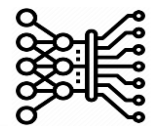


Regression



온도	판매량
20	40
21	42
22	44
23	46

양적



회귀
regression

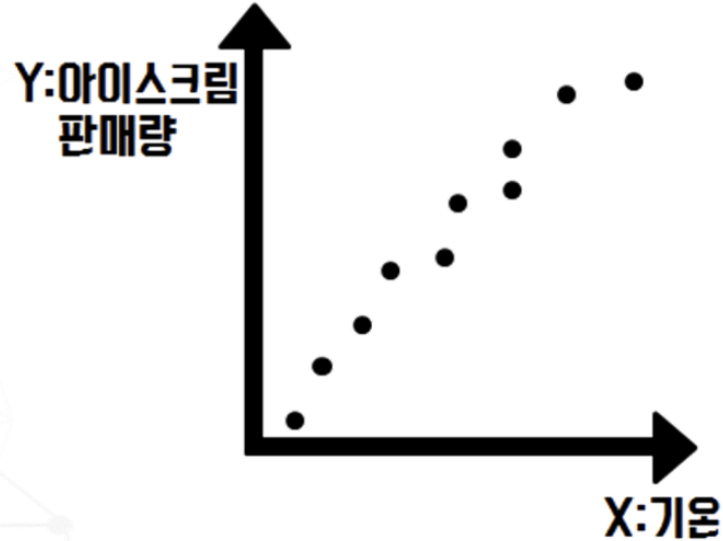
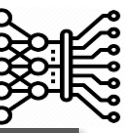
공부시간	시험결과
20	불합격
21	불합격
22	합격
23	합격

범주형

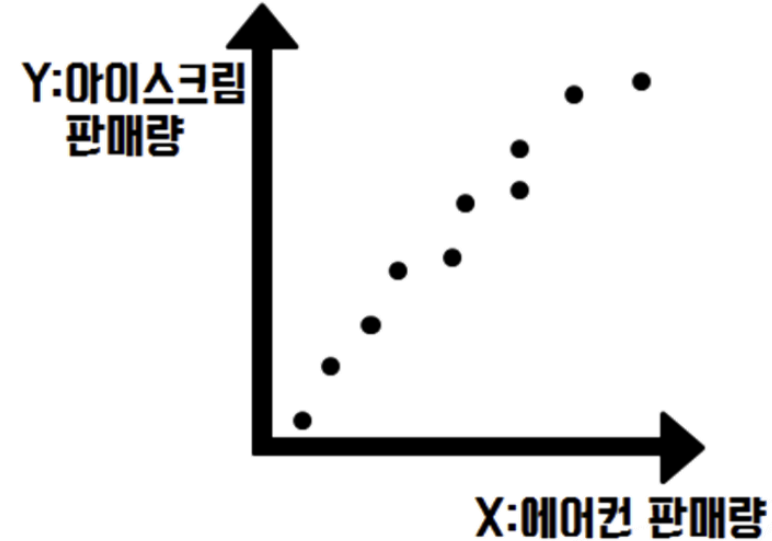


분류
classification

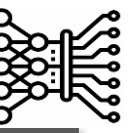
Causal Relationship vs. Correlation



자 봐. 기온이 높아짐에 따라 아이스크림 판매가 증가하는 건 기온과 아이스크림 판매가 **인과적인 관계**에 있는 거야



무슨 소리야 에어컨 판매량이 증가함에 따라 아이스크림 판매량이 증가하는 거야말로 **인과적인 관계**지 바보



❖ Covariance

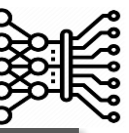
- It is a measure of how two probability variables x and y are related to each other and change, indicating the direction of the relationship between the two probability variables

$$\begin{aligned} \text{Cov}(X, Y) &= E(X - \mu_X)(Y - \mu_Y) \\ &= E(XY) - \mu_X \mu_Y \end{aligned}$$

$$\text{Cov}(aX, bY) = ab \text{Cov}(X, Y)$$

X	1	2	3	4	5
Y	10000	20000	30000	40000	50000

$$\begin{aligned} \text{Cov}(X, Y) &= E\{(1-3)(10000-30000) + (2-3)(20000-30000) + (4-3)(40000-30000) + 5 \\ &\quad - 3(50000-30000)\} = E(40000 + 10000 + 0 + 10000 + 40000) = 20000 \end{aligned}$$

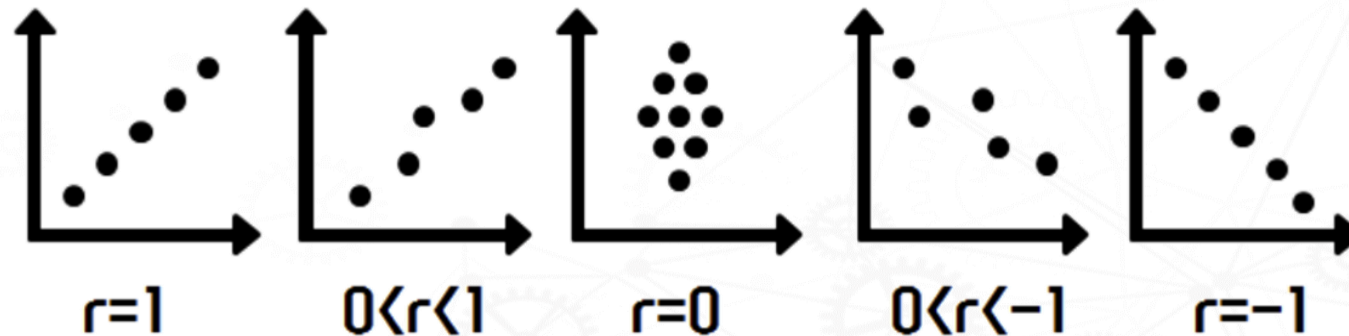


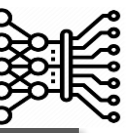
❖ Correlation coefficient

- A numerical value (in coefficients) indicating the degree of correlation between two variables X and Y
- It has a value between -1 and 1, and the closer the absolute value is to 1, the higher the degree of correlation between the two variables

$$\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$$

$$\text{Corr}(aX, bY) = \frac{ab}{|ab|} \text{Corr}(X, Y)$$

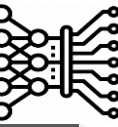




❖ Regression

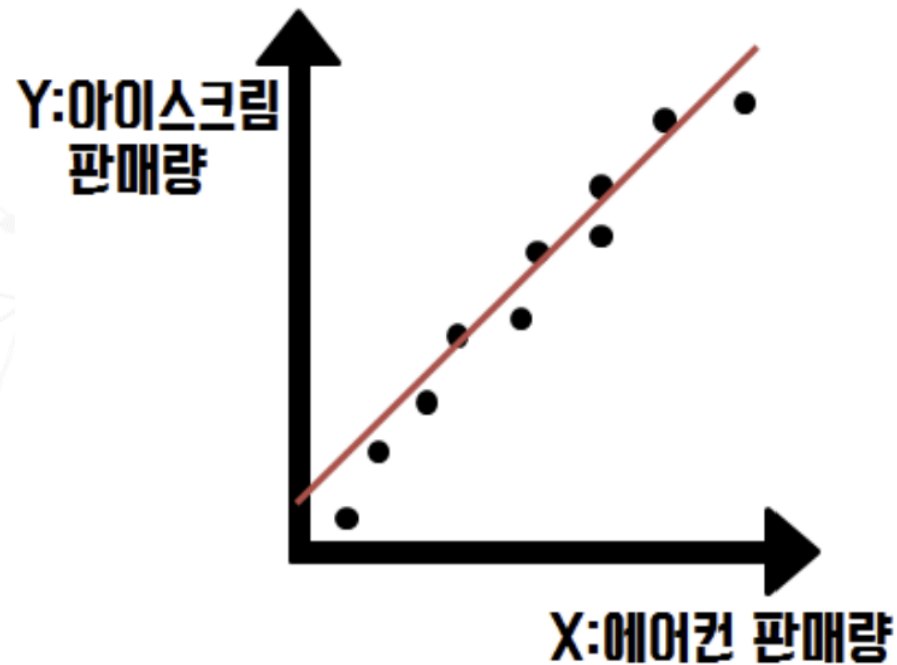
- Analysis method for testing causality
- One or more cause variables (independent variables) affect other variables (dependent variables)
- Independent Variables: Variables that affect dependent variables
- Dependent Variables: Variables affected by other variables

$$X(IV) \rightarrow f(\text{process}) \rightarrow Y(DV)$$

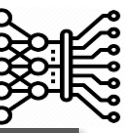


❖ Regression

- A statistical method of assuming a mathematical model and estimating this model from the data of measured variables to determine the relevance of variables
- Predicting the value of a dependent variable based on the value of an independent variable

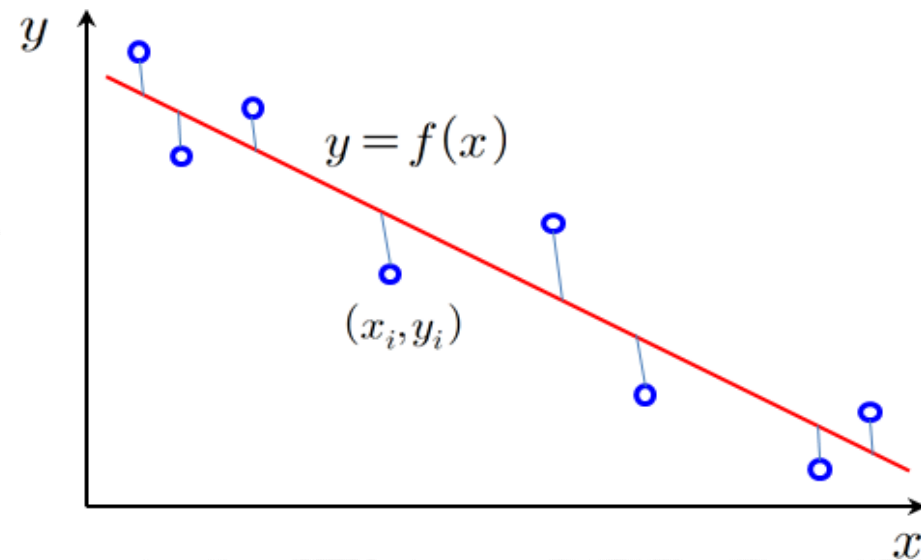


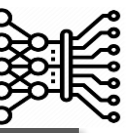
$$\begin{aligned}\text{Cov}(x, y) &= 2000000 \\ r &= 0.8467 \\ Y &= 0.02 + 10.67x\end{aligned}$$



❖ Regression

- An analysis technique that explains the change of dependent variables as a linear combination of independent variables
- A technique for estimating the statistical relationship between dependent and independent variables
- Analyze interrelationships between variables and predict changes in other variables from changes in certain variables

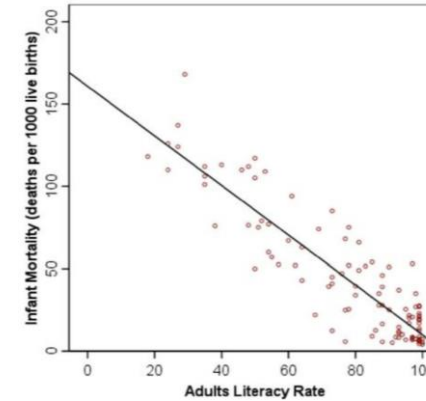




❖ Simple Regression vs. Multiple Regression

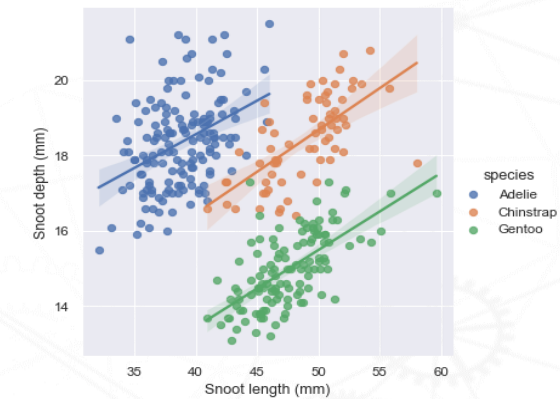
- **Simple Regression:** one independent variable

$$y = a + bx + e$$

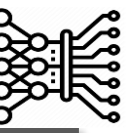


- **Multiple Regression:** Two or more independent variables

$$y = a + b_1 x_1 + b_2 x_2 + b_n x_n + e$$

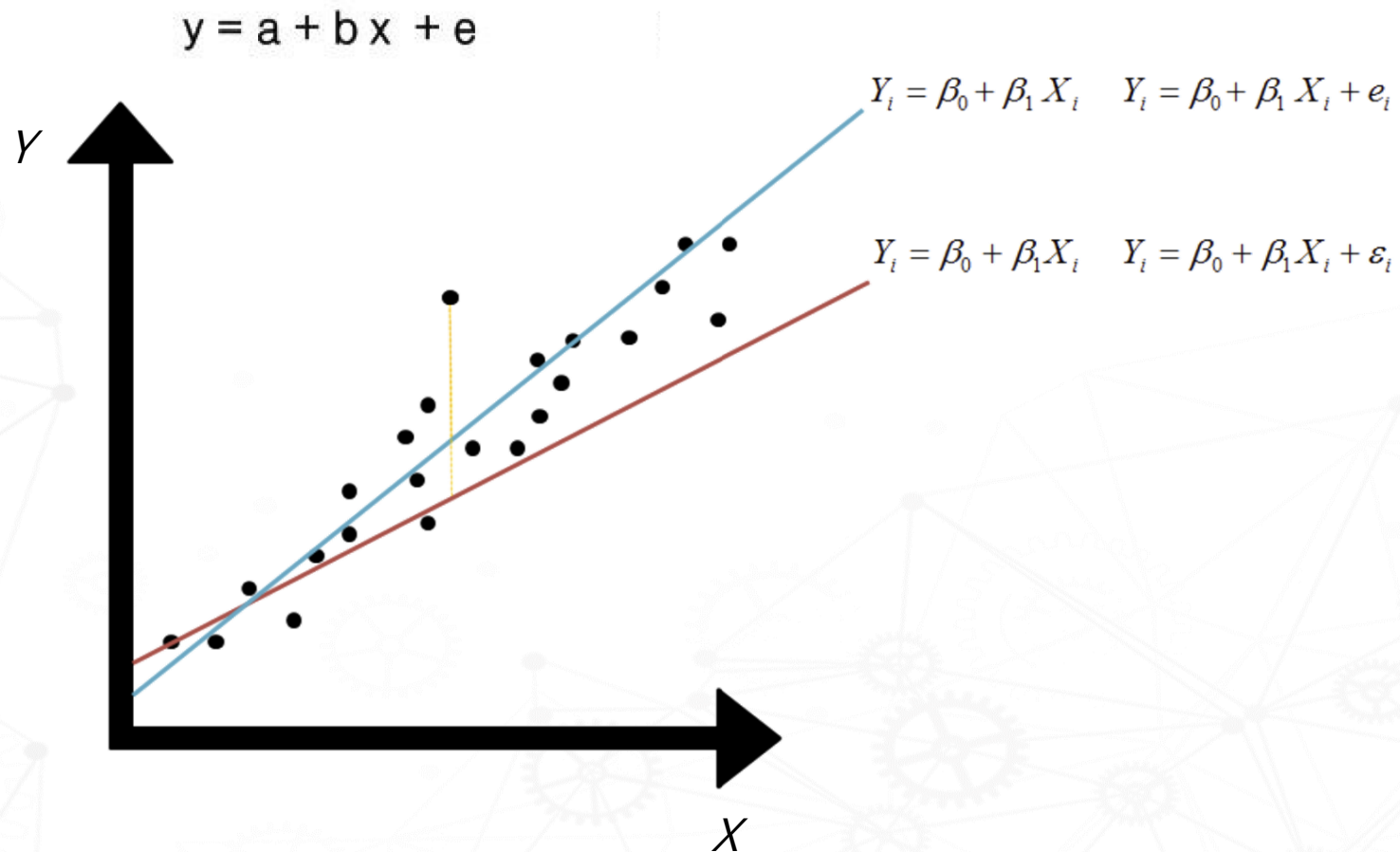


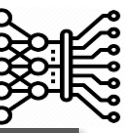
- **e (error term)** : Effects of variables other than independent variables on dependent variables



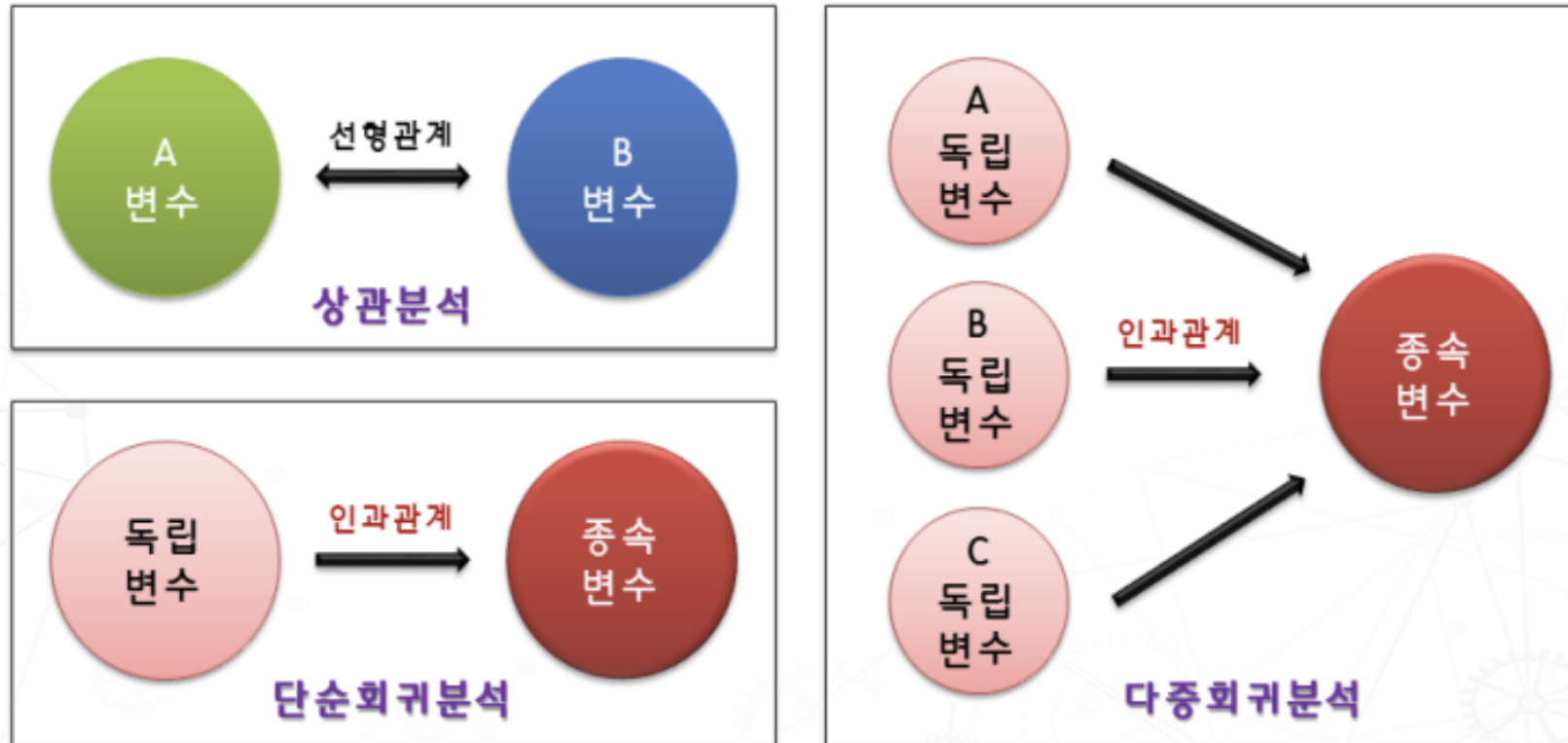
❖ Simple Regression vs. Multiple Regression

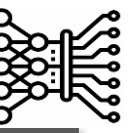
- **Simple Regression:** one independent variable





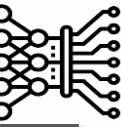
❖ Simple Regression vs. Multiple Regression





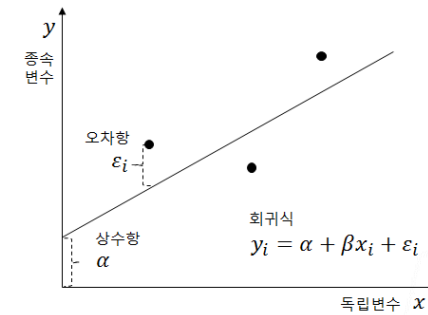
❖ The purpose of regression

- Estimate the value of a dependent variable to the values of an independent variable
- Review relationships between dependent and independent variables
- Review the suitability of regression applications
- Verification of the statistical significance of predictions using regression analysis



❖ Prerequisites for regression assumptions

- Linearity between independent and dependent variables
- The change in the value of the dependent variable according to the change in the value of the independent variable is constant

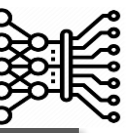


❖ Normality and homoskedasticity error term

- 오차항: 종속변수의 관측치와 예측치 간의 차이
- 오차항의 기대값은 0이며, 일정한 분산을 갖는 정규분포를 이룬다고 가정

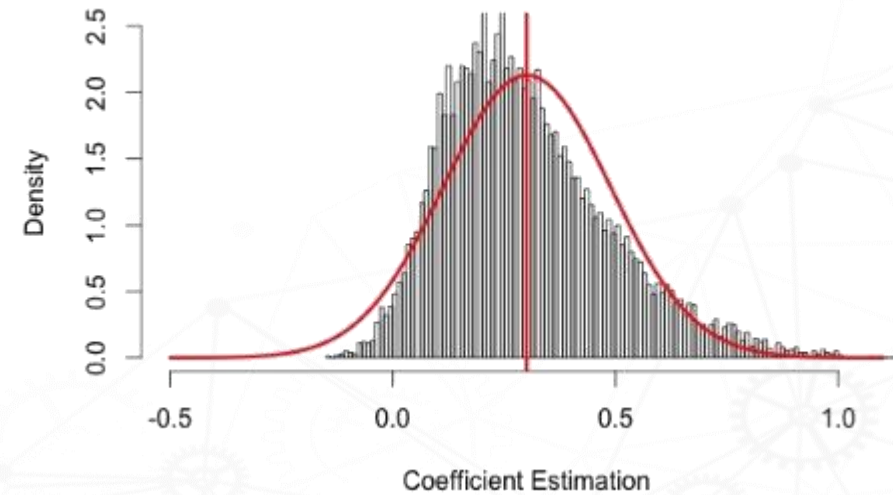
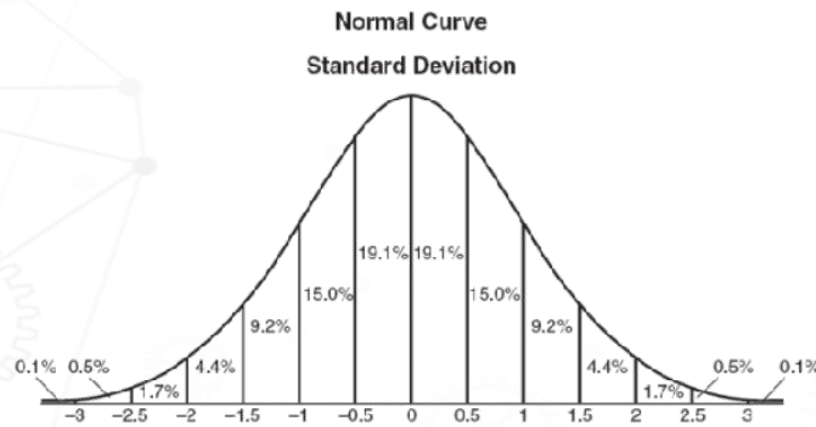
❖ Independency of error term

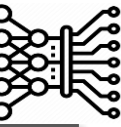
- 예측의 오차값들은 서로 독립적이라는 가정이 필요
- y 의 변화에 따라 오차항이 어떤 패턴을 가져서는 안됨
- y 가 커짐에 따라 오차값이 커지면 가정에 위배



❖ Normality

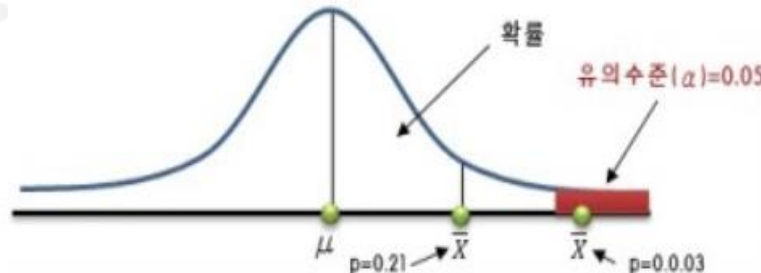
- Assume that the distribution of a population in a continuous variable should be established as a normal distribution
- Normal distribution: Symmetric form relative to the center (average) when plotting the percentage of appearance of a particular value



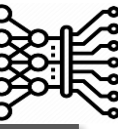


❖ Normality test

- H_0 (귀무가설): 표본의 모집단이 정규분포를 이룸
- H_1 (대립가설): 표본의 모집단이 정규분포를 이루지 않음
- 정규성을 만족하기 위해서는 귀무가설을 채택해야 하며, 대립가설을 기각해야함
 - 유의수준 95% 신뢰구간에서 유의확률 p 값이 0.05보다 크게 나타나야 귀무가설을 채택
 - 유의확률(p -value): 귀무가설 하에서 검정통계량의 값이 나타날 가능성을 측정하는 확률값
 - 유의수준(α): 귀무가설이 참인데, 대립가설을 선택하는 오류의 최대 허용 단계를 의미

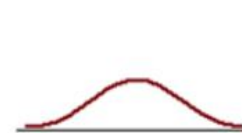
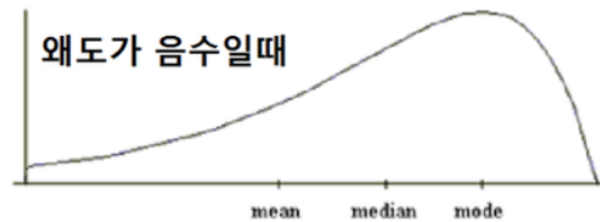
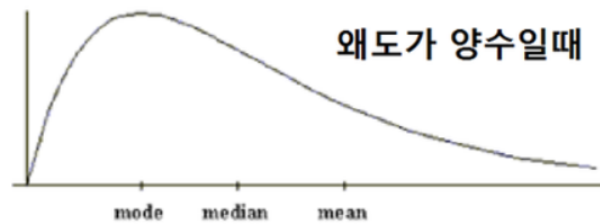


- ✓ 유의확률 < 유의수준 → 귀무가설 기각
- ✓ 유의확률 > 유의수준 → 귀무가설 채택

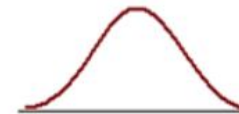


❖ Normality test

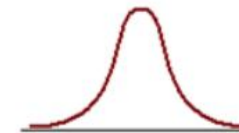
- Central limit theorem
- Kolmogorov-Smirnov & Shapiro-Wilk
- Skewness and kurtosis by Snedecor & Cochran



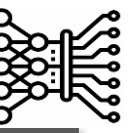
첨도가 음수일때



첨도가 0일때

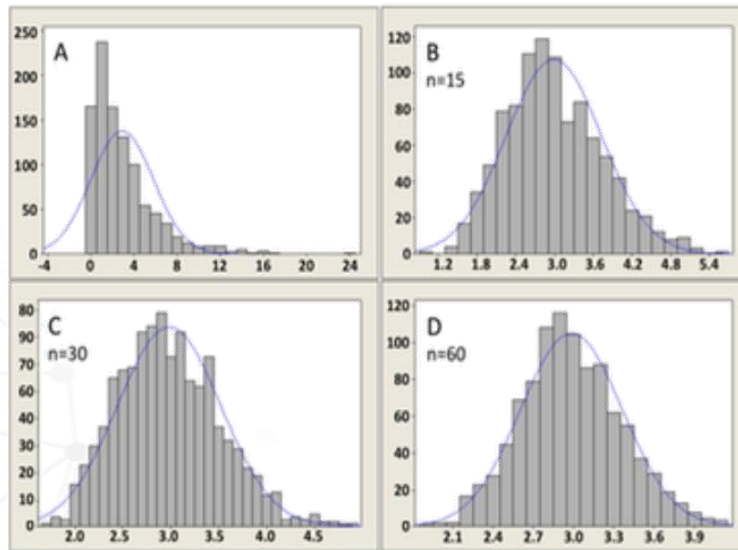


첨도가 양수일때



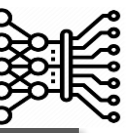
❖ Normality test

■ Central limit theorem



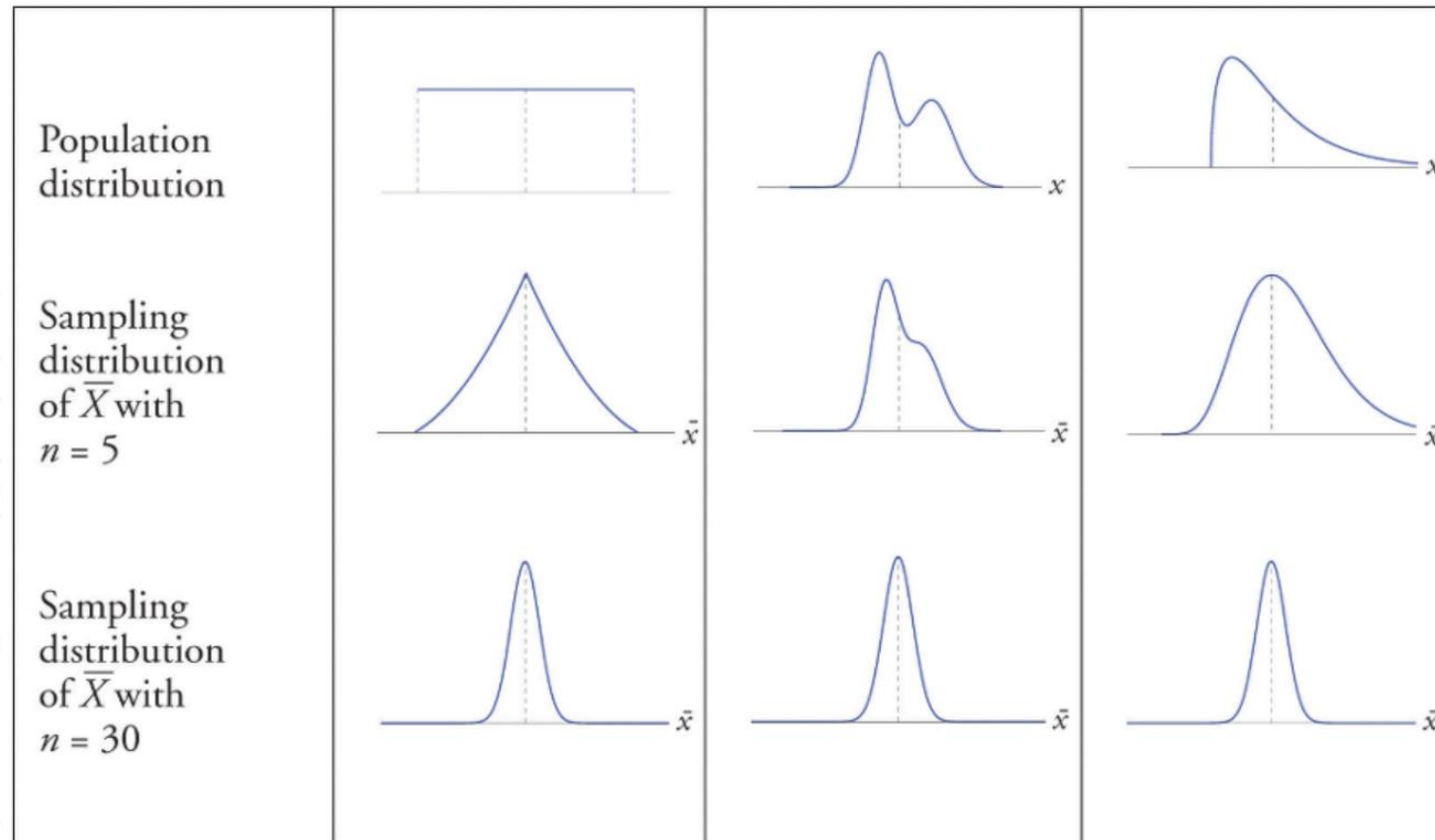
- A는 모집단의 분포를 나타냄(양의 왜도)
- B-C-D는 샘플링 접근법 통해 계산된 그래프로서 표본 수가 증가할 수록 정규 분포에 수렴됨

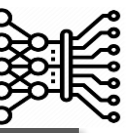
- 표본의 크기가 커질수록 표본 평균의 분포는 모집단의 분포 모양과는 관계 없이 정규 분포에 근접
- 표본 평균의 평균은 모집단의 모 평균과 동일
- 표본 평균의 표준 편차는 모집단의 모 표준 편차를 표본 크기의 제곱근으로 나눈 것
- Data가 대표본(일반적으로, 범주 별 **30개 이상**)이면 대표본 근사(중심극한정리)에 의해 정규성을 만족



❖ Normality test

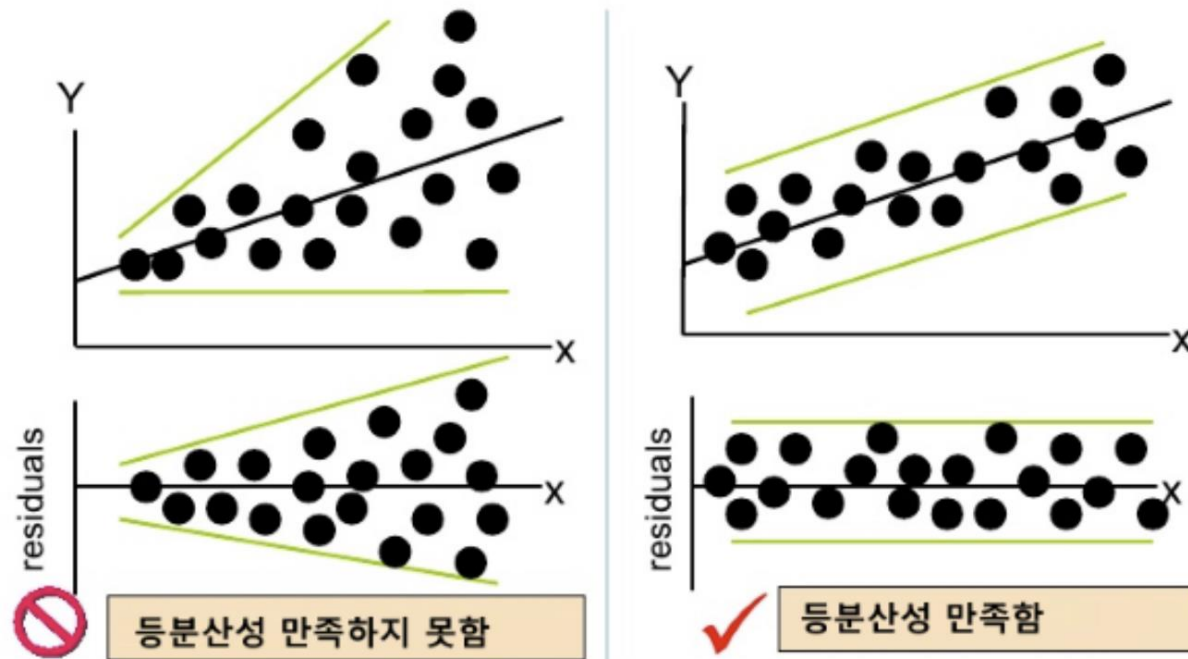
- Central limit theorem

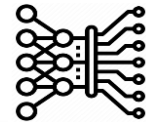




❖ Homoskedasticity

- The condition that two or more different groups must satisfy the same variance through ANOVA
 - The data should evaluate for normality beforehand





❖ Basic principle

$$y = a + bx + e \text{ (회귀식)}$$

$$\hat{y} = \hat{a} + \hat{b}x \text{ (추정식)} \rightarrow \hat{a} \text{ 과 } \hat{b} \text{ 추정}$$

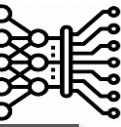
$$e = y - \hat{y} \text{ (관찰치와 예측치의 차이)}$$

❖ Definition of e (error)

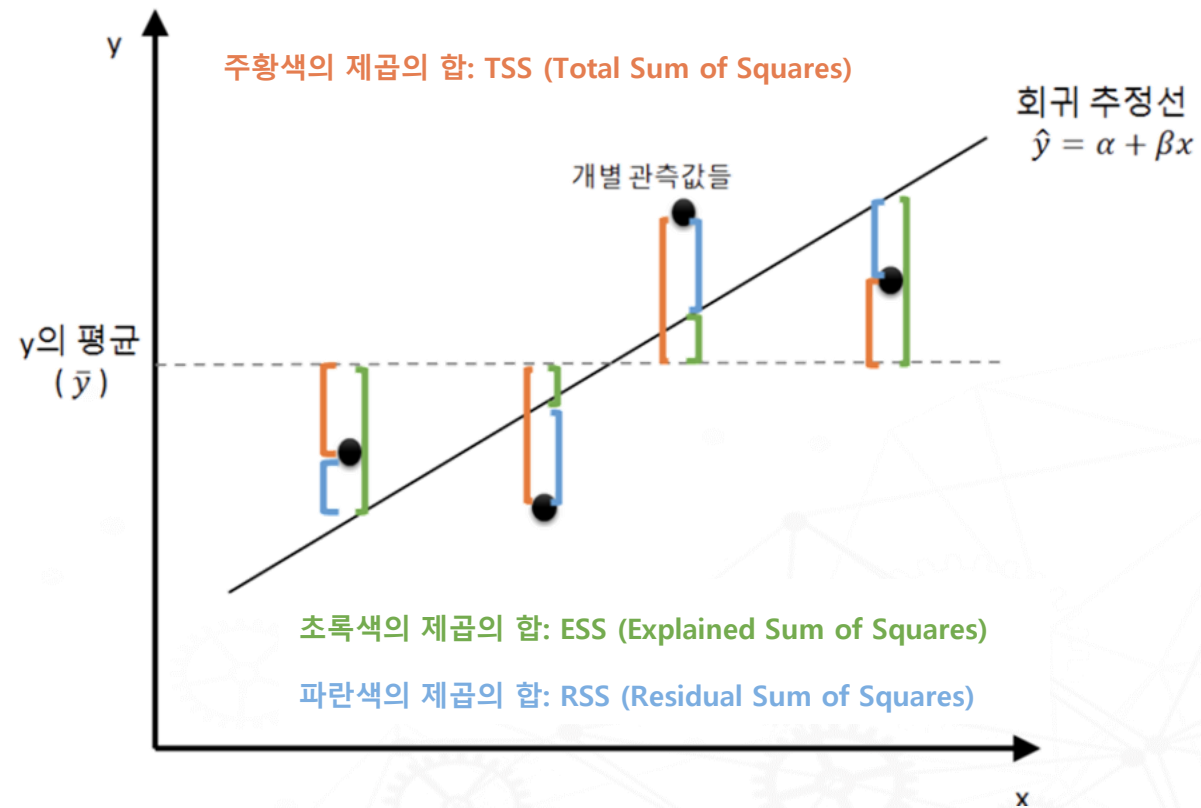
TSS (Total Sum of Square) : $\sum(y - \hat{y})^2$
실제치(y)와 추정치(\hat{y})의 차이의 제곱의 합

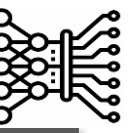
RSS (Residual Sum of Square) : $\sum(y - \bar{y})^2$
실제치(y)와 y의 평균(\bar{y})와의 차이의 제곱의 합
Unexplained Error (회귀선으로 설명이 안 되는 분산)

ESS (Explained Sum of Square) : $\sum(\bar{y} - \hat{y})^2$
y의 평균(\bar{y})와 추정치(\hat{y})의 차이의 제곱의 합

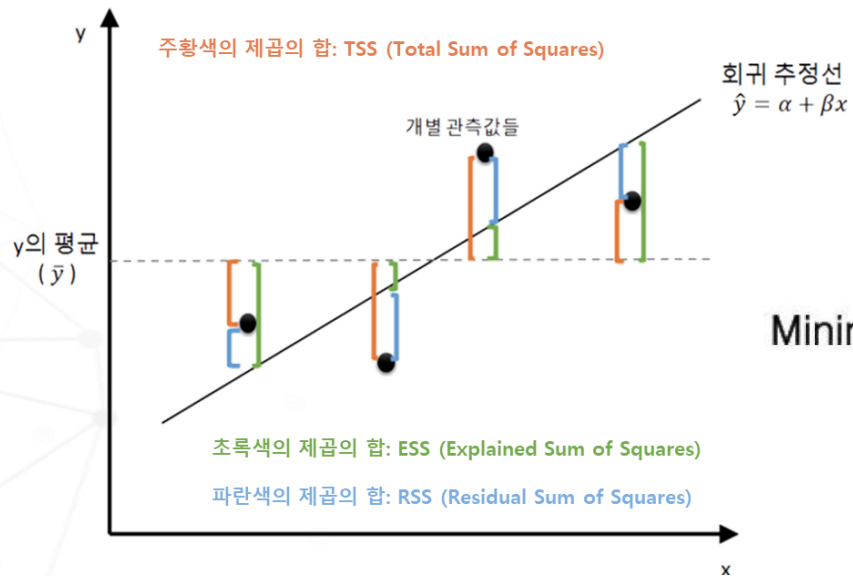


- ❖ Least square method (Ordinary Least Square: OLS)
 - The least squares method that minimizes the sum of the squares of the errors





- ❖ Least square method (Ordinary Least Square: OLS)
 - The least squares method that minimizes the sum of the squares of the errors

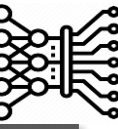


$$\text{Minimize } \sum e_i = \sum [y_i - (a + b x)]^2$$

$$e = y - \hat{y} = (y - \bar{y}) + (\bar{y} - \hat{y})$$
$$\sum (y - \hat{y})^2 = \sum (y - \bar{y})^2 + \sum (\bar{y} - \hat{y})^2$$

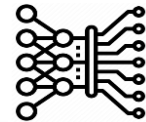
$$\text{TSS} = \text{RSS} + \text{ESS}$$

총 변동 = 설명된 변동 + 설명 안된 변동



❖ Approach

- Simultaneous (or direct) method: “Enter”
 - Derive a regression model from a complete set of independent variables
- Stepwise Method: “Stepwise”
 - A method of sequentially including independent variables one by one in the regression model based on the explanatory power of each independent variable



❖ Goodness of Fit

▪ Multiple R (correlation coefficient)

- 종속변수와 독립변수의 상관관계
- 두 변수의 상관성을 나타내는 척도

▪ 결정계수 R^2 (coefficient of determination)

- 상관계수의 제곱
- x 와 y 간의 상관관계가 클수록 R^2 는 1에 가까워짐
- 회귀식이 자료를 얼마나 잘 설명하고 있는가를 나타냄
- 일반적으로 $R^2 > 0.65$ 일 경우 회귀식이 잘 설명한다고 판단

▪ 수정된 결정계수 (Adjust R^2)

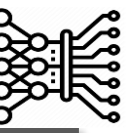
- 독립변수의 수와 데이터 수를 고려한 결정계수
- 변수의 수가 증가할수록 결정계수가 높아지는 단점 존재
- 다중회귀분석에서는 주로 사용하고, 표본의 크기와 독립변수의 수를 고려하여 계산

$$R^2 = \frac{ESS}{TSS} = 1 - \frac{RSS}{TSS}$$

회귀선에 의해 설명되는 변동

$$= \frac{\sum(\bar{y} - \hat{y})^2}{\sum(y - \hat{y})^2} = 1 - \frac{\sum(\bar{y} - \hat{y})^2}{\sum(\hat{y} - \hat{y})^2}$$

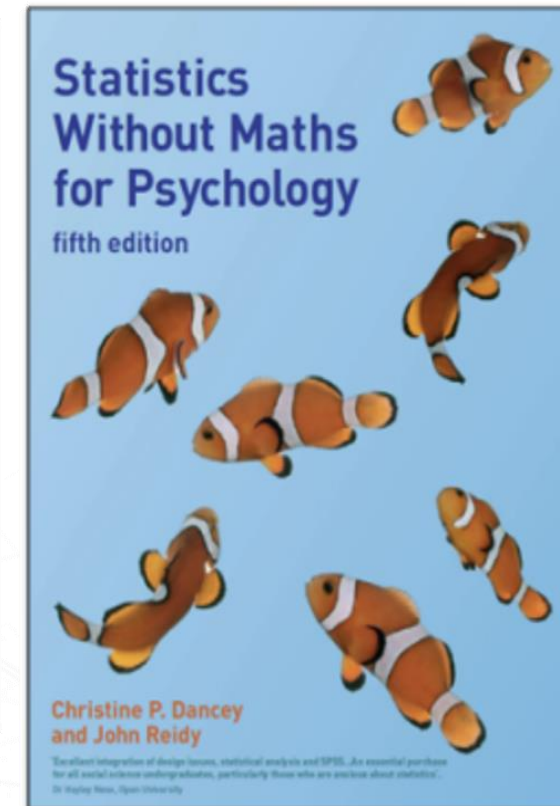
전체 변동

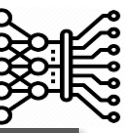


❖ Goodness of Fit

Perfect	+1	-1
Strong	+0.9	-0.9
	+0.8	-0.8
	+0.7	-0.7
Moderate	+0.6	-0.6
	+0.5	-0.5
	+0.4	-0.4
Weak	+0.3	-0.3
	+0.2	-0.2
	+0.1	-0.1
Zero	0	

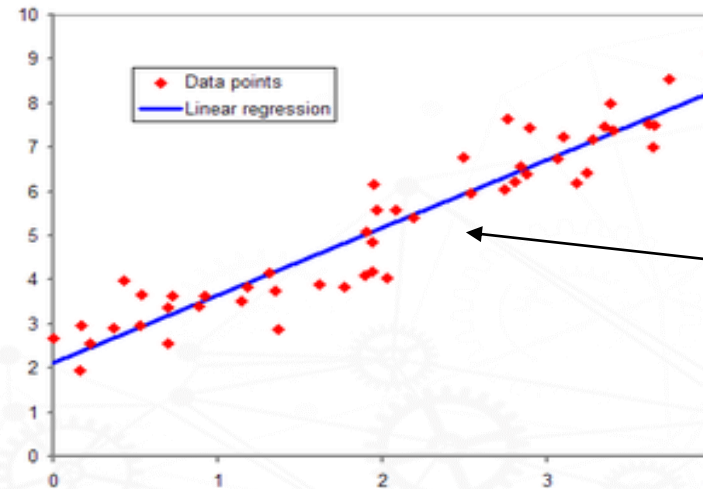
*Christine Dancey and John Reidy, Statistics Without Maths for Psychology, p. 175.
Prentice Hall, 5th edition, 2011.*





❖ Linear regression

- Model the linear correlation between dependent variable Y and one or more independent variable (or explanatory variable) X
- When the parameters of the regression are linear
- Regression coefficients are optimized to minimize RSS (Residual Sum of Squares) between predicted and actual values, and regulations are not applied



$$y = Wx + b$$