

Models for Matched Pairs

(Models for Square Tables)

Situation: Categorical analogue to dependent samples tests/models for continuous data.

“**Matched pairs**” are two samples that are statistically dependent.

1. The two samples that have the same respondents/individuals.
e.g.,
 - Same individuals respond to 2 questions.
 - Same individuals respond to 1 question at two time points.
“Panel data”.
2. The two samples are matched, a natural pairing. e.g.,
 - Ask a husband and wife the same question.
 - Have 2 people rate the same object/individual.
 - Education or occupation of parent and child.
“Mobility tables”.

Frequently matched pairs data yield Square Tables.

“**Square Tables**” are ones in which the row and column classifications (categories) are the same.

Examples of Matched Pairs/Square Tables

Redelmeier, D.A. & Tibshirani, R.J. (1997). Is using a car phone like driving drunk? *Chance*, 10, 5–9.

Data from case-crossover study using non-injury car accident (collision) data from Toronto (July 1994 – August 1995). Each individual acts as their own control.

Cellular telephone call

		during the control interval		
		yes	no	total
during the hazard interval	yes	13	157	170
	no	24	505	529
total		37	662	699

Hazard interval is the 10 minute period prior to accident.

Control interval is the 10 minute period at the same time of the accident but on the day before.

Note:

$$170/699 = .24 \quad \text{versus} \quad 37/699 = .05$$

In honor of Roger Ebert's Overlooked film festival ...

Agresti & Winner (1997) Evaluating agreement and disagreement among movie reviewers. *Chance*, 10, 10–14.

Data are from April 1995 through September 1996.

		Ebert			
		Con	Mixed	Pro	
Siskel	Con	24	8	13	45
	Mixed	8	13	11	32
	Pro	10	9	64	83
		42	30	88	160

Question: Do Siskel and Ebert really disagree?

A Classic Data Set: the Coleman (1964) Panel data

Responses to two items made by 3398 boys:

- Attitude toward the leading crowd (positive, negative).
- Self-perception of membership in the leading crowd (yes, no).

Question: Are they measuring the same thing?

		Attitude	
		positive	negative
Membership	yes	757	496
	no	1071	1074

The responses were actually collected at 2 time points. The above responses are from time point 1. We could look at consistency of responses over time or whether the marginal distributions changed or not:

Attitude over time (“panel data”)

		Time 2	
		positive	negative
Time 1	positive	1283	545
	negative	650	920

Question: Change in attitude over time?

Note: For the Coleman data, it is best to express them as a 4-way table

Membership at		Attitude time 1			
		positive		negative	
		Attitude time 1		Attitude time 2	
time 1	time 2	positive	negative	positive	negative
yes	yes	458	140	171	182
	no	110	49	56	87
no	yes	184	75	85	97
	no	531	281	338	554

A good fitting model is homogeneous association loglinear model.

What do you suppose is the weakest association(s)?

What do you suppose is the strongest association(s)?

A mobility table where a single observational unit is measured at 2 time points.

Migration data comparing the region of residence in the U.S. in 1980 with 1985. These data are from the US Bureau of the Census (Agresti, 1990).

Residence in 1980	Residence in 1985				Total
	NorthEast	Midwest	South	West	
NorthEast	11,607	100	366	124	12,197
Midwest	87	13,677	515	302	14,581
South	172	225	17,819	270	18,486
West	63	176	286	10,192	10,717
Total	11,929	14,178	18,986	10,888	55,981

Many of the models for matched pairs data use methodology for structurally incomplete tables.

Type of analyses/models that may be warranted when studying matched pairs data:

1. Compare the margins of the table (dependent proportions).
 - Marginal homogeneity.
 - McNemar's test.
 - Estimating the difference between proportions (& confidence intervals for them).
2. For binary responses, logistic regression (for matched pairs).
 - McNemar's test.
 - Logit model with subject specific effects.
3. Loglinear models for comparing margins.
 - Conditional likelihood ratio test (symmetry minus quasi-symmetry).
4. Measuring agreement between 2 judges/observers who rate common set of stimuli/subjects/individuals.
 - Quasi-independence.
 - Cohen's Kappa
5. Evaluating preferences between pairs of treatments.
 - BTL model

We'll talk 1, 3 and 4.

Comparing Dependent Proportions

Cellular telephone call

		during the control interval		
		yes	no	total
during the hazard interval	yes	13	157	170
	no	24	505	529
total		37	662	699

Question: Is the probability of a car accident larger when the driver uses a cell-phone than when the driver is not using a phone?

This is the same as asking whether the margins of the table the same?

Compare: $170/699 = .24$ and $37/699 = .05$

Difference in proportions: $.24 - .05 = .19$

Problem: These proportions are statistically dependent.

Solution.....

Cell counts (observed data):

n_{11}	n_{12}	n_{1+}
n_{21}	n_{22}	n_{2+}
n_{+1}	n_{+2}	n_{++}

In terms of probabilities,

π_{11}	π_{12}	π_{1+}
π_{21}	π_{22}	π_{2+}
π_{+1}	π_{+2}	π_{++}

Want to know whether

$$\pi_{1+} - \pi_{+1} = 0$$

Note:

$$\begin{aligned}\pi_{1+} - \pi_{+1} &= (\pi_{11} + \pi_{12}) - (\pi_{11} + \pi_{21}) \\ &= \pi_{12} - \pi_{21}\end{aligned}$$

or

- “Marginal Homogeneity”
- “Symmetry” across the “main diagonal”.

McNemar's Test for (2×2) tables

$H_o : \pi_{1+} = \pi_{+1}$ or equivalently $\pi_{12} - \pi_{21}$

Define $n^* = n_{12} + n_{21}$.

Consider the binomial variable with n^* trials that has its two possible outcomes cells (1,2) and (2,1) in the (2×2) table.

If H_o is true, then

- Expect $n_{12} \sim n_{21}$.
- Probability of cell (1,2) equals .5, and
Probability of cell (2,1) equals .5.

For “small” n^* , just compute the exact probability (p -value).

For “large” n^* ($n^* > 10$), use the normal approximation:

$$z = \frac{n_{12} - .5n^*}{\sqrt{n^*(.5)(.5)}} = \frac{n_{12} - n_{21}}{\sqrt{n_{12} + n_{21}}}$$

where

- $.5n^*$ = the expected count (mean) for the (1,2) cell if H_o is true.
- $n^*(.5)(.5)$ = the variance of the count.

Compare z to the standard normal distribution, or z^2 to the chi-square distribution with $df = 1$.

Example: Cell phones

H_o : marginal homogeneity or $\pi_{12} = \pi_{21}$

H_a : $\pi_{12} > \pi_{21}$

Test statistic:

$$z = \frac{157 - 24}{\sqrt{157 + 24}} = 9.89$$

p -value $\ll .001$.

Estimated difference of proportions:

$$p_{1+} - p_{+1} = \frac{170}{699} - \frac{37}{699} = .24 - .05 = .19$$

To form a confidence interval for $\pi_{1+} - \pi_{+1}$, we need the standard error of the estimated difference of the proportions.

The estimated variance of the difference,

$$\frac{p_{1+}(1 - p_{1+})}{n_{++}} + \frac{p_{+1}(1 - p_{+1})}{n_{++}} - 2 \frac{(p_{11}p_{22} - p_{12}p_{21})}{n_{++}}$$

For a $(1 - \alpha)100\%$ confidence interval for $\pi_{1+} - \pi_{+1}$,

$$(p_{1+} - p_{+1}) \pm z_{\alpha} \sqrt{\frac{p_{1+}(1 - p_{1+}) + p_{+1}(1 - p_{+1}) - 2(p_{11}p_{22} - p_{12}p_{21})}{n_{++}}}$$

In our example, the estimated variance of the difference is

$$\frac{.24(1 - .24) + .05(1 - .05) - 2 \left(\frac{13(505) - 24(157)}{699} \right)}{699} = .0003$$

and the standard error is

$$\sqrt{.0003} = .0683$$

and 95% confidence interval for $(\pi_{1+} - \pi_{+1})$ is

$$.19 \pm 1.96(.0683) \implies (.056, .324)$$

Notes regarding the study:

- The authors varied their choice of control interval, and arrive at the same general conclusion.
- The risk of a collision:

$$\frac{157}{24} = 6.5$$

Drivers have a 6.5-fold increased risk of being in a collision when using a cell-phone compared to when they were not using a phone. Note: 95% CI for risk is (4.5,10.0).

- Comparison to drunk driving:

	Risk
Blood alcohol at legal limit	4
50% alcohol above legal limit	10

Is there anything mis-leading in the comparison of risk while driving drunk?

McNemar's test using SAS:

```
DATA phones;
INPUT hazard $ control $ count @@;
DATALINES;
  yes yes 13    yes no 157
    no yes 24    no no 505 ;
PROC FREQ;
WEIGHT count;
TABLES hazard*control / CHISQ AGREE;
```

The **AGREE** option gives you McNemar's test if you have a (2×2) table. It also gives other stuff.

Loglinear Models for Square Tables

Residence in 1980	Residence in 1985				Total
	NorthEast	Midwest	South	West	
NorthEast	11,607	100	366	124	12,197
Midwest	87	13,677	515	302	14,581
South	172	225	17,819	270	18,486
West	63	176	286	10,192	10,717
Total	11,929	14,178	18,986	10,888	55,981

Note: Relatively few people changed regions —
95% of the observations fell on the main diagonal.

Test of independence:

$$df = (4 - 1)(4 - 1) = 9, G^2 = 125,923 \text{ and } p < .00001.$$

If we disregard the diagonal, does independence hold for the off diagonal cells?

We would represent such a structure by the following loglinear model, which is known as the

“Quasi-Independence Loglinear Model”

$$\log(\mu_{ij}) = \begin{cases} \lambda + \lambda_i^{1980} + \lambda_j^{1985} & \text{for } i \neq j \\ n_{ij} & \text{for } i = j \end{cases}$$

or using indicator variables,

$$\log(\mu_{ij}) = \lambda + \lambda_i^{1980} + \lambda_j^{1985} + \delta_i I(i = j)$$

where

$$I(i = j) = \begin{cases} 1 & \text{for } i = j \\ 0 & \text{otherwise} \end{cases}$$

So there are 5 indicator variables, which we treat as numerical variables.

The δ_i parameters ensure that the diagonal cells are fit perfectly.

$$\begin{aligned} df &= (\text{usual } df) - (\# \text{ diagonals fit perfectly}) \\ &= (I - 1)(I - 1) - I \end{aligned}$$

Migration Example:

H_o : For people who moved (the movers), residence in 1985 is independent of region of residence in 1980.

$$df = 5, G^2 = 69.51, p < .001.$$

The quasi independence model fits much better than independence, primarily because the diagonals are fit perfectly (and this is where most of the observations are).

Quasi independence still is missing something in the data.

Fitting quasi-independence as a generalized linear models using SAS/GENMOD.

You need to create (or enter) a separate indicator (dummy) variable for each diagonal cell. The indicator is treated as a numerical variable in the model.

```
DATA migrate;
  INPUT y1980 $ y1985 $ count;
  ine=0; if y1980='NorthEast' AND y1980=1985 THEN
  ine=1;
  imw=0; if y1980='Midwest' AND y1980=1985 THEN imw=1;
  iso=0; if y1980='South' AND y1980=1985 THEN iso=1;
  iwe=0; if y1980='West' AND y1980=1985 THEN iwe=1;
DATALINES;
  NorthEast NorthEast 11607
  NorthEast Midwest    100
  :
  West          West    10192

PROC GENMOD ORDER=data;
  CLASS y1980 y1985;
  MODEL count = y1980 y1985 ine imw iso iwe
    / LINK=log DIST=Poisson;
```

Short-cut to fitting the quasi-independence model

You first create (or enter) a variable that takes on a unique value for each of the diagonal cells and a common value for all of the off diagonal cells.

For example,

$$qi = \begin{cases} 1 & i = j = 1 \\ 2 & i = j = 2 \\ 3 & i = j = 3 \\ 4 & i = j = 4 \\ 5 & i \neq j \end{cases}$$

This new variable is treated as a nominal/classification variable.

In *SAS* code:

```
DATA migrate;
  INPUT y1980 $ y1985 $ count;
  IF y1980='NorthEast' AND y1980=1985 THEN qi=1;
  ELSE IF y1980='Midwest' AND y1980=1985 THEN qi=2;
  ELSE IF y1980='South' AND y1980=1985 THEN qi=3;
  ELSE IF y1980='West' AND y1980=1985 THEN qi=4;
  ELSE qi=5;
DATALINES;
:
PROC GENMOD ORDER=data;
  CLASS y1980 y1985 qi ;
  MODEL count = y1980 y1985 qi / LINK=log DIST=poi;
```

Symmetry

This model states that

$$\mu_{ij} = \begin{cases} \mu_{ji} & \text{for } i \neq j \\ n_{ii} & \text{for } i = j \end{cases}$$

(i.e., disregard the diagonal).

This is only applicable to square tables.

Example of a perfectly symmetric table:

	1	2	3	total
1	100	20	40	160
2	20	100	30	150
3	40	30	100	170
total	160	150	170	

MLE of μ_{ij} from the symmetry model:

$$\hat{\mu}_{ij} = \hat{\mu}_{ji} = \frac{n_{ij} + n_{ji}}{2}$$

Degrees of freedom:

There are $I(I - 1)$ off diagonal cells.

$I(I - 1)/2$ parameters estimated (unique fitted values).

So,

$$\begin{aligned}df &= (\# \text{ cells}) - (\# \text{ non-redundent parameters}) \\&= \# \text{ off diagonal cells} - \# \text{ unique parameters} \\&= I(I - 1) - I(I - 1)/2 \\&= I(I - 1)/2\end{aligned}$$

The symmetry model can also be written as a loglinear model, which will help to show what the implications of the model are for the structure in the table.

$$\log(\mu_{ij}) = \lambda + \lambda_i + \lambda_j + \lambda_{ij}$$

where $\lambda_{ij} = \lambda_{ji}$.

There are no superscripts on the main/marginal effect terms because they are the same for the rows and columns, i.e.,

$$\lambda_i = \lambda_j \quad \text{when } i = j$$

In other words, the row and column margins are equal, i.e.,

$$\mu_{i+} = \mu_{+i}$$

The symmetry model is

1. A very restrictive, because it has implications for both the association between the variables and the margins of the table.

The symmetry model rarely fits data very well.

2. An important model because testing symmetry is often an important preliminary analysis for other analyses which require symmetric tables.

Example: Does the migration table exhibit symmetry?

Let's first just "look" at the table?

Residence in 1980	Residence in 1985				Total
	NorthEast	Midwest	South	West	
NorthEast	11,607	100	366	124	12,197
Midwest	87	13,677	515	302	14,581
South	172	225	17,819	270	18,486
West	63	176	286	10,192	10,717
Total	11,929	14,178	18,986	10,888	55,981

Symmetry Model fit statistics:

$$df = 4(4 - 1)/2 = 6, G^2 = 243.35, p < .001.$$

Doesn't support the symmetry hypothesis;
symmetry is too simple for these data.

Example 2: Siskell & Ebert data

		Ebert			
		Con	Mixed	Pro	
Siskell	Con	24	8	13	45
	Mixed	8	13	11	32
	Pro	10	9	64	83
		42	30	88	160

Summary of Models fit to the data:

Model	df	G^2	p -value
Independence	4	43.23	< .001
Quasi-independence	1	.01	.92
Symmetry	3	.59	.90

So what would you say about whether Siskell & Ebert really disagree?

Note: For symmetry, you can either

- Fit the symmetry model using PROC GENMOD, as described on the following pages, or
- Use PROC FREQ with the “AGREE” option on the TABLES command. For tables where $I > 2$, this will generate the df , G^2 and p -value for the symmetry model.

Fitting the symmetry model (obtaining parameter estimates of model written as a loglinear model):

There are (at least) two ways.

Method I (Agresti):

You need to create a variable that takes on a unique value for each diagonal cell and a unique value of each pair of cells.

For example,

$$\text{symm} = \left\{ \begin{array}{l} 1 \quad i = j = 1 \\ 2 \quad i = j = 2 \\ 3 \quad i = j = 3 \\ 4 \quad i = j = 4 \\ 5 \quad (i, j) = (1, 2) \text{ or } (2, 1) \\ 6 \quad (i, j) = (1, 3) \text{ or } (3, 1) \\ 7 \quad (i, j) = (1, 4) \text{ or } (4, 1) \\ 8 \quad (i, j) = (2, 3) \text{ or } (3, 2) \\ 9 \quad (i, j) = (2, 4) \text{ or } (4, 2) \\ 10 \quad (i, j) = (3, 4) \text{ or } (4, 3) \end{array} \right.$$

This new variable is treated as a nominal variable in fitting the model.

In *SAS/GENMOD*,

```
PROC GENMOD;  
  CLASS symm  
  MODEL count = symm / LINK=log DIST=poisson;
```

Method II.

This method involves a little “trick” and uses standard loglinear models.

The trick is rewriting the 2–way table (with frequencies n_{ij}) as a 3–way table (with frequencies n_{ijk}^*) as follows.

Create a new (conditioning) variable with 2 levels. Let’s call this variable Z and index it using k , then the entries of the 3–way table equal

$$n_{ijk}^* = \begin{cases} n_{ij} & \text{for } k=1 \\ n_{ji} & \text{for } k=2 \end{cases}$$

That is,

		$Z = 1$			
		Y			
X		1	2	...	I
1		n_{11}	n_{12}	...	n_{1I}
2		n_{21}	n_{22}	...	n_{2I}
\vdots			\vdots		
I		n_{I1}	n_{I2}	...	n_{II}

		$Z = 2$			
		Y			
		1	2	...	I
	1	n_{11}	n_{21}	...	n_{I1}
	2	n_{12}	n_{22}	...	n_{I2}
	\vdots		\vdots		\vdots
	I	n_{1I}	n_{2I}	...	n_{II}

The symmetry model corresponds to the joint independence loglinear model (XY, Z) .

To see why this works, take a table that exhibit perfect symmetry,

$$\{n_{ij}\} = \begin{pmatrix} 100 & 20 & 40 \\ 20 & 100 & 30 \\ 40 & 30 & 100 \end{pmatrix}$$

Then

$$\{n_{ij1}^*\} = \begin{pmatrix} 100 & 20 & 40 \\ 20 & 100 & 30 \\ 40 & 30 & 100 \end{pmatrix} \quad \text{and} \quad \{n_{ij2}^*\} = \begin{pmatrix} 100 & 20 & 40 \\ 20 & 100 & 30 \\ 40 & 30 & 100 \end{pmatrix}$$

or we can write is as Z crossed with XY

	$X = 1$			$X = 2$			$X = 3$		
Z	$Y = 1$	2	3	1	2	3	1	2	3
1	100	20	40	20	100	30	40	30	100
2	100	20	40	20	100	30	40	30	100

When XY is (jointly) independent of Z , then the 2-way table of X crossed with Y is symmetric.

Note that if you use this method, you need to adjust the fit statistics and degrees of freedom. The computer gives you G^2 for a 3-way table. Every cells is counted twice instead of just once, so to get the correct G^2 and df , just divide by 2.

Quasi-Symmetry

Since symmetry is so restrictive, we can remove the restriction that the margins must be equal (i.e., “marginal homogeneity”).

In other words, we will drop the requirement that the main/marginal effects for the two variables are equal.

The quasi-symmetric loglinear model

$$\log(\mu_{ij}) = \lambda + \lambda_i^{1980} + \lambda_j^{1985} + \lambda_{ij}$$

where $\lambda_{ij} = \lambda_{ji}$.

Degrees of freedom,

$$\begin{aligned} df &= (\# \text{ of cell}) - (\# \text{ nonredundant parameters}) \\ &= I^2 - [1 + (I - 1) + (I - 1) + I(I - 1)/2] \\ &= (I - 2)(I - 1)/2 \end{aligned}$$

Migration Example:

$$df = (4 - 2)(4 - 1)/2 = 2(3)/2 = 3,$$

$$G^2 = 2.99,$$

$$p = .39 \quad .$$

To fit the quasi-symmetry model, modify methods for fitting symmetry. The modification needed for

- I. (Agresti) Where you create a “symm” (symmetry) variable, add the row and column variable to the model.

In our example,

```
PROC GENMOD;  
  CLASS 1980 1985 symm;  
  MODEL count = 1980 1985 symm  
    /LINK=log DIST=poi;
```

- II. Where data are reformatted as a 3-way table, fit the homogeneous loglinear model.

Marginal Homogeneity

Are the row and column distributions (of a square table) the same?

The null hypothesis is

$$H_o : \mu_{i+} = \mu_{+i}$$

This is a simple hypothesis, but it difficult to test, because there is no way to use loglinear models to directly fit/test this model.

Options:

1. Do not use loglinear models.
2. Use generalized least squares instead of maximum likelihood estimation.
3. Indirectly test it using loglinear models (i.e., conditional likelihood ratio test).

We'll discuss (3): a contextual/comparison test.

Symmetry has two components:

$$\text{marginal homogeneity} \quad + \quad \text{quasi-symmetry}$$

Symmetry is a special case of quasi-symmetry.

If quasi-symmetry holds, the a test of marginal homogeneity is

$$G^2(\text{marginal homogeneity}) = G^2(\text{quasi symmetry}) - G^2(\text{symmetry})$$

with $df = I - 1$.

Mirgration example:

$$df = 4 - 1 = 3$$

$$G^2 = 240.56,$$

$$p < .001$$

Summary of analyses for the migration example:

Model	df	G^2	p
Independence	9	125,926.00	< .001
Quasi Independence	5	69.51	< .001
Symmetry	6	243.55	< .001
Marginal homogeneity	3	240.56	< .001
Quasi-symmetry	3	2.99	.393
Saturated loglinear	0	0.00	1.000

Residence in 1980	Residence in 1985				Total
	NorthEast	Midwest	South	West	
NorthEast	11,607	100	366	124	12,197
Midwest	87	13,677	515	302	14,581
South	172	225	17,819	270	18,486
West	63	176	286	10,192	10,717
Total	11,929	14,178	18,986	10,888	55,981

To get a better idea what quasi-symmetry means (and that the data are well described by this model. . .

The model:

$$\log(\mu_{ij}) = \lambda + \lambda_i^{1980} + \lambda_j^{1985} + \lambda_{ij}$$

or

$$\mu_{ij} = \exp[\lambda + \lambda_i^{1980} + \lambda_j^{1985} + \lambda_{ij}]$$

where $\lambda_{ij} = \lambda_{ji}$.

Re-arranging terms in the model,

$$\frac{\mu_{ij}}{\exp[\lambda + \lambda_i^{1980} + \lambda_j^{1985}]} = \exp[\lambda_{ij}]$$

Using our parameter estimates and data, let's compute

$$\frac{n_{ij}}{\exp[\hat{\lambda} + \hat{\lambda}_i^{1980} + \hat{\lambda}_j^{1985}]} \sim \text{symetric association}$$

Residence in 1980	Residence in 1985			
	NorthEast	Midwest	South	West
NorthEast	—	.809	1.027	.971
Midwest	.885	—	1.834	1.002
South	.944	1.905	—	1.034
West	1.055	.996	.970	—

The relationship among the models we've discussed for square tables:

- The most general (complex) model is quasi-symmetry.
- Symmetry is a special case of quasi-symmetry. i.e.,

$$\mu_{ij} = \lambda + \lambda_i^X + \lambda_j^Y + \lambda_{ij}$$

where $\lambda_{ij} = \lambda_{ji}$, and $\lambda_i^X = \lambda_i^Y$.

- Quasi-independence is a special case of quasi-symmetry.

Model of quasi symmetry is

$$\mu_{ij} = \lambda + \lambda_i^X + \lambda_j^Y + \lambda_{ij}$$

where $\lambda_{ij} = \lambda_{ji}$.

The model of quasi independence is

$$\mu_{ij} = \lambda + \lambda_i^X + \lambda_j^Y + \delta_i I(i = j)$$

where

$$I(i = j) = \begin{cases} 1 & \text{for } i = j \text{ (diagonals)} \\ 0 & \text{for } i \neq j \text{ (off diagonals)} \end{cases}$$

To see that quasi-independence is a special case of quasi-symmetry note that in the quasi-independence model, for the off diagonal cells,

$$\lambda_{ij} = \lambda_{ji} = 0$$

- Symmetry is **not** a special case of quasi-independence.
- Quasi-independence is **not** a special case of symmetry.