

Chapter 6 Cross Table

- Contingency table, joint, marginal, conditional probabilities

교차표

- π_{ij} : (X, Y) 결합밀도함수
- π_{i+} : (X) 주변밀도함수

Homogeneity (동질성)

- 각 행에 대해 열의 분포가 동일한가?

$$H_0 : \pi_{ij} = \pi_{kj} \text{ for } j = 1, 2, \dots, c \text{ and } k \neq i$$

Independence (독립성)

- (X, Y)는 서로 독립인가? $H_0 : \pi_{ij} = \pi_{i+} \pi_{+j}$

검정통계량

$$\sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i} \sim \chi^2(k-1)$$

\ y	1	2	...	c	Total
x	π_{11}	π_{12}	...	π_{1c}	π_{1+}
2	π_{21}	π_{22}	...	π_{2c}	π_{2+}
...
R	π_{r1}	π_{r2}	...	π_{rc}	π_{r+}
Total	π_{+1}	π_{+2}	...	π_{+c}	π_{++}

Chi-Square Test

1/4

Example: Major vs. Department

Major	Oil	Chemical	Electrical	Computer
Business	30	15	15	40
Engineering	30	30	20	20

```
ds=c(30, 30, 15, 30, 15, 20, 40, 20)
rn=c("B", "E")
cn=c("O", "C", "E", "Com")
CT=matrix(ds, nrow=2, ncol=4, dimnames=list(rn,cn))
```



```
> CT
  O C E Com
B 30 15 15 40
E 30 30 20 20

> prop.table(CT, 1) # 2=열%, 1=행%, none=셀%
  O C E Com
B 0.3 0.15 0.15 0.4
E 0.3 0.30 0.20 0.2
```

```
prop.table(CT, 1) # 2=열%, 1=행%, none=셀%
chisq.test(CT)
```

```
data: CT
X-squared = 12.381, df = 3, p-value = 0.006186
```

o j, Yi† 0.62%† ð; Ô|ö (·ÓÕ ·WÞ × x Ã ÷ x) }{
o x "[·Óu ÷ ^ WÞ, Ó ·Óu Š, X WÞ ·† é, (- Ê •)



in Wñ

Data on the marital status of men and women ages 20 to 29 were obtained as part of a national survey. The results from a sample of 350 men and 400 women follow. These data are representative of results published in the *U.S. Current Population Report (The Statistical Abstract of the United States, 1999)*.

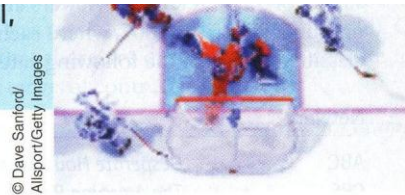
Gender	Marital Status		
	Never Married	Married	Divorced
Men	234	106	10
Women	216	168	16

- Use $\alpha = .01$ and test for independence between marital status and gender. What is your conclusion?
- Summarize the percent in each marital status category for men and for women.



CASE 15.1

Predicting the Outcomes of Basketball, Baseball, Football, and Hockey Games from Intermediate Results*



Some basketball fans generally believe that it doesn't pay to watch an entire game because the outcome is determined in the last few minutes (some say the last 2 minutes) of the game. Is this really true, and, if so, is basketball different in this respect from other professional sports played in North America? For example, is it true that the team that leads a baseball game after seven innings almost always wins the game? To address these questions, three researchers tracked basketball, baseball, football, and hockey games. The results (whether the early-game leader and whether the late-game leader won) of games during the 1990 season (for baseball and football) and during the 1990–1991 season (for basketball and

hockey) were recorded. Early-game leaders are defined as the teams that are ahead after one quarter of basketball and football, one period of hockey, or three innings of baseball. Late-game leaders are defined as the teams that are ahead after three quarters of basketball and football, two periods of hockey, or seven innings of baseball.

The data were recorded in the following way:

- Column 1: Results of games where 2 = early-game leader wins and 1 = early-game leader loses
- Column 2: Early-game leader game where 1 = basketball, 2 = baseball, 3 = football, and 4 = hockey games
- Column 3: Results of games where 2 = late-game leader wins and 1 = late-game leader loses
- Column 4: Late-game leader game where 1 = basketball, 2 = baseball, 3 = football, and 4 = hockey games

Column 3: Results of games where 2 = late-game leader wins and 1 = late-game leader loses

Column 4: Late-game leader game where 1 = basketball, 2 = baseball, 3 = football, and 4 = hockey games

Can we infer from these data that all four professional sports experience the same proportion of early-game leaders winning the game?

Can we infer from these data that all four professional sports experience the same proportion of late-game leaders winning the game?

```
table.sp=table(ds.sp$Game..Late,ds.sp$Late) #교차표작성
```



교차표 그래프 표현 (Simpson's Paradox)

(예제 데이터 중심) Google : Berkeley's Gender discrimination

여학생들이 주장했다. 남학생 지원자 8,442 명 중 44% (2,691 명)가 합격, 4,321 명 중 35% (1,835)가 합격 => 이는 성 차별이라고 주장

(학과별로 작성된 표)

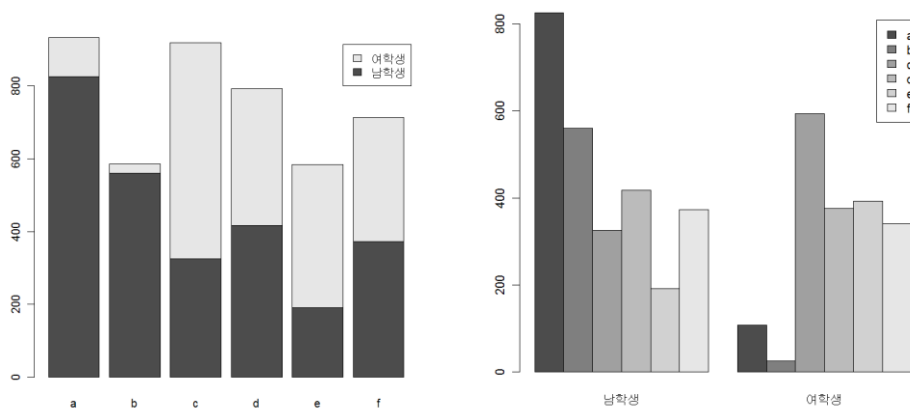
```
> data=c(825,108,560,25,325,593,417,375,191,393,373,341)
> cn=c("a","b","c","d","e","f")
> rn=c("남학생","여학생")
> table=matrix(data,nrow=2,ncol=6,dimnames=list(rn,cn))
> table
      a  b  c  d  e  f
남학생 825 560 325 417 191 373
여학생 108  25 593 375 393 341
> chisq.test(table)
```

Pearson's Chi-squared test

```
data: table
X-squared = 1068.372, df = 5, p-value < 2.2e-16
```

```
> prop.table(table,1)
      a          b          c          d          e          f
남학생 0.30657748 0.20810108 0.1207729 0.1549610 0.07097733 0.1386102
여학생 0.05885559 0.01362398 0.3231608 0.2043597 0.21416894 0.1858311
> prop.table(table,2)
      a          b          c          d          e          f
남학생 0.8842444 0.95726496 0.3540305 0.5265152 0.3270548 0.522409
여학생 0.1157556 0.04273504 0.6459695 0.4734848 0.6729452 0.477591
```

```
barplot(table, legend=rownames(table))
barplot(t(table), beside=T, legend=colnames(table))
```



Department	Men		Women	
	Applicants	Admitted	Applicants	Admitted
A	825	62%	108	82%
B	560	63%	25	68%
C	325	37%	593	34%
D	417	33%	375	35%
E	191	28%	393	24%
F	272	6%	341	7%

빈도 분할

교차표 셀 빈도

$$f_{ij} = \mu + R_i + C_j + e_{ij}, \quad e_{ij} \sim N(0, \sigma^2) \quad \text{가법모형}$$

추정치

- $\hat{\mu} = \bar{y}_{..}$, $\hat{R}_i = \bar{y}_{i.} - \bar{y}_{..}$, $\hat{C}_j = \bar{y}_{.j} - \bar{y}_{..}$
- 위의 추정치는 이상값에 영향을 받으므로 적절한 추정치가 아님

Median Polish

가법모형 Robust 추정방법

단계 1 : 행 효과 (R_i)

- \tilde{y}_i : i 행의 관측치 중앙값을 행 효과의 추정치로 사용
- $y_{ij}^{(1)} = y_{ij} - \tilde{y}_i$
- $y_{ij}^{(1)}$ 에는 $\mu + C_j$ 만 남아 있음

단계 2 : 열 효과 (C_j)

- $\tilde{y}_j^{(1)}$: j 열의 관측치 중앙값을 열 효과의 추정치로 사용
- $y_{ij}^{(2)} = y_{ij}^{(1)} - \tilde{y}_j^{(1)}$



모형 적합성 검토

- 가법모형 vs. 승법모형 $f_{ij} = \mu * R_i * C_j + e_{ij}$
- 두 모형의 잔차 관계 : $e_{ij} = \frac{\text{행효과}(R_i) * \text{열효과}(C_j)}{\text{공통효과}(\mu)} + e_{ij}$
- 비교 comparison = $\frac{\text{행효과}(R_i) * \text{열효과}(C_j)}{\text{공통효과}(\mu)}$
- 승법모형과 비교의 산점도가 패턴이 있다면 승법모형은 적합하지 않다. 왜냐하면, 만약 승법모형이 적합하다면 오차항만 남아 있어야 하는데...
- 산점도의 기울기가 양의 기울기 (+1)를 가지면 승법모형이 적합 => 승법모형이 적합한 경우에는 데이터를 로그 변환하여 승법모형을 적합시키면 된다.



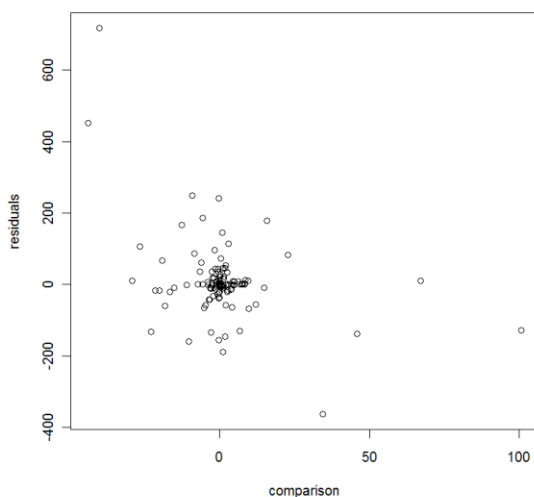
2011 대도시 강우량 데이터

```
rainfall=read.csv("rainfall.csv")
data=as.matrix(data.frame(rainfall[,2:11]))

cn=c("서울", "인천", "강릉", "대전", "전주", "광주", "제주", "부산", "울산", "대구")
rn=c("1월", "2월", "3월", "4월", "5월", "6월", "7월", "8월", "9월", "10월", "11월")
table.rainfall=matrix(data,nrow=12,ncol=10,dimnames=list(rn,cn))

ct.fit<-medpolish(table.rainfall)
attach(ct.fit)
comparison=matrix(row,ncol=1)%*%matrix(col,nrow=1)/overall
plot(residuals~comparison, main="2011년 월 강우량")
```

2011년 월 강우량



(일정한 패턴을 보이지 않음)



Overall: 80.01875

Row Effects:

1월	2월	3월	4월	5월	6월	7월
-73.59375	-21.06875	-56.86875	21.06875	40.81875	227.88125	342.61875
8월	9월	10월	11월	12월		
156.29375	-31.19375	-37.81875	28.40625	-65.06875		

Column Effects:

서울	인천	강릉	대전	전주	광주	제주
-9.37500	-10.27500	23.53750	-2.93750	-0.10000	0.40000	8.03750
부산	울산	대구				
3.44375	-5.34375	-0.96875				

- o 행 효과 (월별) 7 월에 비가 많이 오고 1 월 강우량이 가장 낮음
- o 열 효과 (도시) 강릉에 비가 많이 오고 서울이 가장 낮음
- o (잔차) 대전 지역 7,8 월에 비가 많이 오음, 제주는 7 월 비가 적게 오음

Residuals:

	서울	인천	강릉	대전	전주	광주	제주
1월	11.850	10.850	-16.9625	0.5125	-1.8250	0.5750	1.1375
2월	-20.475	-16.175	61.5125	-11.2125	-0.4500	0.4500	-31.8875
3월	0.825	1.625	-20.1875	-1.2125	0.4500	-0.4500	1.2125
4월	18.387	36.788	8.3750	-27.1500	-20.1875	1.5125	-57.8250
5월	-58.063	-66.362	-56.3750	44.1000	2.0625	21.6625	-62.7750
6월	105.975	9.975	11.5625	86.6375	-156.0000	-188.3000	82.6625
7월	717.737	451.838	-127.9750	167.6000	44.4625	-145.5375	-362.3750
8월	-60.137	-17.637	-138.6500	186.9250	241.5875	145.7875	178.7500
9월	-13.850	-12.150	248.5375	45.8125	22.4750	-35.7250	-9.0625
10월	-0.825	-1.625	-0.9375	-2.2625	-7.5000	-22.1000	-1.1375
11월	-42.850	-42.450	0.9375	-2.2875	10.5750	27.9750	114.8375
12월	1.525	2.425	66.5125	-0.5125	-4.3500	-2.0500	36.9125



(1990~2011 년) 소비지출

가계수지항목별	2004			2005	
	전가구	근로자가구	근로자외가구	전가구	근로자가구
□ 소비지출 (원)	1,963,316	2,005,758	1,906,569	2,035,256	2,035,256
▣ 식료품 (원)	532,452	546,711	513,309	539,260	539,260
▣ 주거 (원)	64,627	63,675	65,902	69,498	69,498
▣ 광열·수도 (원)	97,477	94,888	100,951	102,561	102,561
▣ 가구집기가사용 (원)	77,928	81,004	73,845	84,950	84,950
▣ 의류및신발 (원)	101,523	104,133	98,047	106,698	106,698
▣ 보건의료 (원)	93,705	91,552	96,609	102,446	102,446
▣ 교육 (원)	219,833	220,422	218,983	229,747	229,747
▣ 교양오락 (원)	93,918	98,053	88,411	98,275	98,275
▣ 교통·통신 (원)	339,047	356,895	315,078	356,166	356,166
▣ 기타소비지출 (원)	342,806	348,427	335,437	345,656	345,656

