45
46     fig.suptitle(tourpoint + '상관관계 분석')
47
48     plt.subplot(1, 3, 1)
49     plt.xlabel('중국인 입국수')
50     plt.ylabel('외국인 입장객수')
51     r = correlation(list(merge_table['china']), list(merge_table['ForNum']))
52     plt.title('r = {:.5f}'.format(r))
53     plt.scatter(list(merge_table['china']), list(merge_table['ForNum']), edgecolor='none', alpha=0.75, s=6, c='black')
54
55     plt.subplot(1, 3, 2)
56     plt.xlabel('일본인 입국수')
57     plt.ylabel('외국인 입장객수')
58     r = correlation(list(merge_table['japan']), list(merge_table['ForNum']))
59     plt.title('r = {:.5f}'.format(r))
60     plt.scatter(list(merge_table['japan']), list(merge_table['ForNum']), edgecolor='none', alpha=0.75, s=6, c='black')
61
62     plt.subplot(1, 3, 3)
63     plt.xlabel('미국인 입국수')
64     plt.ylabel('외국인 입장객수')
65     r = correlation(list(merge_table['usa']), list(merge_table['ForNum']))
66     plt.title('r = {:.5f}'.format(r))
67     plt.scatter(list(merge_table['usa']), list(merge_table['ForNum']), edgecolor='none', alpha=0.75, s=6, c='black')
68
69     plt.tight_layout()
70
71     # 이미지 저장
72     #fig = matplotlib.pyplot.gcf()
73
74     #fig.set_size_inches(10, 7)
75
76     #fig.savefig(tourpoint+'.png', dpi=300)
77
78     plt.show()

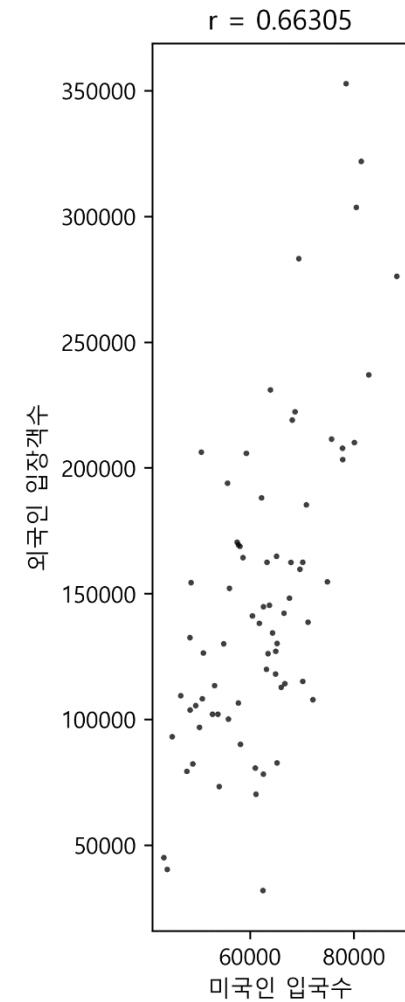
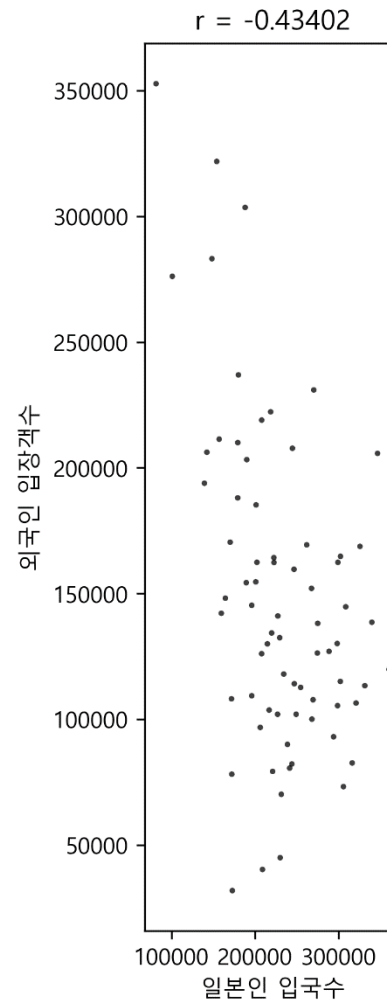
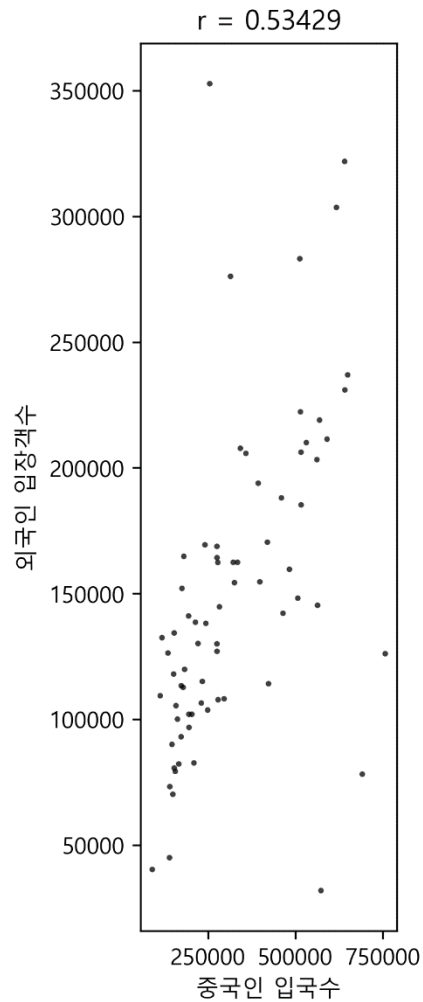
```

```

79 def main():
80
81     font_location = "c:/Windows/fonts/malgun.ttf"
82     font_name = font_manager.FontProperties(fname=font_location).get_name()
83     matplotlib.rc('font', family=font_name)
84
85     tpFileName = '서울특별시_관광지입장정보_2011_2016.json'
86     jsonTP = json.loads(open(tpFileName, 'r', encoding='utf-8').read())
87     tour_table = pd.DataFrame(jsonTP, columns=('yyymm', 'resNm', 'ForNum')) # 필요한 데이터만 추출
88     tour_table = tour_table.set_index('yyymm')
89
90     ''' ['창덕궁', '운현궁', '경복궁', '창경궁', '종묘', '국립중앙박물관', '서울역사박물관', '덕수궁',
91         '서울시립미술관 본관', '태릉·강릉·조선왕릉전시관', '서대문형무소역사관', '서대문자연사박물관',
92         '트릭아이미술관', '현릉·인릉', '선릉·정릉', '롯데월드'] '''
93
94     resNm = tour_table.resNm.unique()
95
96     fv_CFileName = '중국(112)_해외방문객정보_2011_2016.json'
97     jsonFV = json.loads(open(fv_CFileName, 'r', encoding='utf-8').read())
98     china_table = pd.DataFrame(jsonFV, columns=('yyymm', 'visit_cnt'))
99     china_table = china_table.rename(columns={'visit_cnt': 'china'})
100     china_table = china_table.set_index('yyymm')
101
102     fv_JFileName = '일본(130)_해외방문객정보_2011_2016.json'
103     jsonFV = json.loads(open(fv_JFileName, 'r', encoding='utf-8').read())
104     japan_table = pd.DataFrame(jsonFV, columns=('yyymm', 'visit_cnt'))
105     japan_table = japan_table.rename(columns={'visit_cnt': 'japan'})
106     japan_table = japan_table.set_index('yyymm')
107
108     fv_UFileName = '미국(275)_해외방문객정보_2011_2016.json'
109     jsonFV = json.loads(open(fv_UFileName, 'r', encoding='utf-8').read())
110     usa_table = pd.DataFrame(jsonFV, columns=('yyymm', 'visit_cnt'))
111     usa_table = usa_table.rename(columns={'visit_cnt': 'usa'})
112     usa_table = usa_table.set_index('yyymm')
113
114     fv_table = pd.merge(china_table, japan_table, left_index=True, right_index=True)
115     fv_table = pd.merge(fv_table, usa_table, left_index=True, right_index=True)
116
117     for tourpoint in resNm:
118         setScatterGraph(tour_table, fv_table, tourpoint)
119
120 if __name__ == "__main__":
121     main()
122 
```

# 서울시 입장객과 각국 입국수에 따른 상관분석 자료

경북공상관관계 분석



```

1 import json
2 import math
3 import numpy as np
4
5 import matplotlib.pyplot as plt
6 import matplotlib
7 from matplotlib import font_manager, rc
8 import pandas as pd
9
10 '''[['창덕궁', -0.058791104060063125, 0.27744435701410114, 0.40281606330501574], ['운현궁', 0.44594488384450376, 0.3026152182879861,
11 0.2812576500158649], ['경복궁', 0.5256734293511214, -0.4352281861341233, 0.42513726387044926], ['창경궁', 0.4512325398089607,
12 -0.16458589402253013, 0.6245403780269381], ['종묘', -0.5834218986767474, 0.5298702802205213, -0.12112666829294959], ['국립중앙박물관'
13 , 0.39663594900292837, -0.06923889417914424, 0.3789788348060077], ['서울역사박물관', 0.4169985898874495, 0.4929777868070643,
14 0.2411976107709704], ['덕수궁', 0.4332132943587757, -0.4326719125679966, 0.4808588469548069], ['서울시립미술관 본관', 0.0, 0.0, 0.0],
15 ['태릉 · 강릉 · 조선왕릉전시관', -0.08179909096513825, 0.0634032985209752, -0.068840
16 0.47262271531670347, 0.006098570233700235, 0.22900879409607508], ['서대문자연사박물관',
17 0.340084882575556, -0.15036015533747007, 0.18094502388483083], ['현릉 · 인릉', -0.581325
18 -0.1853887818740637], ['선릉 · 정릉', -0.5715258789199192, 0.38806730592260075, -0.12494
19 0.23511773800458452, -0.12673869767365747]]'''
20
21 f = open("rlist.txt", "r")
22 r_table = f
23 f.close()
24
25
26 r_table = pd.DataFrame(r_list, columns=('tourpoint', 'china', 'japan', 'usa'))
27 r_table = r_table.set_index('tourpoint')
28 r_table

```

	china	japan	usa
tourpoint			
창덕궁	-0.058791	0.277444	0.402816
운현궁	0.445945	0.302615	0.281258
경복궁	0.525673	-0.435228	0.425137
창경궁	0.451233	-0.164586	0.624540
종묘	-0.583422	0.529870	-0.121127
국립중앙박물관	0.396636	-0.069239	0.378979
서울역사박물관	0.416999	0.492978	0.241198
덕수궁	0.433213	-0.432672	0.480859
서울시립미술관 본관	0.000000	0.000000	0.000000
태릉 · 강릉 · 조선왕릉전시관	-0.081799	0.063403	-0.068840
서대문형무소역사관	0.472623	0.006099	0.229009
서대문자연사박물관	0.000000	0.000000	0.000000
트릭아미미술관	0.340085	-0.150360	0.180945
현릉 · 인릉	-0.581325	0.464530	-0.185389
선릉·정릉	-0.571526	0.388067	-0.124945
롯데월드	0.510559	0.235118	-0.126739

```
1 # 상관계수가 없는 경우에는 삭제
2 r_table.drop('서울시립미술관 본관')
3 r_table.drop('서대문자연사박물관')
4 r_table = r_table.sort_values('china', ascending=False)
5 r_table.head() # default 가 5개
```

	china	japan	usa
tourpoint			
경복궁	0.525673	-0.435228	0.425137
롯데월드	0.510559	0.235118	-0.126739
서대문형무소역사관	0.472623	0.006099	0.229009
창경궁	0.451233	-0.164586	0.624540
운현궁	0.445945	0.302615	0.281258

```

1 # 중국인 입국수 대비 관광객 입장객수의 상관계수가 높은 순서대로 3개국의 비교
2 font_location = "c:/Windows/fonts/malgun.ttf"
3 font_name = font_manager.FontProperties(fname=font_location).get_name()
4 matplotlib.rc('font', family=font_name)
5 r_table.plot(kind='bar', rot=70)
6
7 plt.show()

```

