

7.5 공분산과 상관계수

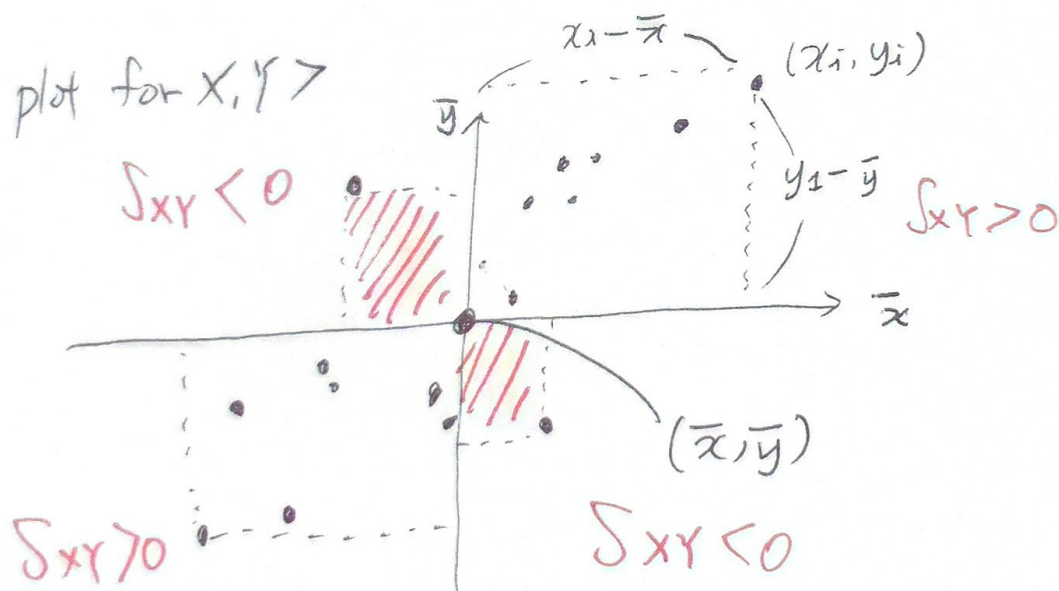
다변수 확률변수 간의 상관관계를 숫자로 나타낸 것이
공분산(covariance)과 상관계수(correlation coefficient)다.

표본 공분산 sample covariance

$$S_{xy} = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})$$

$X \longleftrightarrow Y$, 얼마나 독립인가? 얼마나 상관성이 존재하는가?

< Scatter plot for X, Y >



- | | | | |
|--|---|---|---------|
| | → | • 1, 3사분면 多, 2, 4사분면 少 → $S_{xy} > 0$ | } 모양 정보 |
| | → | • 2, 4사분면 多, 1, 3사분면 少 → $S_{xy} < 0$ | |
| | → | • 2, 4사분면 \cong 1, 3사분면 data → $S_{xy} \cong 0$ | |

→ $|S_{xy}|$ 작음 } 퍼져있는 정도에 대한 정보 (분산 크기)
→ $|S_{xy}|$ 큼

표본상관계수

표본공분산은 평면을 중심으로 각 자료들이 어떻게 분포되어 있는지 크기와 방향성을 같이 보여준다. 그런데 분포의 크기는 공분산이 아닌 분산만으로도 알 수 있기 때문에 대부분의 경우 자료 분포의 방향성만 보아도 되는 것이 유용하다.

이때 필요한 것이 표본상관계수 (sample correlation coefficient) 다. 표본상관계수는 아래와 같이 공분산을 각각의 표본표준편차값으로 나누어 정규화 (normalize) 하여 정의한다.

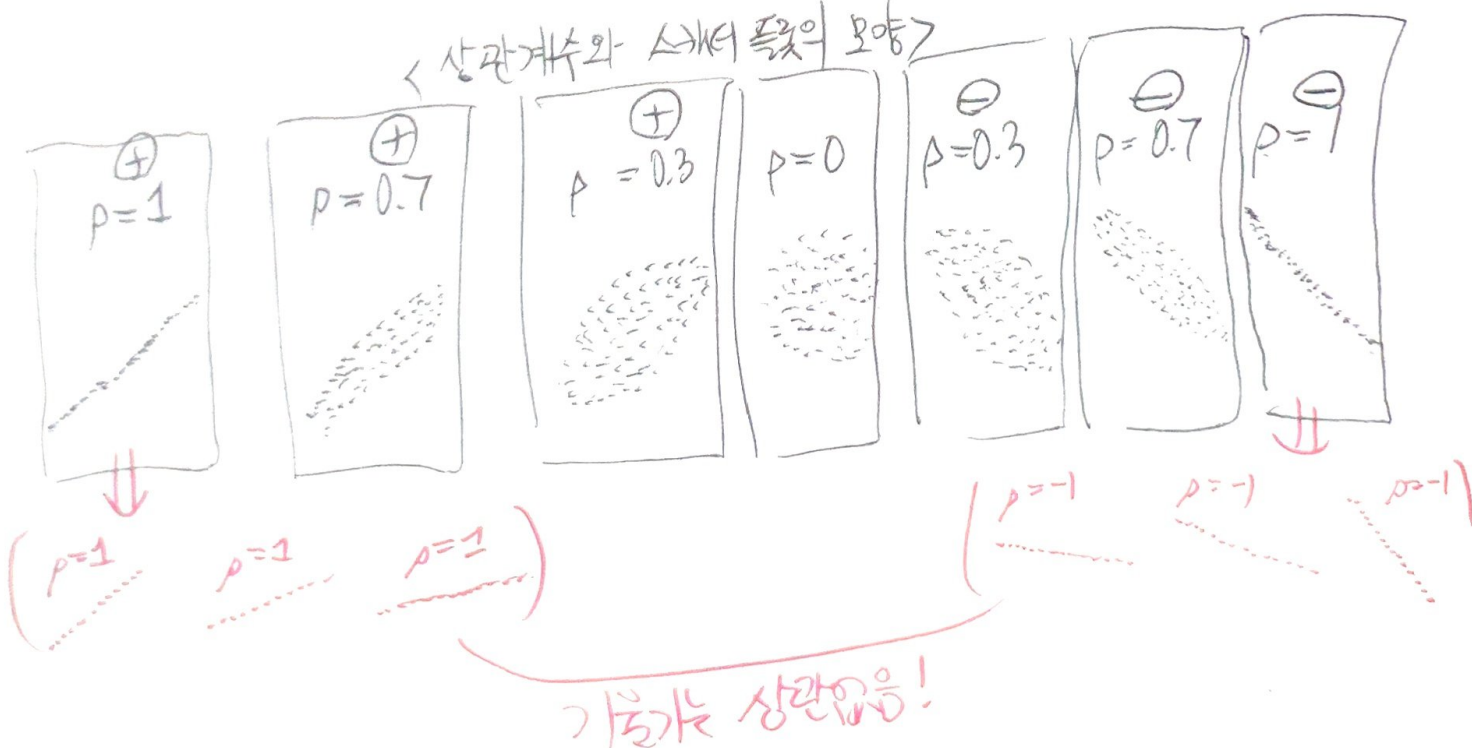
$$\text{Pearson } r = r_{xy} = \frac{\overset{\text{공분산}}{S_{xy}}}{\sqrt{\underset{\text{X분산}}{S_x^2} \cdot \underset{\text{Y분산}}{S_y^2}}} = \frac{\overset{\text{공분산}}{S_{xy}}}{\sqrt{\Pi \text{ 표본표준편차}^2}}$$

이와 다르게 정의한 상관계수도 있기 때문에 다른 종류의 상관계수와 비교하여 말하는 경우 피어슨(Pearson) 상관계수라고 하기도 한다.

사이파이 stats 서브패키지는 피어슨 상관계수를 계산하는 `pearsonr()` 함수를 제공한다. `pearsonr()` 함수는 상관계수와 유의확률을 반환한다. 유의확률에 대해서는 아래에서 공부한다.

`scipy.stats.pearsonr(x1, x2)[0] = r`
`[1] = 유의p`

< 상관계수와 스캐터 플롯의 모양 >



확률변수의 공분산과 상관계수

두 확률변수 X 와 Y 의 공분산은 기대값 연산자를 사용하여 다음과 같이 정의된다.

$$Cov[X, Y] = E[(X - E[X])(Y - E[Y])]$$

마찬가지로 두 확률변수 X 와 Y 의 상관계수도 다음과 같이 정의한다.

$$\rho[X, Y] = \frac{Cov[X, Y]}{\sqrt{Var[X] \cdot Var[Y]}}$$

↑ 이분점 확률, ho

확률변수의 상관계수는 다음과 같은 성질을 가진다.

$$-1 \leq \rho \leq 1$$

또한 특수한 ρ 에 대하여 다음과 같이 부른다.

- $\rho = 1$) 완전선형 상관관계
- $\rho = 0$) 무상관 (\neq 독립)
- $\rho = -1$) 완전선형 반상관관계

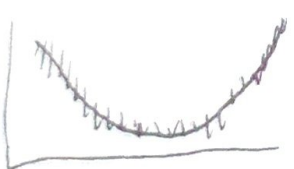
독립 $\Rightarrow \rho = 0$
 $\rho = 0 \nRightarrow$ 독립

비선형 상관관계

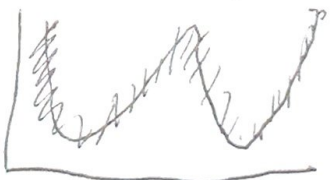
두 확률변수가 상관관계가 있으면, 두 확률변수의 값 중 하나를 알았을 때 다른 확률변수의 값에 대한 정보를 알 수 있다. < 반대로 정확한 값을 알 수 있어야 하는 것은 아니다. >
hint를 얻을 수 있다는 의미.

* 피어슨 상관계수는 두 확률변수의 관계가 선형적일 때만 상관관계를 제대로 계산할 수 있다.

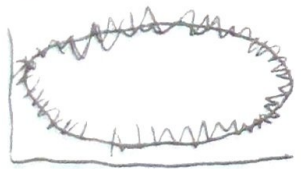
(비선형 상관관계 1)
 $r = 0.0$



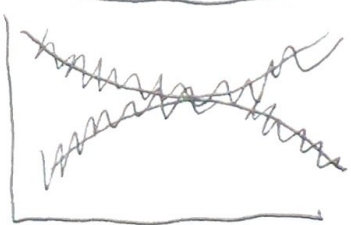
(비선형 상관관계 2)
 $r = 0.0$



(비선형 상관관계 3)
 $r = 0.0$



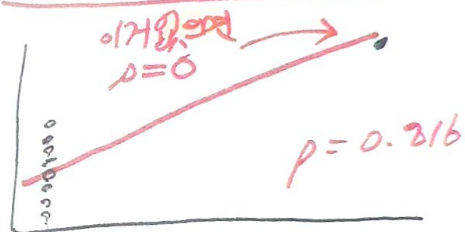
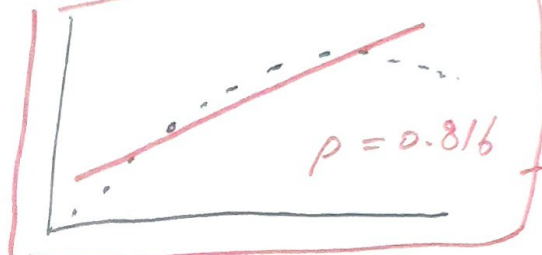
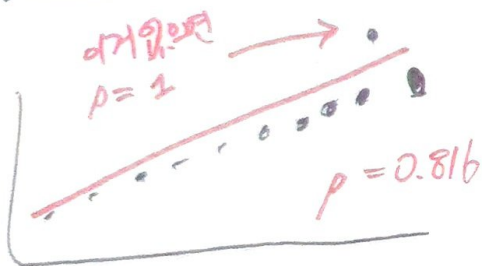
(비선형 상관관계 4)
 $r = 0.0$



앤스콤 데이터

상관계수로 분포 모양을 추측할 때, 개별 자료-상관계수에 미치는 영향력이 각각 다른 점에 유의해야 한다. 다음은 프랭크 앤스콤 (Frank Anscombe) 의 문제에 제시된 데이터다. 이 데이터는 서로 다른 4종류의 2차원 데이터셋을 포함하는데, 4종류 데이터셋 모두 상관계수가 약 0.816으로 동일하다.

(Anscombe Data)



매우 개량한 (x,y) 2차-상관계수로 시뮬레이션함
But $\rho = 0.816$

Conclusion

pearson 상관계수는 2차식 등 비선형 상관관계를 catch하지 못하여 outlier이 민감/취약하다.

다변수 확률변수의 표본공분산

- 스칼라가 아닌 벡터 형태를 가지는 다변수 확률변수의 공분산
- X_1, X_2, \dots, X_M 이라는 서로 다른 M 개의 확률변수
- \hookrightarrow 확률변수마다 각각 N 개의 표본
 - $\hookrightarrow j$ ($j = 1, \dots, M$) 번째 확률변수의 i ($i = 1, \dots, N$) 번째 데이터
 $= x_{i,j}$

"하나로 묶으면 다음과 같은 특징행렬이 됨."

$$X = \begin{bmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,M} \\ x_{2,1} & x_{2,2} & \dots & x_{2,M} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N,1} & x_{N,2} & \dots & x_{N,M} \end{bmatrix}$$

M 개의 서로 다른 확률변수의 모든 조합에 대한 공분산을 한꺼번에 표기하기 위해,
 다음처럼 **표본 공분산 행렬 (Sample Covariance Matrix)**을 정의한다.

- 대각성분 = 각 확률변수의 분산으로 정의
- 비대각성분 = **서로 다른 두 확률변수의 공분산**으로 정의

오개씩 계산한다.

$$S = \begin{bmatrix} S_{x_1}^2 & S_{x_1 x_2} & \dots & S_{x_1 x_M} \\ S_{x_1 x_2} & S_{x_2}^2 & \dots & S_{x_2 x_M} \\ \vdots & \vdots & \ddots & \vdots \\ S_{x_1 x_M} & S_{x_2 x_M} & \dots & S_{x_M}^2 \end{bmatrix}$$

위 행렬의 값은 다음처럼 구한다.

(1) 각 확률변수 x_j ($j=1, \dots, M$)의 표본평균을 계산한다.

$$\bar{x}_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

(2) 각 확률변수 x_j ($j=1, \dots, M$)의 분산을 계산한다.

$$s_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \bar{x}_j)^2$$

(3) 두 확률변수 x_j, x_k 의 공분산을 계산한다.

$$s_{j,k} = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \bar{x}_j)(x_{i,k} - \bar{x}_k)$$

만약 x_i ($i=1, \dots, N$)가 다음과 같은 M 차원 표본 벡터를 정의하면

$$x_i = \begin{bmatrix} x_{i,1} \\ \vdots \\ x_{i,M} \end{bmatrix}$$

표본 공분산 행렬 S 는 다음 식으로 구할 수 있다.

$$S = \frac{1}{N} \sum_{i=1}^N (\overset{\text{column vector}}{x_i - \bar{x}})(\overset{\text{row vector}}{x_i - \bar{x}})^T$$



$$S = \frac{1}{N} \underbrace{X_0^T}_{\text{평균값을 뺀 행}} \underbrace{X_0}_{\text{표본 행렬}} = \text{공분산 행렬}$$

고정값 0

대칭행렬의
특성을 갖는다

다변수 확률변수의 공분산

M개의 다변수 확률변수 벡터

$$X = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_M \end{bmatrix}$$

의 이차적 공분산행렬은 Σ 로 표기하며 다음처럼 정의한다.

$$\Sigma = \text{Cov}[X] = \begin{bmatrix} \sigma_{x_1}^2 & \sigma_{x_1 x_2} & \sigma_{x_1 x_3} & \dots & \sigma_{x_1 x_M} \\ \sigma_{x_1 x_2} & \sigma_{x_2}^2 & \sigma_{x_2 x_3} & \dots & \sigma_{x_2 x_M} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \sigma_{x_1 x_M} & \sigma_{x_2 x_M} & \sigma_{x_3 x_M} & \dots & \sigma_{x_M}^2 \end{bmatrix}$$

$$= \begin{bmatrix} (X_1 - E[X_1])^2 & \dots & (X_1 - E[X_1])(X_M - E[X_M]) \\ (X_1 - E[X_1])(X_2 - E[X_2]) & \dots & (X_2 - E[X_2])(X_M - E[X_M]) \\ \vdots & \ddots & \vdots \\ (X_1 - E[X_1])(X_M - E[X_M]) & \dots & (X_M - E[X_M])^2 \end{bmatrix}$$

다음과 같이 표기할 수도 있다.

$$\Sigma = E[(X - E[X])(X - E[X])^T]$$

$$= E \left[\begin{bmatrix} X_1 - E[X_1] \\ X_2 - E[X_2] \\ \vdots \\ X_M - E[X_M] \end{bmatrix} [X_1 - E[X_1] \ X_2 - E[X_2] \ \dots \ X_M - E[X_M]] \right]$$