

- From Agresti (2013, 3rd ed.)

Table 2.11 Data for Exercise 2.15 on Graduate Admissions

3-way contingency table

check table in LNp.2-4

Department	Whether Admitted			
	Male		Female	
	Yes	No	Yes	No
A	512	313	89	19
B	353	207	17	8
C	120	205	202	391
D	138	279	131	244
E	53	138	94	299
F	22	351	24	317
Total	1198	1493	557	1278

Fixed? random?

resp. (Admitted) expl. (Gender) (Department)

	Yes	No	M	F	A	F
1st	Yes	No	M	F	A	F
2nd	Yes	No	F	F	A	F
⋮	⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	M	F	F	F
⋮	⋮	⋮	F	F	F	F

sufficient statistics

count resp. (Gender) (Department) (Admitted) expl. (X)

	Yes	No	M	F	A	F
512 = N_{1A1}	Yes	No	M	F	A	F
313 = N_{1A2}	Yes	No	M	F	A	F
89 = N_{2A1}	Yes	No	F	F	A	F
19 = N_{2A2}	Yes	No	F	F	A	F
⋮	⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮	⋮

Note. $P_{G,D,1} = \lambda_{G,D,1}$
 $\lambda_{G,D,1} + \lambda_{G,D,2}$

$N_{G,D,A} \sim \text{Poisson}(\lambda_{G,D,A})$ parameter

Q: how (G,D,A) influence $\lambda_{G,D,A}$'s?
 or how (G,D) influence

binomial? Poisson (LNp.2-7~8)

- From Agresti (2013, 3rd ed.)

Table 2.5 Cross-Classification of Smoking by Lung Cancer

retrospective study

random

Smoker	Lung Cancer	
	Cases	Controls
Yes	688	650
No	21	59
Total	709	709

Fixed

- In data collection, collect information of $X|Y$

resp. (smoker) expl. (cancer) Suff. Stat. count resp. expl.

	Yes	No	Yes	No
1st	Yes	No	Yes	No
2nd	No	Yes	Yes	No
⋮	⋮	⋮	⋮	⋮
⋮	Yes	No	Yes	No
⋮	No	No	Yes	No

Yes(1)
No(2)

statistical modeling:

$$N_{11} (Y=1) \sim \text{binomial}(709, P_1) \rightarrow P(X=\text{Yes}|Y=1)$$

$$N_{21} (Y=2) \sim \text{binomial}(709, P_2) \rightarrow P(X=\text{Yes}|Y=2)$$

→ can study problem related to $X|Y$

pro: X is fixed in advance, like in DOE
 retro: Y (categorical response) is fixed.
 Under what circumstances, we can get useful information about $Y|X$ from $X|Y$?
 But, in data analysis, interested in problems related to $Y|X$, e.g., how smoking ($X=\text{expl}$) influence the probability of getting cancer ($Y=\text{resp}$)

Table 3.9 Data for Fisher's Tea-Tasting Experiment

Poured First (PE)	Guess Poured First		Total
	Milk	Tea	
Milk	3	1	4
Tea	1	3	4
Total	4	4	8

fixed

resp. (Guess) expl. (PF) Suff. Statistics

	M	T	M	T
1st	M	T	M	T
2nd	T	M	M	T
⋮	⋮	⋮	⋮	⋮
7th	M	T	M	T
8th	T	T	M	T

always, 4 M, 4 T

count resp. expl. (Guess) (PF)

N_{11}, N_{12} M (1)
 N_{21}, N_{22} T (2)

statistical modeling under random guess:

$$N_{11} \sim \text{hypergeometric}(4, 4, 4)$$

of balls drawn (claim M)

of black balls drawn

of white balls (T)

of black balls (M)

If N_{11} is known, the rest counts are known.

• From Agresti (2002, 2nd ed.)

TABLE 7.1 Primary Food Choice of Alligators

Lake	Gender	Size (m)	Primary Food Choice								
			Fish	Invertebrate	Reptile	Bird	Other				
Hancock	Male	≤ 2.3 1	7	+	1	+	0	+	0	+	5
		> 2.3 2	4		0		0		1		2
	Female	≤ 2.3 3	16		3		2		2		3
		> 2.3 4	3		0		1		2		3
Oklawaha	Male	≤ 2.3 5	2		2		0		0		1
		> 2.3 6	13		7		6		0		0
	Female	≤ 2.3 7	3		9		1		0		2
		> 2.3 8	0		1		0		1		0
Trafford	Male	≤ 2.3 9	3		7		1		0		1
		> 2.3 10	8		6		6		3		5
	Female	≤ 2.3 11	2		4		1		1		4
		> 2.3 12	0		1		0		0		0
George	Male	≤ 2.3 13	13		10		0		2		2
		> 2.3 14	9		0		0		1		2
	Female	≤ 2.3 15	3		9		1		0		1
		> 2.3 16	8	+	1	+	0	+	0	+	1

Total (K) fixed ②?
or random ③?

rowsum(K_x)
fixed ①? or
random ②?

Information: ① < ② < ③

count resp. (X) (Food)
 $N_{x_1, F}$ x_1 F

$N_{x_1, 0}$ x_1 0

$N_{x, Food} \sim \text{Poisson}(\lambda_{x, Food})$

Q: how $(X, Food)$ influence $\lambda_{x, Food}$'s?

② $N_{x, Food} \sim \text{multinomial}$

80×1 vector $(K, P_{x, Food})$
 $\sum_x \sum_{Food} P_{x, Food} = 1$

- data collection: if data from survey sampling, X & Y are random
- data analysis: in regression-type analysis, Y : random; X : fixed (conditioned)

Suff. Stat.

N → count resp. (Food)

nominal

expl. (X)
(Lake)(Gender)(Size)

$N_{x_1} = (N_{1F}, N_{1I}, N_{1R}, N_{1B}, N_{1O})$
 $H \quad M \quad \leq \quad X_1$

① $N_x \sim \text{multinomial}(K_x, P_x \equiv (P_{xF}, P_{xI}, P_{xR}, P_{xB}, P_{xO}))$

5 × 1 vector $P_x = P_{xF} + P_{xI} + P_{xR} + P_{xB} + P_{xO} = 1$ (X fixed)
Q: how $X, Food$ influence $P_{x, Food}$?

Y → resp. (Food) (Lake) (Gender) (Size)
1st F R nominal
2nd R G nominal
2 levels
2 levels

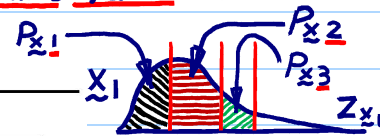
• From Agresti (2002, 2nd ed.)

TABLE 7.7 Life-Length Distribution of U.S. Residents (Percent),^a 1981

Life Length	Gender			
	Males		Females	
	White ← race → Black		White	Black
R1 0-20	2.4 (2.4)	3.6 (4.4)	1.6 (1.2)	2.7 (2.3)
R2 20-40	3.4 (3.5)	7.5 (6.4)	1.4 (1.9)	2.9 (3.4)
R3 40-50	3.8 (4.4)	8.3 (7.7)	2.2 (2.4)	4.4 (4.3)
R4 50-60	17.5 (16.7)	25.0 (26.1)	9.9 (9.6)	16.3 (16.3)
R5 Over 65	72.9 (73.0)	55.6 (55.4)	84.9 (84.9)	73.7 (73.7)

assume column sums are fixed. Q: What if they are random? check LNp2-12

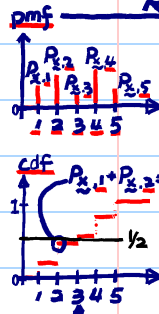
Z_x : continuous lifetime distribution



Y → resp. (Life) (Gender) (Race)
1st R3 disc. inte. M 2 levels W 2 levels
2nd R5 F 2 levels W 2 levels
R4 F B
R1 M B

Suff. Stat. → N → count resp. (Life) disc. inte. expl. (X) (Gender)(Race)

R_1, R_2, R_3, R_4, R_5
 $N_{x_1} = (N_{11}, N_{12}, N_{13}, N_{14}, N_{15})$
 $M \quad W \leftarrow X_1$
 $F \quad W \leftarrow X_2$
 $M \quad B \leftarrow X_3$
 $F \quad B \leftarrow X_4$



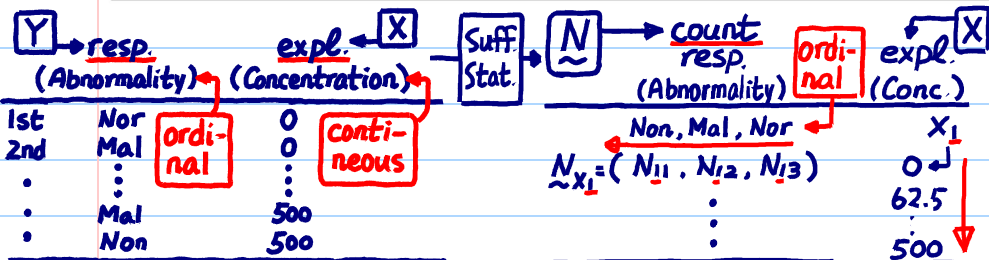
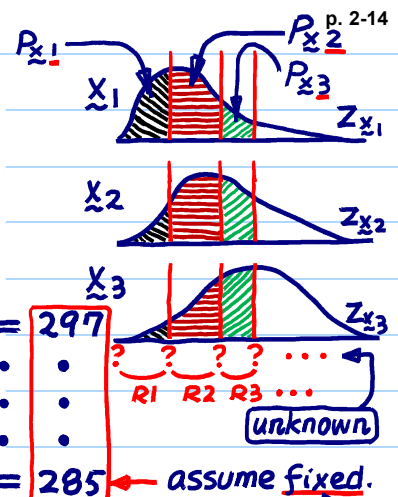
statistical modeling:
 $N_x \sim \text{multinomial}(K_x, P_x \equiv (P_{x1}, P_{x2}, P_{x3}, P_{x4}, P_{x5}))$

has some structure ∴ discrete interval categories =
Q: how X influence P_x (pmf) or cdf of P_x ?
continuous, ordinal, nominal
e.g. median is meaningful

- From Agresti (2013, 3rd ed.)

Table 8.7 Outcomes for Pregnant Mice in Developmental Toxicity Study

Concentration (mg/kg per day)	Response (abnormality)				
	Nonlive		Malformation		Normal
0 (controls)	15	+	1	+	281
62.5	17	.	0	.	225
125	22	.	7	.	283
250	38	.	59	.	202
500	144	+	132	+	9



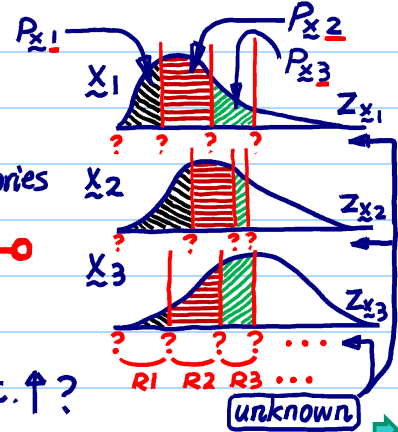
Q: What if they're random? Check LNp.2-12

- statistical modeling: has some structure \therefore ordinal categories

$$\underline{N}_X \sim \text{multinomial}(K_X, \underline{P}_X \equiv (P_{X1}, P_{X2}, P_{X3}))$$

Q: how X influence \underline{P}_X or cdf of \underline{P}_X ?

say, whether abnormality \uparrow when conc. \uparrow ?

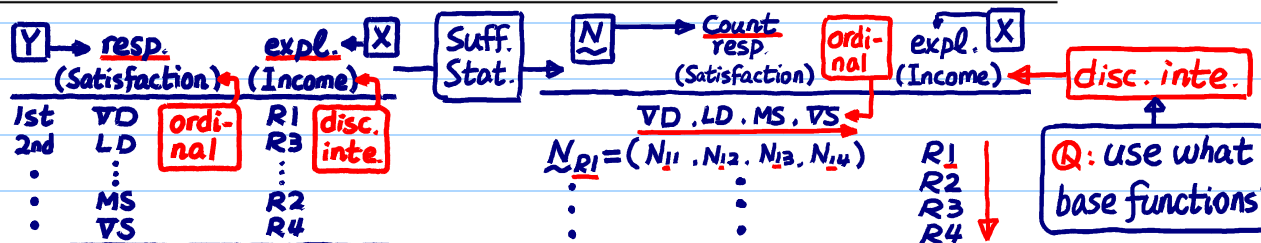


- From Agresti (2002, 2nd ed.)

TABLE 2.8 Cross-Classification of Job Satisfaction by Income

Income (dollars)	Job Satisfaction			
	Very Dissatisfied	Little Dissatisfied	Moderately Satisfied	Very Satisfied
R1 < 15,000	1	+	3	+
R2 15,000-25,000	2	.	3	.
R3 25,000-40,000	1	.	6	.
R4 > 40,000	0	+	1	+

Q: What if they're random? Check LNp.2-12



- statistical modeling:

$$\underline{N}_X \sim \text{multinomial}(K_X, \underline{P}_X \equiv (P_{X1}, P_{X2}, P_{X3}, P_{X4}))$$

Q: how X influence \underline{P}_X (pmf) or cdf of \underline{P}_X ?

say, whether Satisfaction \uparrow when Income \uparrow ?

