

일반통계학

제 4장

표본분포

통계학 2016.1학기 정혜영

1 베르누이 분포와 정규분포

1.1 베르누이 시행

-어떤 실험의 결과를 오직 두 가지 중의 하나로 생각하는 시행

(1) 베르누이 시행의 표본공간 : $S=\{s \text{ ('성공')}, f \text{ ('실패')} \}$

(2) 성공확률 $p=P(\{s\})$, 실패확률 $q=1-p$

1.2 베르누이 확률변수

- 베르누이 표본공간 S 에서 실수로 가는 함수 $X(s)=1, X(f)=0$ 인 X

1.3 베르누이 분포

베르누이 확률변수 X 의 확률분포

(1) 특성값이 이원적인 모집단의 분포를 나타냄

x	0	1	합계
$P(X=x)$	$1-p$	p	1

(2) 기호 $X \sim B(1,p)$ or $Ber(p)$

(3) 평균과 분산

$$E(X) = p$$

$$V(X) = p(1 - p)$$

1 베르누이 분포와 정규분포

1.4 베르누이 시행이 독립적으로 반복되는 실험

(1) 기호 : $X_1, \dots, X_n \sim B(1, p)$ iid (Independent identically distributed)

(2) 10개 중 3개의 당첨제비가 있을 때, 2개의 제비를 단순랜덤 **복원**추출하는 경우
 X_1 ='첫번째 시행의 결과', X_2 ='두번째 시행의 결과'

p

.

※ n 개의 제비를 단순랜덤 **복원**추출하는 경우에 X_1, \dots, X_n 을 각 시행의 결과라고 하면, X_1, \dots, X_n 은 서로 독립이며 동일한 분포를 갖는다.

1 베르누이 분포와 정규분포

(예) 10개 중 3개의 당첨제비 중 2개의 제비를 단순랜덤비복원추출하는 경우
 X_1 ='첫번째 시행의 결과', X_2 ='두번째 시행의 결과'
 X_1 과 X_2 는 서로 독립은 아니지만 동일한 분포를 갖는다.

$$P(X_1) = 3/10 \text{ and } P(X_2) = 3/10 = 3/10 * 2/9 + 7/10 * 3/9$$

(예) N 개 중 M 개의 당첨제비 중 2개의 제비를 단순랜덤비복원추출하는 경우
 N 이 M 에 비해 충분히 클 경우

$$P(X_1=1) \cong P(X_2=1|X_1=1) \cong P(X_2=1|X_1=0)$$

$$(\because \frac{M}{N} \cong \frac{M-1}{N-1} \cong \frac{M}{N-1})$$

즉, X_1 과 X_2 는 독립에 가깝게 된다. 모집단의 크기가 표본의 크기에 비해 충분히 클 때, 유한모집단에서 표본을 단순랜덤비복원추출하는 것은 단순랜덤복원추출하는 것과 별 차이가 없다.

1 베르누이 분포와 정규분포

1.5 정규분포(normal distribution) - 가우스 분포

- 특성값이 연속적인 무한모집단의 분포로서 가장 대표적인 분포.
- 정규분포를 따르는 무한모집단에서 하나의 값을 관측하는 실험을 할 때, 이를 X 라고 하면 확률변수 X 의 확률분포는 정규분포임.

(1) 평균 μ , 표준편차 σ 인 정규분포의 확률밀도 함수

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, -\infty < x < \infty$$

(2) 평균 μ 를 중심으로 대칭

(3) 표준정규분포 : $Z \sim N(0,1)$, $X \sim N(\mu, \sigma^2)$ 일 때, $Z = \frac{X-\mu}{\sigma} \sim N(0,1)$

(4) $z_\alpha : P(Z \geq z_\alpha) = \alpha$ 인 값, 즉 $100(1-\alpha)\%$ 백분위수

(5) $P(a \leq X \leq b) = P\left(\frac{a-\mu}{\sigma} \leq Z \leq \frac{b-\mu}{\sigma}\right)$

1 베르누이 분포와 정규분포

$$P(0 \leq Z \leq z)$$



TABLE 1 Normal Curve Areas

z	.00	.01	.02	.03	.04	.05	.06	.07	.08	.09
0.00	.0000	.0040	.0080	.0120	.0160	.0199	.0239	.0279	.0319	.0359
0.10	.0398	.0438	.0478	.0517	.0557	.0596	.0636	.0675	.0714	.0753
0.20	.0793	.0832	.0871	.0910	.0948	.0987	.1026	.1064	.1103	.1141
0.30	.1179	.1217	.1255	.1293	.1331	.1368	.1406	.1443	.1480	.1517
0.40	.1554	.1591	.1628	.1664	.1700	.1736	.1772	.1808	.1844	.1879
0.50	.1915	.1950	.1985	.2019	.2054	.2088	.2123	.2157	.2190	.2224
0.60	.2257	.2291	.2324	.2357	.2389	.2422	.2454	.2486	.2517	.2549
0.70	.2580	.2611	.2642	.2673	.2704	.2734	.2764	.2794	.2823	.2852
0.80	.2881	.2910	.2939	.2967	.2995	.3023	.3051	.3078	.3106	.3133
0.90	.3159	.3186	.3212	.3238	.3264	.3289	.3315	.3340	.3365	.3389
1.00	.3413	.3438	.3461	.3485	.3508	.3531	.3554	.3577	.3599	.3621
1.10	.3643	.3665	.3686	.3708	.3729	.3749	.3770	.3790	.3810	.3830
1.20	.3849	.3869	.3888	.3907	.3925	.3944	.3962	.3980	.3997	.4015
1.30	.4032	.4049	.4066	.4082	.4099	.4115	.4131	.4147	.4162	.4177
1.40	.4192	.4207	.4222	.4236	.4251	.4265	.4279	.4292	.4306	.4319
1.50	.4332	.4345	.4357	.4370	.4382	.4394	.4406	.4418	.4429	.4441
1.60	.4452	.4463	.4474	.4484	.4495	.4505	.4515	.4525	.4535	.4545

1 베르누이 분포와 정규분포

- 표준정규분포의 확률 계산
 - 예) Z가 0과 1.21 사이에 포함될 확률

z	.00	.01	.02
⋮			
.9	.3159	.3186	.3212
1.0	.3413	.3438	.3461
1.1	.3643	.3665	.3688
1.2	.3849	.3869	.3888
⋮			

- 예) Z가 구간 -1.00과 0.00사이에 있을 확률

$$P(-1.00 \leq Z \leq 0) = P(0 \leq Z \leq 1.00) = 0.3413$$

⇒ 표준정규확률밀도함수는 0을 중심으로 대칭

1 베르누이 분포와 정규분포

(예) 우체국에서 소포 무게의 상한선을 설정하고자 하여 기존 고객이 부치는 짐의 무거운 5% 정도를 제한하고자 한다. 만약 기존 고객의 소포 무게의 분포가 평균 5kg, 표준편차 1kg인 정규분포를 따른다면 상한선은 얼마로 해야 할까?

$$(a-5)/1 = 1.65$$

1 베르누이 분포와 정규분포

1.6 정규분포의 성질

(1) $X \sim N(\mu, \sigma^2)$ 일 때, 임의의 상수 a, b 에 대하여 $aX + b \sim N(a\mu + b, a^2\sigma^2)$

(2) $X_1 \sim N(\mu_1, \sigma_1^2)$, $X_2 \sim N(\mu_2, \sigma_2^2)$ 이고 X_1 과 X_2 가 서로 독립일 때,
 $a_1X_1 + a_2X_2 \sim N(a_1\mu_1 + a_2\mu_2, a_1^2\sigma_1^2 + a_2^2\sigma_2^2)$

(예) 한 가정에서 아침식사로 빵과 우유 한 잔을 먹는다고 하자. 빵에서 섭취하는 칼로리 $\sim N(200, 225)$, 우유에서 섭취하는 칼로리 $\sim N(80, 25)$ 일 때, 300칼로리 이상을 섭취할 확률은 얼마인가?

$$\sum (X_i - \bar{X}) = 0$$

- $$\sum (X_i - \mu) (X) (\mu)$$
- 가

가

Q)통계량은 확률변수인가? YES

$$\text{예)) } \hat{\mu} = \bar{X}, \quad \hat{\sigma}^2 = S^2$$

• 표본분포는 추정의 정확도를 나타내는 중요한 도구

(Why? 참값인 모수를 모를 때, 표본으로부터의 추정값이 모수를 얼마나 정확하게 추정한 것인지 알 수가 없다. 표본분포를 이용해서 정확도에 대한 확률적 접근이 가능해짐)

(예) 유한모집단 $\{2,2,4,5\}$ 에서 크기 2인 표본을 단순랜덤비복원추출하였을 때, 표본평균의 표본분포를 구하여라. 4.6

- 통계량의 확률분포(표본분포)는 **모집단의 분포와 표본의 추출방법에** 의하여 결정된다.
- 랜덤표본 (Random sample)
 - 유한 모집단인 경우 단순랜덤비복원추출로 뽑은 확률변수의 모임 $\{X_1, X_2, \dots, X_n\}$
 - 무한모집단인 경우 다음의 두 조건을 만족하는 확률변수의 모임 $\{X_1, X_2, \dots, X_n\} \Rightarrow$
 - (i) X_1, X_2, \dots, X_n 각각의 확률분포는 모집단이 분포와 동일
 - (ii) X_1, X_2, \dots, X_n 은 서로 독립
- 유한모집단에서의 랜덤표본인 경우 표본분포의 유도는 매우 복잡하고 어렵다. 따라서 모집단의 크기가 큰 경우 흔히 무한모집단에서의 랜덤표본으로 간주하여 표본분포를 구한 다음 이를 실제 표본분포의 근사분포로 사용한다.
- 이후의 **모든 추정과 검정에서는 무한모집단에서의 랜덤표본을 가정한다.**

3.1 초기하분포 (hypergeometric distribution)

두가지 특성값($\{s, f\}$)만 가지는 유한 모집단의 모비율을 추정하기 위해 표본비율을 사용할 때 표본비율의 확률분포를 나타내기 위해 사용한다.

특성값 1의 개수가 D , 0의 개수가 $N-D$ 인 크기 N 의 유한모집단에서 크기 n 인 랜덤표본을 뽑을 때, 표본 1의 개수의 확률분포를 초기하 분포라고 한다.

$$p(x) = \frac{\binom{D}{x} \binom{N-D}{n-x}}{\binom{N}{n}}, \quad x = 0, 1, \dots, n \quad (\text{단, } n \leq D, n \leq N - D)$$
$$= np, \quad = np(1-p)(N-n)/(N-1)$$

(예) 신세대 의식구조 조사에서 모집단 크기가 10이며, 이 중 찬성자의 수가 6명, 반대자의 수가 4명이라고 하자. 크기 3인 랜덤표본에서 찬성자수의 평균과 분산을 구하여라.

(풀이) X =찬성자의 수, $p(x) = \frac{\binom{6}{x} \binom{10-6}{3-x}}{\binom{10}{3}} : x = 0, 1, 2, 3$

$$u = 3 * 6/10, \quad V = n * 6/10 * 4/10 * (10 - 3)/(10-1)$$

3 초기하분포와 이항분포

- 초기하 분포의 평균과 분산
 - 평균= np , 단, $p = D/N$
 - 분산= $np(1 - p) \cdot \frac{N-n}{N-1}$

3.2 이항분포 (binomial distribution)

이제, 특성값이 1또는 0과 같이 이원적으로 분류되고 1의 비가 p ($0 < p < 1$)인 무한모집단을 생각해보자. 랜덤포본 X_1, X_2, \dots, X_n 은 서로 독립이며 각각 1의 확률이 p 인 베르누이 분포를 따른다. 1을 성공, 0을 실패로 생각하면 X_1, X_2, \dots, X_n 은 성공률 p 인 베르누이 시행을 n 번 반복 시행할 때, 시행의 결과이며 $X =$ 표본에서 1의 개수 = 성공의 수 $= X_1 + X_2 + \dots + X_n \sim B(n, p)$

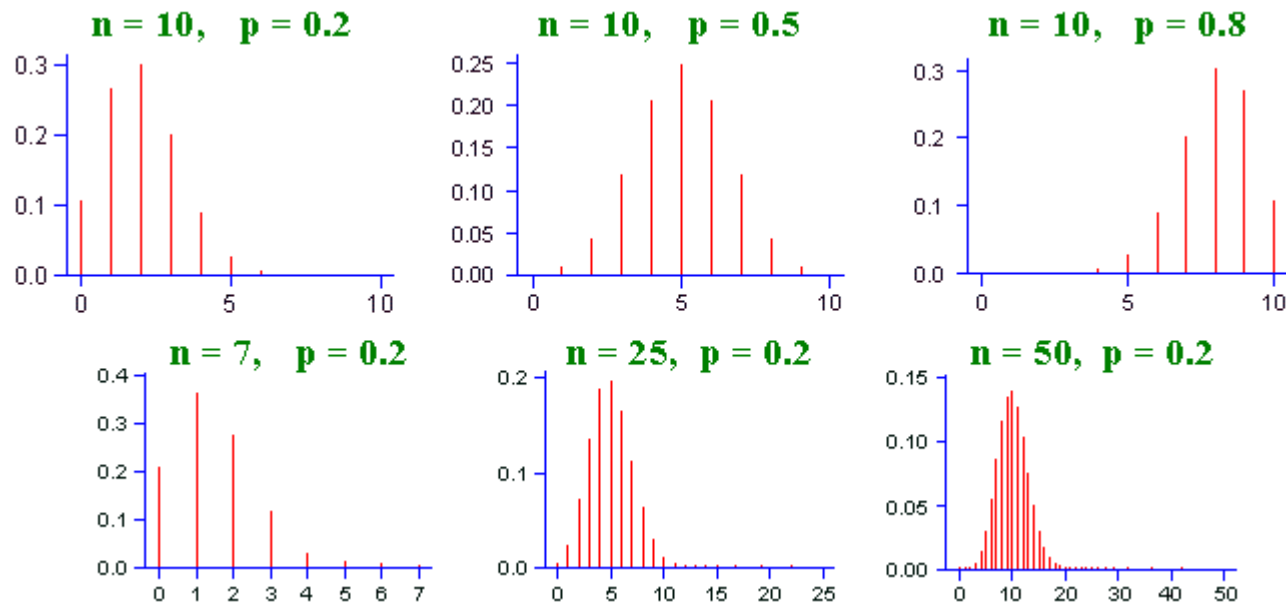
이항분포 n 가 .
시행 횟수가 n 이고 성공의 확률이 p 일 때 성공의 수에 대한 확률분포

$$p(x) = \binom{n}{x} p^x (1 - p)^{n-x}, \quad x = 0, 1, 2, \dots, n$$

- 이항 분포의 평균과 분산
 - 평균= np $E(X) = E(\text{sum } X_i) = np$
 - 분산= $np(1 - p)$ $V(X) = V(\text{sum } X_i) = np(1-p)$ (* 가)



3 초기하분포와 이항분포



(예) 5개 중 하나를 택하는 선다형 문제가 20문항 있는 시험에서 랜덤하게 답을 써 넣는 경우에 다음 물음에 답하여라. 4.8

- (a) 정답이 하나도 없을 확률은 얼마인가?
- (b) 8개 이상의 정답을 맞힐 확률은 얼마인가?
- (c) 4개부터 6개 사이의 정답을 맞힐 확률은 얼마인가?

3 초기하분포와 이항분포

1의 비 D/N 이 p 에 무한히 가까워질 때, 초기하분포는 이항분포 $B(n, p)$ 에 근사한다.

(예) 어떤 제품을 생산하는 공정의 불량률은 5%로 알려져 있다. 오늘 생산한 10000개의 제품 중에서 20개를 단순랜덤추출하여 조사할 때 표본의 불량률이 10%이상일 확률을 구하여라.

(풀이) $X=20$ 개 중 불량품의 수, X 는 초기하 분포를 따르며 근사적으로 $B(20, 0.05)$ 를 따른다.

가

.

4 표본평균의 분포

- X_1, X_2, \dots, X_n 이 모평균 μ , 모분산 σ^2 인 무한모집단으로부터의 랜덤표본인 경우 표본평균 $\bar{X} = \frac{1}{n}(X_1 + X_2 + \dots + X_n)$ 의 기대값과 분산은

$$-E(\bar{X}) = \frac{1}{n} E(\sum X_i) = \frac{1}{n} \cdot n\mu = \mu$$

$$(*)$$

$$-Var(\bar{X}) = \frac{1}{n^2} V(\sum X_i) = \frac{1}{n^2} \cdot n \sigma^2 = \sigma^2 / n$$

$$-sd(\bar{X}) = \sigma / \sqrt{n}$$

- 모집단의 분포가 정규분포 $N(\mu, \sigma^2)$ 일 때, 표본평균 \bar{X} 는 정규분포 $N(\mu, \frac{\sigma^2}{n})$ 을 따른다.

4 표본평균의 분포

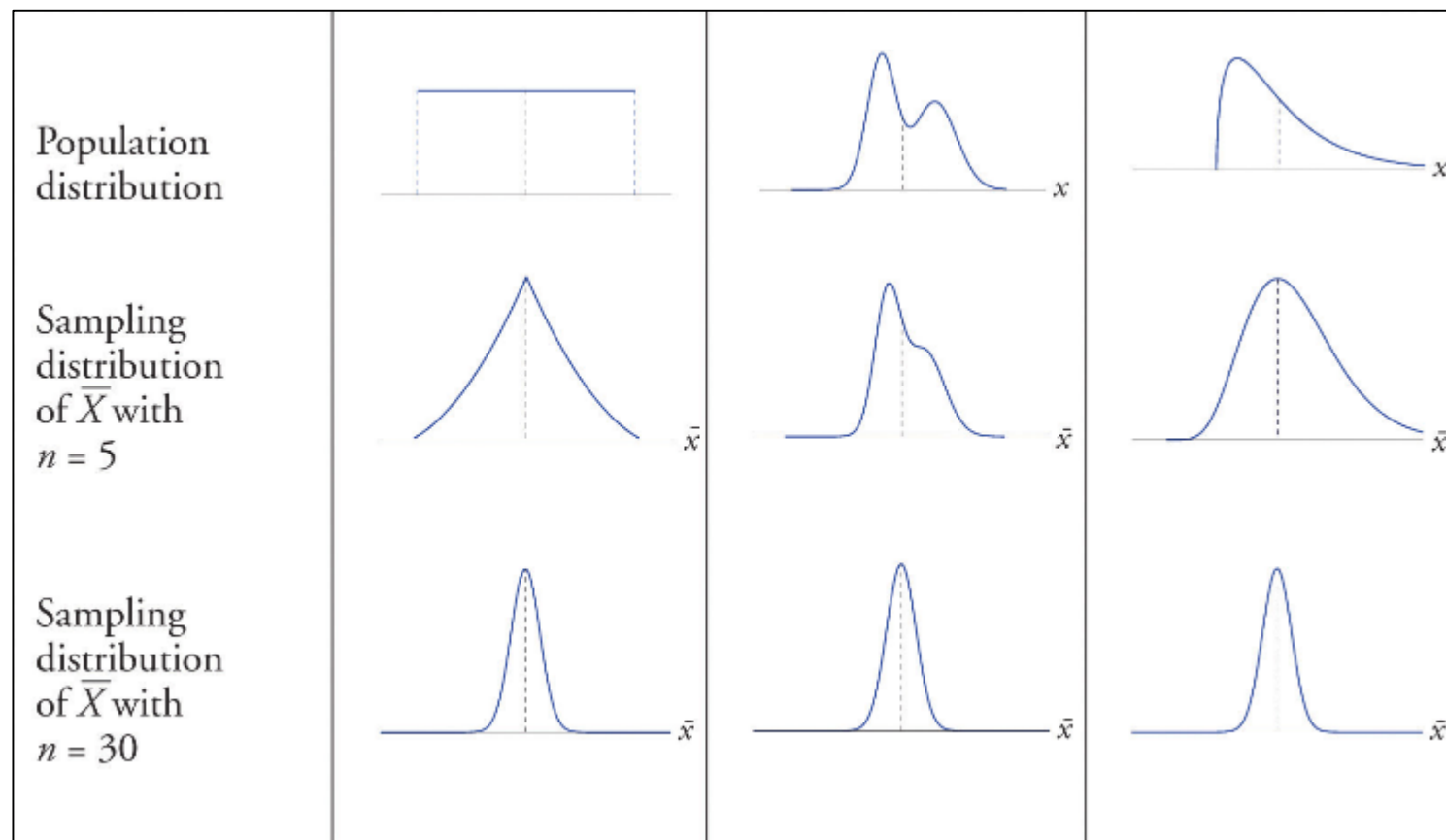
중심극한정리

X_1, X_2, \dots, X_n 이 모평균 μ , 모분산 σ^2 인 무한모집단으로부터의 랜덤표본인 경우, n 이 충분히 크면 (일반적으로 30이상)

$$\bar{X} \approx N\left(\mu, \frac{\sigma^2}{n}\right)$$
$$\text{즉, } \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \approx N(0, 1)$$

(예) 대학 신입생 신장의 평균이 168cm이고 표준편차가 6cm라고 알려져 있다. 100명의 신입생을 단순랜덤추출하는 경우 표본평균이 167cm이상 169cm이하일 확률을 구하여라. 4.10)

4 표본평균의 분포



n

,

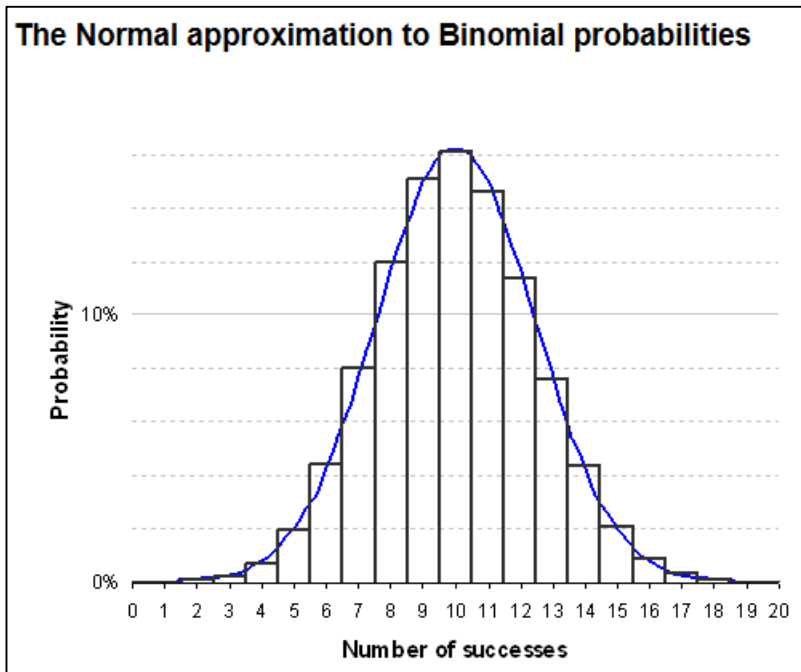
,

.

4 표본평균의 분포

- 이항분포의 정규근사 ($np > 5, n(1-p) > 5$)
 $X \sim B(n, p)$ 이고 n 이 충분히 클 때, X 는 근사적으로 정규분포 $N(np, np(1-p))$ 를 따른다. 즉, $\frac{\bar{X} - np}{\sqrt{np(1-p)}} = \frac{\hat{p} - p}{\sqrt{p(1-p)/n}} \approx N(0,1)$ (단, $\hat{p} = X/n$)

- 연속성 수정계수



$$P(X < 5) = P(X \leq 4) = N(np, np(1-p)) \quad 4.5$$

$$. (4 \quad 1/2 \quad !)$$

$$P(3 \leq X \leq 5) = N(np, np(1-p)) \quad 2.5$$

$$5.5$$

$$4.11$$

$$!$$

5 카이제곱, t, F 분포 (필기노트)

- 카이제곱분포

Z_1, Z_2, \dots, Z_v 가 $N(0,1)$

$$\chi^2(v) = \sum Z_i^2$$

가 v

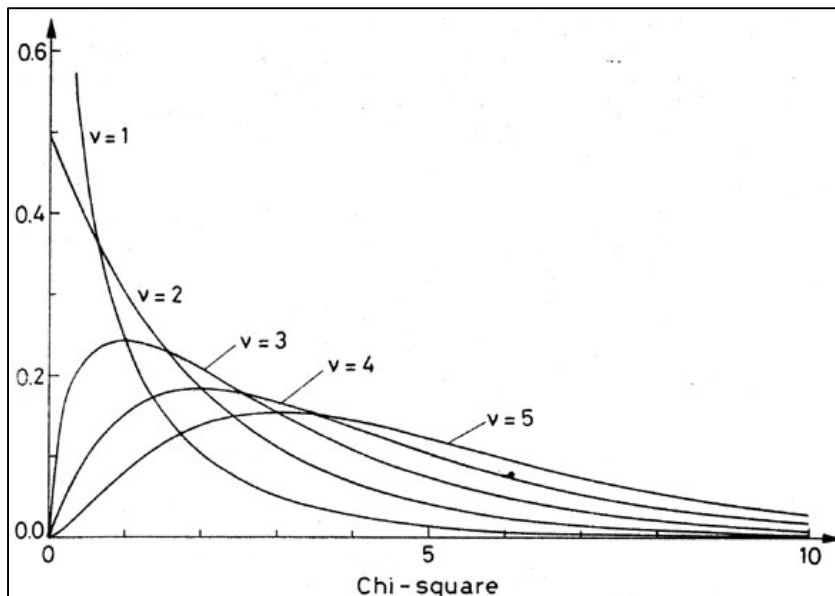
.

가

$X_1 \sim \chi^2(v_1), X_2 \sim \chi^2(v_2)$

X_1, X_2 가

$X_1 + X_2 \sim \chi^2(v_1 + v_2)$



5 카이제곱, t, F 분포 (필기노트)

- 정규모집단에서의 표본분산의 분포

X_1, X_2, \dots, X_n 이 $N(\mu, \sigma^2)$ 으로부터의 랜덤표본이라 할 때,
표본분산 $S^2 = \sum_{i=1}^n (x - \bar{x})^2 / (n - 1)$ 에 대하여 다음이 성립한다.

$$\frac{(n-1)S^2}{\sigma^2} \sim \chi^2(n-1)$$

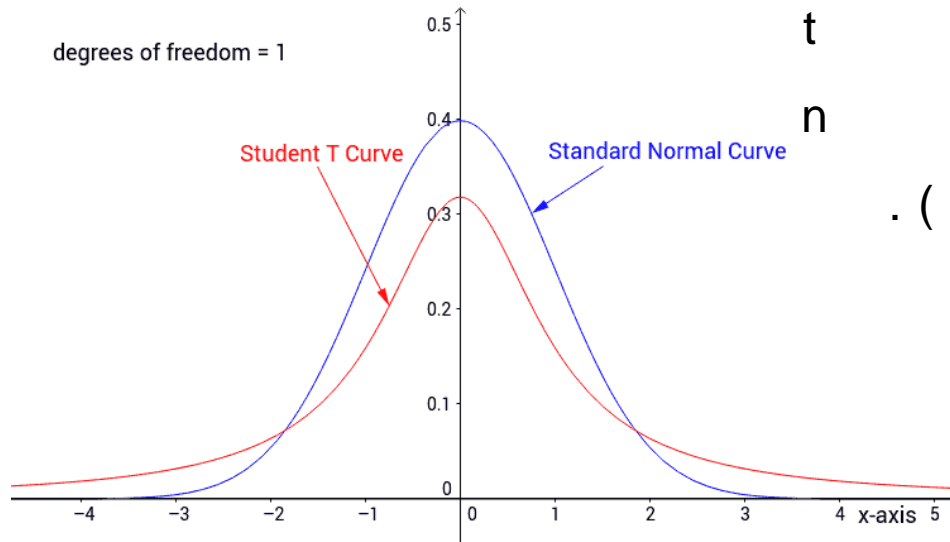
5 카이제곱, t, F 분포 (필기노트)

- **분산이 동일한** 정규모집단에서의 표본분산의 분포

$$\frac{(n_1+n_2-2)S_p^2}{\sigma^2} \sim \chi^2(n_1 + n_2 - 2), \quad \text{단, } S_p^2 = \frac{(n_1-1)S_1^2 + (n_2-1)S_2^2}{n_1+n_2-2}$$

- t 분포

5 카이제곱, t, F 분포 (필기노트)



• 스튜던트화된 표본평균의 분포

$$(\bar{X} - \mu) / (S / \sqrt{n}) \sim t(n - 1)$$

$$(\bar{X} - \mu) / (\sigma / \sqrt{n}) \sim N(0, 1)$$

$$(n-1)S^2 / \sigma^2 \sim \chi^2(n-1)$$

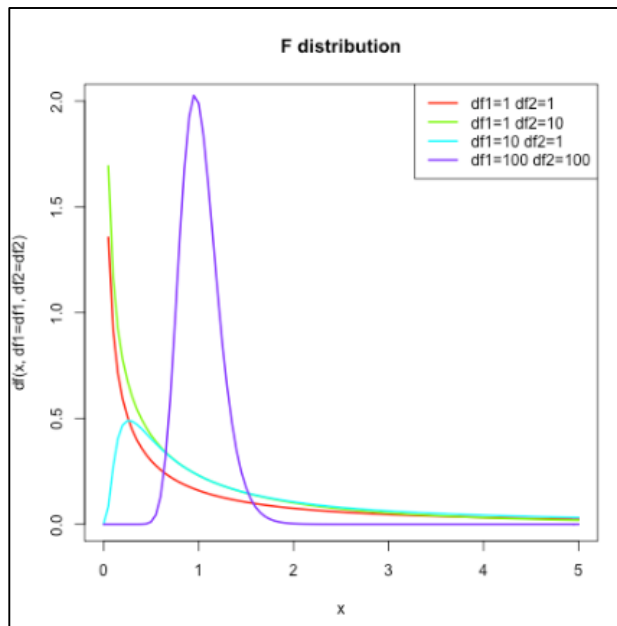
$$\bar{X} \pm S \cdot t$$

5 카이제곱, t, F 분포 (필기노트)

- 분산이 동일한 두 정규모집단에서의 t분포
- 분산이 동일하지 않은 두 정규모집단에서의 t분포

5 카이제곱, t, F 분포 (필기노트)

- F 분포



5 카이제곱, t, F 분포 (필기노트)

- F 분포의 성질

- F 분포와 t분포의 관계

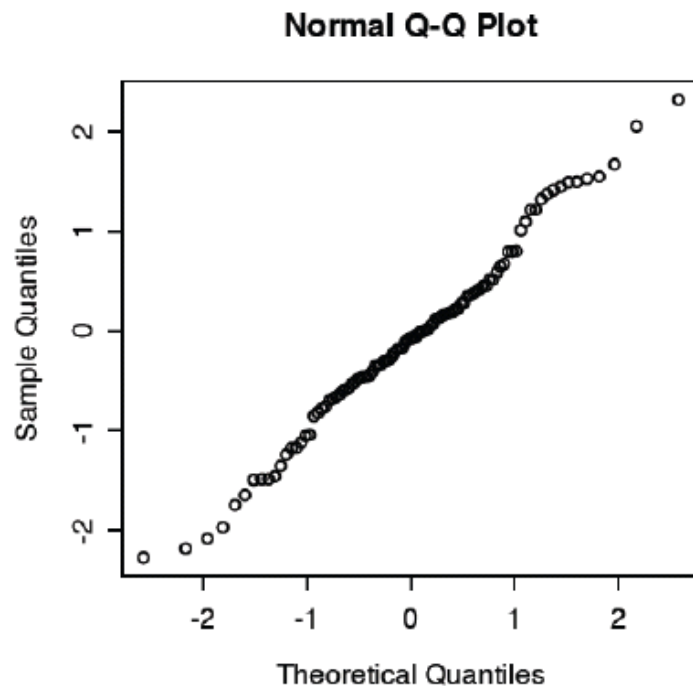
5 카이제곱, t , F 분포 (필기노트)

- 두 정규모집단에서의 표본분산의 비에 대한 분포

5 카이제곱, t, F 분포 (필기노트)

- 정규분포 분위수 대조도

- 모집단의 분포가 정규분포라는 가정을 검토하기 위한 방법
- 정규분포의 분위수와 이에 대응하는 자료 분포의 분위수를 좌표평면에서 각각 수평축과 수직축의 좌표로 하여 점을 찍은 것
- 정규분포에 가까우면 직선의 모양



가

.

가

.

<-

.

(P.177)