

SOC 3811/5811:  
BASIC SOCIAL STATISTICS

Associations Between Categorical Variables

# Associations Between Variables

Between now and the 3<sup>rd</sup> exam we will focus on measuring the association between two variables, X & Y

1. When X is discrete and Y is continuous, we will use “analysis of variance” techniques (Tuesday)
2. When X and Y are both discrete, we will use cross-tabular and  $\chi^2$  analyses (Today)
3. When X and Y are both continuous, we will use correlation & regression analyses (Next Week)

# Associations Between Discrete Variables

## How Important Are “Loose Morals and Drunkenness” in Explaining Poverty

Source: GSS

	<i>Very Important</i>	<i>Somewhat Important</i>	<i>Not Important</i>	Row Total
<i>Democrat</i>	164	161	131	456
<b>Political Party</b> <i>Independent</i>	179	139	100	418
<i>Republican</i>	173	154	103	430
Column Total	516	454	334	N=1,304

# Associations Between Discrete Variables

**Is there an association in the population?**

$\chi^2$  analysis

**How strong is the association? In what direction is the association?**

Gamma, relative risk, odds ratios

# Associations Between Discrete Variables

## How Important Are “Loose Morals and Drunkenness” in Explaining Poverty

Source: GSS

	<i>Very Important</i>	<i>Somewhat Important</i>	<i>Not Important</i>	Row Total
<i>Democrat</i>	164	161	131	456
<b>Political Party</b> <i>Independent</i>	179	139	100	418
<i>Republican</i>	173	154	103	430
Column Total	516	454	334	N=1,304

# The Logic of $\chi^2$ Tests

Even if there is no association between two discrete variables  $X$  and  $Y$  in the population, we may observe an association between  $X$  and  $Y$  in sample data because of random error or sampling variability

# The Logic of $\chi^2$ Tests

X and Y in a Complete Population

Y	1	2	3	Total
1	10,000	10,000	10,000	30,000
2	10,000	10,000	10,000	30,000
Total	20,000	20,000	20,000	60,000

X and Y in a Random Sample of 100 cases

Y	X			Total
	1	2	3	
1	20	19	18	57
2	13	13	17	43
Total	33	32	35	100

# The Logic of $\chi^2$ Tests

Even if there is no association between two discrete variables  $X$  and  $Y$  in the population, we may observe an association between  $X$  and  $Y$  in sample data because of random error or sampling variability

How can we tell whether the association observed between  $X$  and  $Y$  in sample data is strong enough to rule out the hypothesis that in the population  $X$  and  $Y$  are statistically independent (or not associated with one another)?



# The Logic of $\chi^2$ Tests

*We begin with the null hypothesis that there is no association between X and Y in the population ... that is, we assume “statistical independence”*

We then compute the cell frequencies that we would expect to observe under the null hypothesis and compare them to the actually observed cell frequencies

The  $\chi^2$  test statistic quantifies the degree to which the observed frequencies differ from the frequencies that we would expect to observe under the null hypothesis

# The Logic of $\chi^2$ Tests

Imagine (1) that 10% of people are left-handed, (2) that 50% of people are male and 50% are female, and (3) that there is no relationship between gender and whether someone is left-handed.

What would you expect in the cells of the table below if you sampled 1,000 people?

	Right	Left	<i>Total</i>
Male			
Female			
<i>Total</i>			1,000

# $\chi^2$ Tests

## Hypothesis Testing in 6 Steps ... Just Like Before

1. State the null ( $H_0$ ) and alternative ( $H_1$ ) hypotheses
2. Check that the sample data conform to basic assumptions; if they do not, then do not go any further
3. Choose an  $\alpha$  probability level ... that is, a probability associated with incorrectly rejecting the null hypothesis
4. Determine the “critical value” ... that is, how large the test statistic must be in order to reject the null hypothesis at the given  $\alpha$  level
5. Calculate the test statistic ...  $\chi^2$
6. Compare the test statistic to the critical value

# $\chi^2$ Tests

State the null ( $H_0$ ) and alternative ( $H_1$ ) hypotheses

We begin with the assumption that there is no association between X and Y

$H_0$ : X and Y are statistically independent  
(i.e., not associated in the population)

$H_1$ : X and Y are not statistically independent  
(i.e., are associated in the population)

# $\chi^2$ Tests

Check that the sample data conform to basic assumptions;  
if they do not, then do not go any further

X and Y must be collected from a random sample of individuals from the population

Standard  $\chi^2$  testing procedures should be used with extreme caution — or not at all — if any cell frequency is less than 5

# $\chi^2$ Tests

Choose an  $\alpha$  probability level ... that is, a probability associated with incorrectly rejecting the null hypothesis

Let's use  $\alpha=0.05$  in our example

# $\chi^2$ Tests

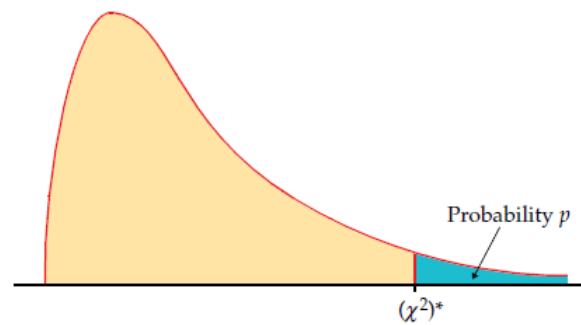
Determine the “critical value” ... that is, how large the test statistic must be in order to reject the null hypothesis at the given  $\alpha$  level

Critical values of  $\chi^2$  corresponding to common levels of  $\alpha$  and to specific numbers of degrees of freedom ( $df$ ) are given in a  $\chi^2$  table

For  $\chi^2$ ,  $df=(R-1)(C-1)$  where  $R$ =number of rows in the crosstable and  $C$ =number of columns in the crosstable

In our example,  $\alpha=0.05$  and  $df=(3-1)(3-1)=4$ , so the critical value of  $\chi^2$  equals 9.49

Table entry for  $p$  is the critical value  $(\chi^2)^*$  with probability  $p$  lying to its right.



**TABLE F**

$\chi^2$  distribution critical values

df	Tail probability $p$											
	.25	.20	.15	.10	.05	.025	.02	.01	.005	.0025	.001	.0005
1	1.32	1.64	2.07	2.71	3.84	5.02	5.41	6.63	7.88	9.14	10.83	12.12
2	2.77	3.22	3.79	4.61	5.99	7.38	7.82	9.21	10.60	11.98	13.82	15.20
3	4.11	4.64	5.32	6.25	7.81	9.35	9.84	11.34	12.84	14.32	16.27	17.73
4	5.39	5.99	6.74	7.78	9.49	11.14	11.67	13.28	14.86	16.42	18.47	20.00
5	6.63	7.29	8.12	9.24	11.07	12.83	13.39	15.09	16.75	18.39	20.51	22.11
6	7.84	8.56	9.45	10.64	12.59	14.45	15.03	16.81	18.55	20.25	22.46	24.10
7	9.04	9.80	10.75	12.02	14.07	16.01	16.62	18.48	20.28	22.04	24.32	26.02
8	10.22	11.03	12.03	13.36	15.51	17.53	18.17	20.09	21.95	23.77	26.12	27.87
9	11.39	12.24	13.29	14.68	16.92	19.02	19.68	21.67	23.59	25.46	27.88	29.67
10	12.55	13.44	14.53	15.99	18.31	20.48	21.16	23.21	25.19	27.11	29.59	31.42
11	13.70	14.63	15.77	17.28	19.68	21.92	22.62	24.72	26.76	28.73	31.26	33.14
12	14.85	15.81	16.99	18.55	21.03	23.34	24.05	26.22	28.30	30.32	32.91	34.82
13	15.98	16.98	18.20	19.81	22.36	24.74	25.47	27.69	29.82	31.88	34.53	36.48
14	17.12	18.15	19.41	21.06	23.68	26.12	26.87	29.14	31.32	33.43	36.12	38.11
15	18.25	19.31	20.60	22.31	25.00	27.49	28.26	30.58	32.80	34.95	37.70	39.72
16	19.37	20.47	21.79	23.54	26.30	28.85	29.63	32.00	34.27	36.46	39.25	41.31
17	20.49	21.61	22.98	24.77	27.59	30.19	31.00	33.41	35.72	37.95	40.79	42.88
18	21.60	22.76	24.16	25.99	28.87	31.53	32.35	34.81	37.16	39.42	42.31	44.43
19	22.72	23.90	25.33	27.20	30.14	32.85	33.69	36.19	38.58	40.88	43.82	45.97
20	23.83	25.04	26.50	28.41	31.41	34.17	35.02	37.57	40.00	42.34	45.31	47.50
21	24.93	26.17	27.66	29.62	32.67	35.48	36.34	38.93	41.40	43.78	46.80	49.01
22	26.04	27.30	28.82	30.81	33.92	36.78	37.66	40.29	42.80	45.20	48.27	50.51
23	27.14	28.43	29.98	32.01	35.17	38.08	38.97	41.64	44.18	46.62	49.73	52.00
24	28.24	29.55	31.13	33.20	36.42	39.36	40.27	42.98	45.56	48.03	51.18	53.48
25	29.34	30.68	32.28	34.38	37.65	40.65	41.57	44.31	46.93	49.44	52.62	54.95
26	30.43	31.79	33.43	35.56	38.89	41.92	42.86	45.64	48.29	50.83	54.05	56.41
27	31.53	32.91	34.57	36.74	40.11	43.19	44.14	46.96	49.64	52.22	55.48	57.86
28	32.62	34.03	35.71	37.92	41.34	44.46	45.42	48.28	50.99	53.59	56.89	59.30
29	33.71	35.14	36.85	39.09	42.56	45.72	46.69	49.59	52.34	54.97	58.30	60.73
30	34.80	36.25	37.99	40.26	43.77	46.98	47.96	50.89	53.67	56.33	59.70	62.16
40	45.62	47.27	49.24	51.81	55.76	59.34	60.44	63.69	66.77	69.70	73.40	76.09
50	56.33	58.16	60.35	63.17	67.50	71.42	72.61	76.15	79.49	82.66	86.66	89.56
60	66.98	68.97	71.34	74.40	79.08	83.30	84.58	88.38	91.95	95.34	99.61	102.7
80	88.13	90.41	93.11	96.58	101.9	106.6	108.1	112.3	116.3	120.1	124.8	128.3
100	109.1	111.7	114.7	118.5	124.3	129.6	131.1	135.8	140.2	144.3	149.4	153.2



# WORKSHEET

What is the critical value of  $\chi^2$  when  $\alpha=0.05$  and...

...the table has 2 rows and 3 columns?

...the table has 5 rows and 3 columns?

What is the critical value of  $\chi^2$  when  $\alpha=0.01$  and...

...the table has 3 rows and 4 columns?

...the table has 4 rows and 4 columns?

# $\chi^2$ Tests

Calculate the test statistic ...  $\chi^2$

The  $\chi^2$  test statistic is based on a comparison of the observed cell frequencies to the cell frequencies that we expect under the null hypothesis

For the cell in row  $i$  and column  $j$ , the “expected” frequency under the null hypothesis is

$$\hat{f}_{ij} = \frac{(f_{i\bullet})(f_{\bullet j})}{N}$$

where  $f_{i\bullet}$  is the number of cases in row  $i$ ,  $f_{\bullet j}$  is the number of cases in column  $j$ , and  $N$  is the sample size

# $\chi^2$ Tests

Calculate the test statistic ...  $\chi^2$

For the cell in row  $i=1$  and column  $j=1$  (the top left cell):

$$\hat{f}_{1,1} = \frac{(f_{1\cdot})(f_{\cdot 1})}{N} = \frac{(456)(516)}{1,304} = 180$$

For the cell in row  $i=1$   
and column  $j=2$ :

$$\hat{f}_{1,2} = \frac{(456)(454)}{1,304} = 159$$

And so forth...

How Important Are "Loose Morals and Drunkenness" in Explaining Poverty

	Very Important	Somewhat Important	Not Important	Row Total
Democrat	164 180	161 159	131 117	456
Independent	179 165	139 146	100 107	418
Republican	173 170	154 150	103 110	430
Column Total	516	454	334	N=1,304

Political Party

# $\chi^2$ Tests

Calculate the test statistic ...  $\chi^2$

The  $\chi^2$  test statistic equals

$$\chi^2 = \text{Sum of } \frac{(\text{Expected} - \text{Observed})^2}{\text{Expected}} = \sum_{i=1}^R \sum_{j=1}^C \frac{(\hat{f}_{ij} - f_{ij})^2}{\hat{f}_{ij}}$$

where

R is the number of rows in the crosstable,

C is the number of columns in the crosstable,

$f_{ij}$  is the observed frequency in the cell in row i and column j, &

$\hat{f}_{ij}$  is the expected frequency in the cell in row i and column j

# $\chi^2$ Tests

Calculate the test statistic ...  $\chi^2$

$$\chi^2 = \frac{(180 - 164)^2}{180} + \frac{(159 - 161)^2}{159} + \dots + \frac{(110 - 103)^2}{110} = 5.767$$

In this example,  $\chi^2 = 5.767$

How Important Are “Loose Morals and Drunkenness” in Explaining Poverty

		<i>Very Important</i>	<i>Somewhat Important</i>	<i>Not Important</i>	Row Total
<b>Political Party</b>	<i>Democrat</i>	164 180	161 159	131 117	456
	<i>Independent</i>	179 165	139 146	100 107	418
	<i>Republican</i>	173 170	154 150	103 110	430
	Column Total	516	454	334	N=1,304

# $\chi^2$ Tests

## Compare the test statistic to the critical value

1. If the test statistic is as large or larger than the critical value, then reject  $H_0$  (with probability of  $\alpha$  of doing so even though  $H_0$  should not actually be rejected)
2. If the test statistic is less than the critical value, then do not reject  $H_0$  (with probability of  $\beta$  of doing so even though  $H_0$  should be rejected)

We can restate the hypotheses as

$H_0$ : X & Y are independent  $\rightarrow$  Fail to Reject if  $\chi^2 \leq 9.49$

$H_1$ : X & Y not independent  $\rightarrow$  Reject if  $\chi^2 > 9.49$

Since  $\chi^2 = 5.767$ , we fail to reject  $H_0$

# $\chi^2$ Tests

## How Important Are “Loose Morals and Drunkenness” in Explaining Poverty

Source: GSS

		<i>Very Important</i>	<i>Somewhat Important</i>	<i>Not Important</i>	Row Total
<b>Political Party</b>	<i>Democrat</i>	164	161	131	456
	<i>Independent</i>	179	139	100	418
	<i>Republican</i>	173	154	103	430
	Column Total	516	454	334	N=1,304

## X & Y in the Population ... $\chi^2=0$

Y	X			Total
	1	2	3	
1	10,000	10,000	10,000	30,000
2	10,000	10,000	10,000	30,000
Total	20,000	20,000	20,000	60,000

## X & Y in a Sample of 100 ... $\chi^2=2.2$

Y	X			Total
	1	2	3	
1	20	17	14	51
2	13	17	19	49
Total	33	34	33	100

## X & Y in Another Sample of 100 ... $\chi^2=1.7$

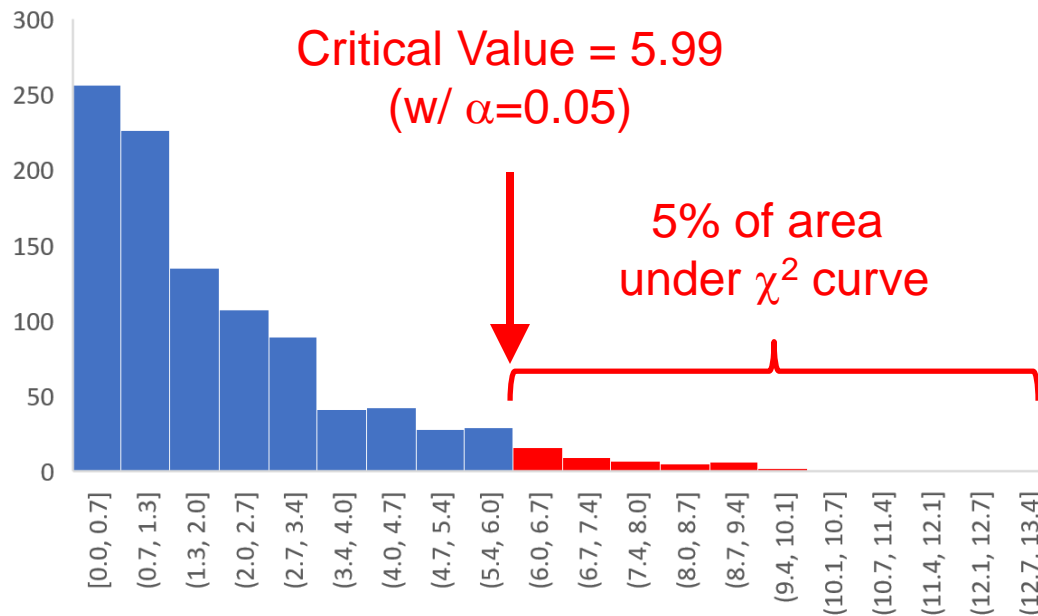
Y	X			Total
	1	2	3	
1	17	22	13	52
2	16	15	17	48
Total	33	37	30	100



# X & Y in the Population ... $\chi^2=0$

Y	X			Total
	1	2	3	
1	10,000	10,000	10,000	30,000
2	10,000	10,000	10,000	30,000
Total	20,000	20,000	20,000	60,000

## Distribution of $\chi^2$ across 1,000 Samples of n=100



# $\chi^2$ Example

Is there an association between people's childhood family income and their income as adults?

Source: GSS

		Family Income in Childhood			Row Total
		<i>Below Average</i>	<i>Average</i>	<i>Above Average</i>	
Family Income as an Adult	<i>Below Average</i>	4,085	4,027	1,054	9,166
	<i>Average</i>	4,796	9,853	2,121	16,770
	<i>Above Average</i>	1,494	2,886	2,075	6,455
Column Total		10,375	16,766	5,250	N=32,391

# $\chi^2$ Example

Is there an association between people's childhood family income and their income as adults?

Source: GSS

		Family Income in Childhood			Row Total
		<i>Below Average</i>	<i>Average</i>	<i>Above Average</i>	
Family Income as an Adult	<i>Below Average</i>	39%	24%	20%	28%
	<i>Average</i>	46%	59%	40%	52%
	<i>Above Average</i>	14%	17%	40%	20%
Column Total		100%	100%	100%	100%

# $\chi^2$ Example

State the null ( $H_0$ ) and alternative ( $H_1$ ) hypotheses

Again, we begin with the assumption that there is no association between X and Y

$H_0$ : Family income in childhood and family income in adulthood are statistically independent

$H_1$ : Family income in childhood and family income in adulthood are not statistically independent

# $\chi^2$ Example

Check that the sample data conform to basic assumptions;  
if they do not, then do not go any further

X and Y must be collected from a random sample of  
individuals from the population (OK)

Standard  $\chi^2$  testing procedures should be used with  
extreme caution — or not at all — if any cell frequency is  
less than 5 (OK)

# $\chi^2$ Example

Choose an  $\alpha$  probability level ... that is, a probability associated with incorrectly rejecting the null hypothesis

Let's use  $\alpha=0.01$  in our example

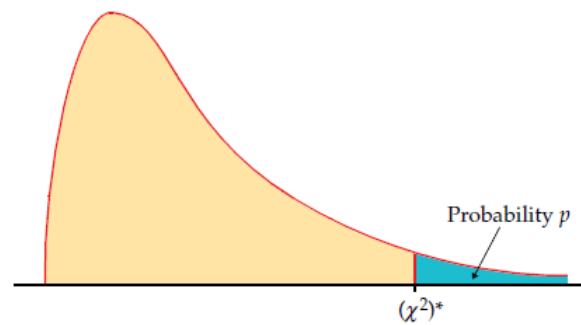
# $\chi^2$ Example

Determine the “critical value” ... that is, how large the test statistic must be in order to reject the null hypothesis at the given  $\alpha$  level

In our example,  $\alpha=0.01$  and  $df=(3-1)(3-1)=4$

According to a  $\chi^2$  table, the critical value of  $\chi^2$  thus equals 13.28

Table entry for  $p$  is the critical value  $(\chi^2)^*$  with probability  $p$  lying to its right.



**TABLE F**

$\chi^2$  distribution critical values

df	Tail probability $p$											
	.25	.20	.15	.10	.05	.025	.02	.01	.005	.0025	.001	.0005
1	1.32	1.64	2.07	2.71	3.84	5.02	5.41	6.63	7.88	9.14	10.83	12.12
2	2.77	3.22	3.79	4.61	5.99	7.38	7.82	9.21	10.60	11.98	13.82	15.20
3	4.11	4.64	5.32	6.25	7.81	9.35	9.84	11.34	12.84	14.32	16.27	17.73
4	5.39	5.99	6.74	7.78	9.49	11.14	11.67	13.28	14.86	16.42	18.47	20.00
5	6.63	7.29	8.12	9.24	11.07	12.83	13.39	15.09	16.75	18.39	20.51	22.11
6	7.84	8.56	9.45	10.64	12.59	14.45	15.03	16.81	18.55	20.25	22.46	24.10
7	9.04	9.80	10.75	12.02	14.07	16.01	16.62	18.48	20.28	22.04	24.32	26.02
8	10.22	11.03	12.03	13.36	15.51	17.53	18.17	20.09	21.95	23.77	26.12	27.87
9	11.39	12.24	13.29	14.68	16.92	19.02	19.68	21.67	23.59	25.46	27.88	29.67
10	12.55	13.44	14.53	15.99	18.31	20.48	21.16	23.21	25.19	27.11	29.59	31.42
11	13.70	14.63	15.77	17.28	19.68	21.92	22.62	24.72	26.76	28.73	31.26	33.14
12	14.85	15.81	16.99	18.55	21.03	23.34	24.05	26.22	28.30	30.32	32.91	34.82
13	15.98	16.98	18.20	19.81	22.36	24.74	25.47	27.69	29.82	31.88	34.53	36.48
14	17.12	18.15	19.41	21.06	23.68	26.12	26.87	29.14	31.32	33.43	36.12	38.11
15	18.25	19.31	20.60	22.31	25.00	27.49	28.26	30.58	32.80	34.95	37.70	39.72
16	19.37	20.47	21.79	23.54	26.30	28.85	29.63	32.00	34.27	36.46	39.25	41.31
17	20.49	21.61	22.98	24.77	27.59	30.19	31.00	33.41	35.72	37.95	40.79	42.88
18	21.60	22.76	24.16	25.99	28.87	31.53	32.35	34.81	37.16	39.42	42.31	44.43
19	22.72	23.90	25.33	27.20	30.14	32.85	33.69	36.19	38.58	40.88	43.82	45.97
20	23.83	25.04	26.50	28.41	31.41	34.17	35.02	37.57	40.00	42.34	45.31	47.50
21	24.93	26.17	27.66	29.62	32.67	35.48	36.34	38.93	41.40	43.78	46.80	49.01
22	26.04	27.30	28.82	30.81	33.92	36.78	37.66	40.29	42.80	45.20	48.27	50.51
23	27.14	28.43	29.98	32.01	35.17	38.08	38.97	41.64	44.18	46.62	49.73	52.00
24	28.24	29.55	31.13	33.20	36.42	39.36	40.27	42.98	45.56	48.03	51.18	53.48
25	29.34	30.68	32.28	34.38	37.65	40.65	41.57	44.31	46.93	49.44	52.62	54.95
26	30.43	31.79	33.43	35.56	38.89	41.92	42.86	45.64	48.29	50.83	54.05	56.41
27	31.53	32.91	34.57	36.74	40.11	43.19	44.14	46.96	49.64	52.22	55.48	57.86
28	32.62	34.03	35.71	37.92	41.34	44.46	45.42	48.28	50.99	53.59	56.89	59.30
29	33.71	35.14	36.85	39.09	42.56	45.72	46.69	49.59	52.34	54.97	58.30	60.73
30	34.80	36.25	37.99	40.26	43.77	46.98	47.96	50.89	53.67	56.33	59.70	62.16
40	45.62	47.27	49.24	51.81	55.76	59.34	60.44	63.69	66.77	69.70	73.40	76.09
50	56.33	58.16	60.35	63.17	67.50	71.42	72.61	76.15	79.49	82.66	86.66	89.56
60	66.98	68.97	71.34	74.40	79.08	83.30	84.58	88.38	91.95	95.34	99.61	102.7
80	88.13	90.41	93.11	96.58	101.9	106.6	108.1	112.3	116.3	120.1	124.8	128.3
100	109.1	111.7	114.7	118.5	124.3	129.6	131.1	135.8	140.2	144.3	149.4	153.2



# $\chi^2$ Example

Calculate the test statistic ...  $\chi^2$

$$\hat{f}_{i,j} = \frac{(f_{i\cdot})(f_{\cdot j})}{N}$$

		Family Income in Childhood			Row Total
		<i>Below Average</i>	<i>Average</i>	<i>Above Average</i>	
Family Income as an Adult	<i>Below Average</i>	4,085 <b>2,935.92</b>	4,027 <b>4,744.44</b>	1,054 <b>1,485.64</b>	9,166
	<i>Average</i>	4,796 <b>5,371.52</b>	9,853 <b>8,680.37</b>	2,121 <b>2,718.12</b>	16,770
	<i>Above Average</i>	1,494 <b>2,067.57</b>	2,886 <b>3,341.19</b>	2,075 <b>1,046.24</b>	6,455
Column Total		10,375	16,766	5,250	N=32,391

# $\chi^2$ Example

Calculate the test statistic ...  $\chi^2$

$$\chi^2 = \text{Sum of } \frac{(\text{Expected} - \text{Observed})^2}{\text{Expected}} = \sum_{i=1}^R \sum_{j=1}^C \frac{(\hat{f}_{ij} - f_{ij})^2}{\hat{f}_{ij}}$$

In this example,  $\chi^2=2,267.59$

# $\chi^2$ Example

## Compare the test statistic to the critical value

1. If the test statistic is as large or larger than the critical value, then reject  $H_0$  (with probability of  $\alpha$  of doing so even though  $H_0$  should not actually be rejected)
2. If the test statistic is less than the critical value, then do not reject  $H_0$  (with probability of  $\beta$  of doing so even though  $H_0$  should be rejected)

We can restate the hypotheses as

$H_0$ : X & Y are independent  $\rightarrow$  Fail to Reject if  
 $\chi^2 \leq 13.28$

$H_1$ : X & Y not independent  $\rightarrow$  Reject if  $\chi^2 > 13.28$

Since  $\chi^2 = 2,267.59$ , we reject  $H_0$

# Worksheet

Are political views related to whether people view the bible as the literal word of God?

## Political Views

		<i>Liberal</i>	<i>Moderate</i>	<i>C'servative</i>	Row Total
Bible is...	<i>Word of God</i>	1,500	2,888	3,234	7,622
	<i>Something Else</i>	4,670	5,961	5,079	15,710
	Column Total	6,170	8,849	8,313	23,332

# $\chi^2$ vs Tests for Two Proportions

Does the proportion living past 60 vary by family SES?

Different methods yield the same results...

Using what we learned earlier about hypothesis tests for differences in proportions:  $Z=-2.453$ ;  $p\text{-value}=0.014$

Using a  $\chi^2$  test, we find:  $\chi^2=6.027$ ;  $p\text{-value}=0.014$

		Lived Past Age 60?	
		Yes=1	No=0
Family SES	Low=0	603	330
	High=1	456	319

# Factors Driving Significance

Two things affect our chances of observing a statistically significant association

1. The strength of the association in the population
2. The size of the sample ... consider the following two cross-tabs/results:

		Voted in 2016?			Voted in 2016?		
		Yes	No	Row Total	Yes	No	Row Total
Political Party Affiliation	Democrat	<b>6</b>	<b>4</b>	<b>10</b>	<b>300</b>	<b>200</b>	<b>500</b>
	Republican	<b>8</b>	<b>2</b>	<b>10</b>	<b>400</b>	<b>100</b>	<b>500</b>
Column Total		<b>14</b>	<b>6</b>	<b>20</b>	<b>700</b>	<b>300</b>	<b>1,000</b>
		$\chi^2 = 0.95$			$\chi^2 = 47.62$		

# Association Between Discrete Variables

We use  $\chi^2$  to assess whether there is any statistically significant association between two categorical variables, X and Y, in the population

We have said nothing about measuring the direction or strength of that association

If we conclude that categorical variables X and Y are associated, how can we quantify the direction and strength of that association?

# Associations Between Discrete Variables

**Is there an association in the population?**

$\chi^2$  analysis

**How strong is the association? In what direction is the association?**

Gamma, relative risk, odds ratios



# Association Between Discrete Variables

## Measures of Association

Statistics that show the direction and/or magnitude of a relationship between pairs of variables

### When X and Y are Both Ordinal:

Gamma

Others (Not discussed today)

### When X and Y are Both Dichotomous:

Gamma

Relative Risk (RR)

Odds Ratio (OR)

Others (Not discussed today)

# Association Between Ordinal Variables

Source: GSS  
1972-2008

		Family Income in Childhood			Row
		<i>Below Average</i>	<i>Average</i>	<i>Above Average</i>	Total
Family Income as an Adult	<i>Above Average</i>	1,494	2,886	2,075	6,455
	<i>Average</i>	4,796	9,853	2,121	16,770
	<i>Below Average</i>	4,085	4,027	1,054	9,166
	Column Total	10,375	16,766	5,250	N=32,391

# Gamma

$\chi^2 = 2,267.59$ ; rejected  $H_0$  ... that is, we rejected the null hypothesis that there is no association in the population

**But how strong is the association? And in which direction?**

**Gamma** is a measure of association for use when both variables are ordinal

		Family Income in Childhood			Row Total
		<i>Below Average</i>	<i>Average</i>	<i>Above Average</i>	
Family Income as an Adult	<i>Above Average</i>	1,494	2,886	2,075	6,455
	<i>Average</i>	4,796	9,853	2,121	16,770
	<i>Below Average</i>	4,085	4,027	1,054	9,166
Column Total		10,375	16,766	5,250	32,391

Source:  
GSS,  
1972-  
2008

# Gamma

Here is this analysis done separately for 1972 and 2008. Are the two variables associated in both years? Has the strength of the association changed over time?

**1972**

$\chi^2=130.79$   
( $p<0.01$ )

		Family Income in Childhood			Row Total
		<i>Below Average</i>	<i>Average</i>	<i>Above Average</i>	
Family Income as an Adult	<i>Above Average</i>	40	72	71	183
	<i>Average</i>	173	591	135	899
	<i>Below Average</i>	181	241	65	495
<b>Column Total</b>		402	904	271	1,577

Source: GSS

**2008**

$\chi^2=181.19$   
( $p<0.01$ )

		Family Income in Childhood			Row Total
		<i>Below Average</i>	<i>Average</i>	<i>Above Average</i>	
Family Income as an Adult	<i>Above Average</i>	87	124	159	370
	<i>Average</i>	263	472	165	900
	<i>Below Average</i>	331	274	100	705
<b>Column Total</b>		681	870	424	1,975

Source: GSS

# Gamma

The possible values of gamma range from -1 to +1

Gamma = -1 = a perfect negative association

Gamma = +1 = a perfect positive association

Gamma = 0 = no association

**Gamma = +0.986**

	x=1	x=2	x=3
y=3	1	1	100
y=2	1	100	1
y=1	100	1	1

**Gamma = -0.986**

	x=1	x=2	x=3
y=3	100	1	1
y=2	1	100	1
y=1	1	1	100

**Gamma = 0.000**

	x=1	x=2	x=3
y=3	1	1	1
y=2	1	1	1
y=1	1	1	1

# Gamma

Calculating gamma is semi-complicated; see the “bonus slides” at the end for instructions and an example

Focus for now on how it is interpreted...

Is the association positive or negative?

Is it weak or strong?

# Gamma

Here is this analysis done separately for 1972 and 2008. Are the two variables associated in both years? Has the strength of the association changed over time?

**1972**

Family Income in  
Childhood

Source:  
GSS

**Gamma=  
0.312**

		<i>Below Average</i>	<i>Average</i>	<i>Above Average</i>	<b>Row Total</b>
Family Income as an Adult	<i>Above Average</i>	40	72	71	183
	<i>Average</i>	173	591	135	899
	<i>Below Average</i>	181	241	65	495
<b>Column Total</b>		402	904	271	1,577

**2008**

Family Income in  
Childhood

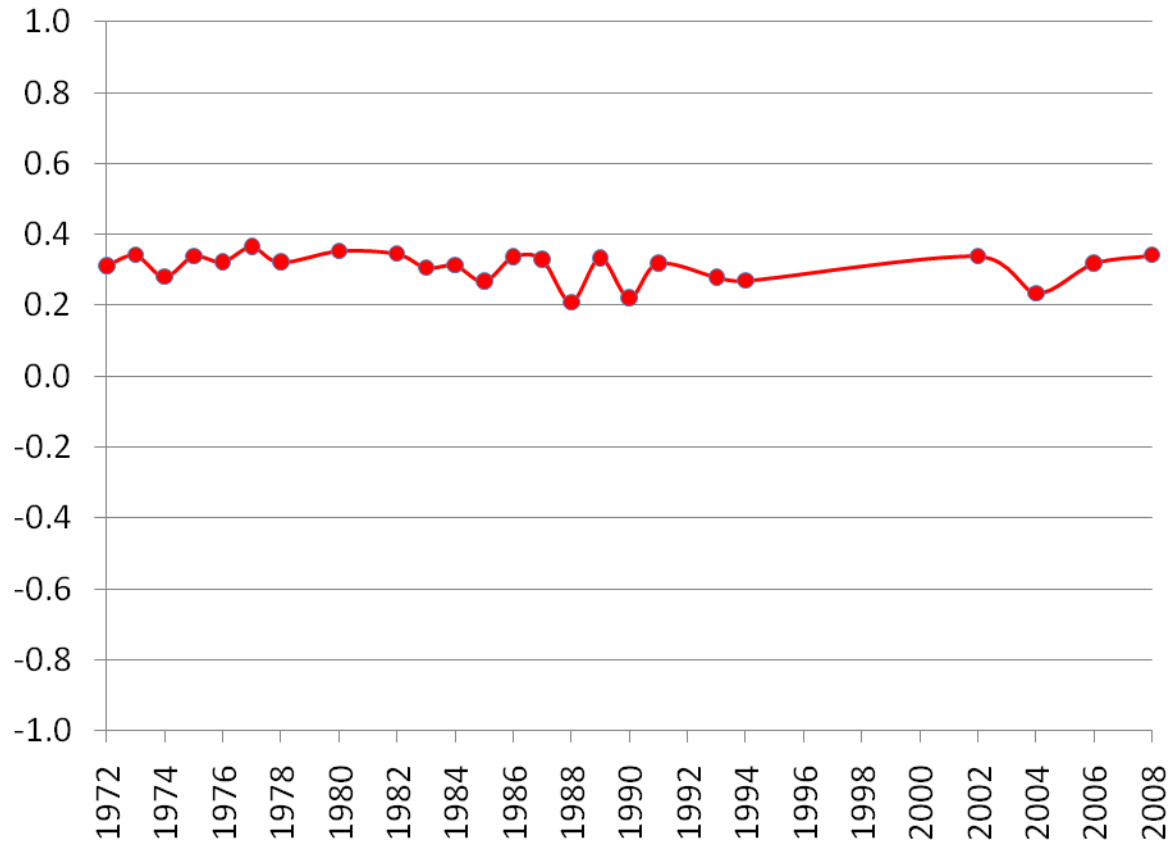
Source:  
GSS

**Gamma=  
0.342**

		<i>Below Average</i>	<i>Average</i>	<i>Above Average</i>	<b>Row Total</b>
Family Income as an Adult	<i>Above Average</i>	87	124	159	370
	<i>Average</i>	263	472	165	900
	<i>Below Average</i>	331	274	100	705
<b>Column Total</b>		681	870	424	1,975

# Gamma

And separately by year...





# Association in 2x2 Tables

Cross-tabulations of categorical variables that each have 2 categories—that is, that are dichotomous—are a special case of cross-tabulations of ordinal variables

We will explore three measures of association that pertain to 2x2 tables

Gamma

Relative Risk

Odds Ratio

# Association in 2x2 Tables

Each measure relies on a particular naming convention for the four cells in the cross-tabulation

		X	
		X=1	X=2
Y	Y=2	Cell a	Cell b
	Y=1	Cell c	Cell d

# Relative Risk

Gamma expresses the association between two dichotomous variables without regard to whether one variable affects change in the other ... they are thus referred to as “symmetric” measures

**Relative Risk** (and later **Odds Ratios**) are “asymmetric” measures, in that they express the change in the chances of being in a particular category of the “dependent” (Y) variable that result from changing categories of the “independent” (X) variable

# Relative Risk

When  $X=1$ , the “Risk” that  $Y=2$  equals  $a/(a+c)$  ... which is the same as  $P(Y=2 | X=1)$

When  $X=2$ , the “Risk” that  $Y=2$  equals  $b/(b+d)$  ... which is the same as  $P(Y=2 | X=2)$

The **relative risk (RR)**

is the ratio of the

two risks:

$$RR = \frac{b/(b+d)}{a/(a+c)}$$

		X	
		X=1	X=2
Y	Y=2	a	b
	Y=1	c	d

# Relative Risk

The **risk** of being in excellent health ( $Y=2$ ) if you **do not** smoke ( $X=1$ ) is  $a/(a+c) = 1,955/(3,838+1,955)=0.337$

The **risk** of being in excellent health ( $Y=2$ ) if you **do** smoke ( $X=2$ ) is  $b/(b+d) = 826/(2,356+826)=0.260$

The Relative Risk (RR)

is the ratio of the

two risks:

$$RR = \frac{0.260}{0.337} = 0.77$$

Health

		Current Smoker? <small>Source: GSS</small>	
		$X=1$ : No	$X=2$ : Yes
Y=2: <i>Excellent</i>	(a)	1,955	(b) 826
	(c)	3,838	(d) 2,356
Y=1: <i>Not Excellent</i>			

# Relative Risk

Relative Risk can range from zero to infinity

**RR 0.00 to 1.00** means that the risk that  $Y=2$  is **reduced** when you move from  $X=1$  to  $X=2$

**RR = 1.00** means that the risk that  $Y=2$  is **unchanged** when you move from  $X=1$  to  $X=2$

**RR > 1.00** means that the risk that  $Y=2$  is **increased** when you move from  $X=1$  to  $X=2$

# Relative Risk

Interpret the value of RR with reference to the number 1

**RR = 0.77:** The risk that Y=2 is **reduced by 23%** when you move from X=1 to X=2

**RR = 1.20:** The risk that Y=2 is **increased by 20%** when you move from X=1 to X=2

**RR = 2.50:** The risk that Y=2 is **increased by 150% (or is 2.5 times greater)** when you move from X=1 to X=2

# Odds Ratio

When  $X=1$ , the “Odds” that  $Y=2$  equals  $a/c$

When  $X=2$ , the “Odds” that  $Y=2$  equals  $b/d$

The “Odds Ratio” (OR) is the ratio of the two odds:

$$OR = \frac{b/d}{a/c} = \frac{bc}{ad}$$

		X	
		X=1	X=2
Y	Y=2	a	b
	Y=1	c	d



# Odds Ratio

The **odds** of being in excellent health (Y=2) if you **do not** smoke (X=1) is  $a/c = 1,955/3,838=0.509$

The **odds** of being in excellent health (Y=2) if you **do** smoke (X=2) is  $b/d = 826/2,356=0.351$

The Odds Ratio (OR)  
is the ratio of the two  
odds:

$$OR = \frac{0.351}{0.509} = 0.69$$

Health

	Current Smoker? <small>Source: GSS</small>	
	X=1: No	X=2: Yes
Y=2: Excellent	(a) 1,955	(b) 826
Y=1: Not Excellent	(c) 3,838	(d) 2,356

# Odds Ratio

Odds Ratios can range from zero to infinity

**OR 0.00 to 1.00** means that the odds that  $Y=2$  is **reduced** when you move from  $X=1$  to  $X=2$

**OR = 1.00** means that the odds that  $Y=2$  is **unchanged** when you move from  $X=1$  to  $X=2$

**OR > 1.00** means that the odds that  $Y=2$  is **increased** when you move from  $X=1$  to  $X=2$

# Odds Ratio

Interpret the value of OR with reference to the number 1

**OR = 0.77**: The odds that  $Y=2$  is **reduced by 23%** when you move from  $X=1$  to  $X=2$

**OR = 1.20**: The odds that  $Y=2$  is **increased by 20%** when you move from  $X=1$  to  $X=2$

**OR = 2.50**: The odds that  $Y=2$  is **increased by 150% (or is 2.5 times greater)** when you move from  $X=1$  to  $X=2$

# Worksheet

How is race (white vs. black) associated with political party preference? ( $\chi^2=2,214.2$  w/ 1 df; Reject  $H_0$  w/  $\alpha=0.01$ ). Compute **and interpret** RR and OR

*Source: GSS*

		Race		Row Total
		<i>X=1: White</i>	<i>X=2: Black</i>	
Party	<i>Y=2: Democrat</i>	15,594	4,401	19,995
	<i>Y=1: Republican</i>	13,288	505	13,793
Column Total		28,882	4,906	N=33,788

# Practical vs. Statistical Significance

Even when we can reject the hypothesis of no association, we should always investigate the direction and magnitude of the association

		Ever Smokes Pot?		Row Total	$\chi^2 = 4.5$ w/ $df=1$ Reject $H_0$ at $\alpha=0.05$
		Yes	No		
Right or Left Handed	<b>Right</b>	100,000	200,000	300,000	Odds Ratio = 1.03
	<b>Left</b>	10,000	20,550	30,550	
Column Total		110,000	220,550	330,550	

Right handed people are 3% more likely to have ever smoked pot ... is this a practically meaningful result?

# Want More?

David Lane's Book

[http://onlinestatbook.com/2/chi\\_square/Chi\\_Square.html](http://onlinestatbook.com/2/chi_square/Chi_Square.html)

Chapter 8 of Richard Lowry's Book

<http://vassarstats.net/textbook/>

Gerard Dallal's Book

<http://www.jerrydallal.com/LHSP/ctab.htm>

Stat Trek

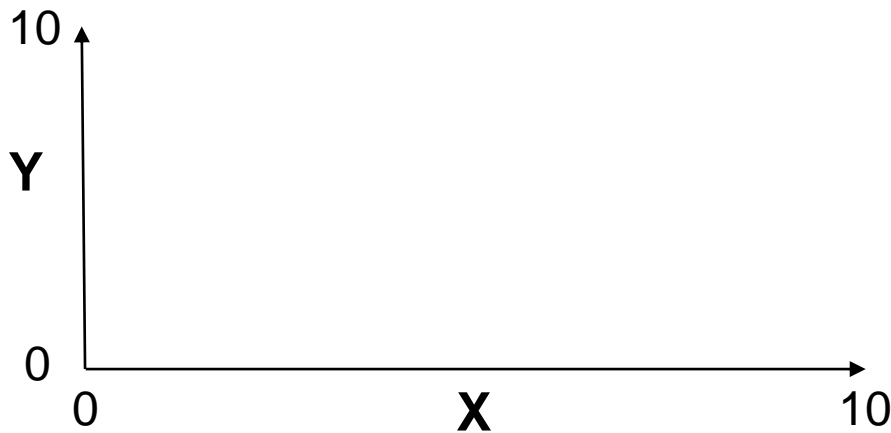
<http://stattrek.com/chi-square-test/homogeneity.aspx?tutorial=ap>

Bonus:  
The Next Slides Show  
How to Compute Gamma  
(If you are interested...)

# Bonus: Computing Gamma

In order to compute Gamma, the crosstable must be arranged such that

- (1) the column variable is arranged from lowest to highest going from left to right
- (2) the row variable is arranged from lowest to highest going from the bottom to the top



		Family Income in Childhood			Row Total
		<i>Below Average</i>	<i>Average</i>	<i>Above Average</i>	
Family Income as an Adult	<i>Above Average</i>	1,494	2,886	2,075	6,455
	<i>Average</i>	4,796	9,853	2,121	16,770
	<i>Below Average</i>	4,085	4,027	1,054	9,166
Column Total		10,375	16,766	5,250	32,391

Source:  
GSS,  
1972-2008



# Bonus: Computing Gamma

The formula for Gamma requires evaluating all pairs of observations in a cross-tabulation, counting the total number that are untied concordant pairs ( $n_c$ ) and the total number that are untied discordant pairs ( $n_d$ )

An untied pair is “one in which both cases have different values on two variables”

		Family Income in Childhood			Row Total
		<i>Below Average</i>	<i>Average</i>	<i>Above Average</i>	
Family Income as an Adult	<i>Above Average</i>	1,494	2,886	2,075	6,455
	<i>Average</i>	4,796	9,853	2,121	16,770
	<i>Below Average</i>	4,085	4,027	1,054	9,166
Column Total		10,375	16,766	5,250	32,391

Source:  
GSS,  
1972-2008

# Bonus: Computing Gamma

## Concordant Pairs

“One observation has a higher rank on both variables than does the other member of the pair”

Example: On both variables, the 2,075 cases in the top right cell have higher rank than the  $4,796+9,853+4,085+4,027=22,761$  cases in the four cells in the bottom left

Thus the observations in the top right cell contribute  
 $(2,075)(22,761)=47,229,075$   
 concordant pairs

*Source:*  
GSS,  
1972-2008

		Family Income in Childhood			Row Total
		<i>Below Average</i>	<i>Average</i>	<i>Above Average</i>	
Family Income as an Adult	<i>Above Average</i>	1,494	2,886	2,075	6,455
	<i>Average</i>	4,796	9,853	2,121	16,770
	<i>Below Average</i>	4,085	4,027	1,054	9,166
Column Total		10,375	16,766	5,250	32,391

# Bonus: Computing Gamma

Cells contribute concordant pairs if there are other cells with lower rank on both variables

For the top-right:  $(2,075)(22,761)=47,229,075$

For the top-middle:  $(2,886)(4,796+4,085)=25,630,566$

For the middle-right:  $(2,121)(4,085+4,027)=17,205,552$

For the middle-middle:  $(9,853)(4,085)=40,249,505$

Summing, there are  $n_c =$

$47,229,075 + 25,630,566 +$

$17,205,552 + 40,249,505 =$

$130,314,698$  concordant pairs

		Family Income in Childhood			Row Total
		<i>Below Average</i>	<i>Average</i>	<i>Above Average</i>	
Family Income as an Adult	<i>Above Average</i>	1,494	2,886	2,075	6,455
	<i>Average</i>	4,796	9,853	2,121	16,770
	<i>Below Average</i>	4,085	4,027	1,054	9,166
Column Total		10,375	16,766	5,250	32,391

Source:  
GSS,  
1972-2008

# Bonus: Computing Gamma

## Discordant Pairs

“One member of a pair of observations ranks higher than the other member on one variable but ranks lower on the other variable”

Example: The 1,494 cases in the top left cell have higher rank than the  $9,853+2,121+4,027+1,054=17,055$  cases in the four cells in the bottom right on one variable,

but lower rank on the other

Thus the observations in the top left cell contribute

$$(1,494)(17,055)=25,480,170$$

discordant pairs

		Family Income in Childhood			Row Total
		<i>Below Average</i>	<i>Average</i>	<i>Above Average</i>	
Family Income as an Adult	<i>Above Average</i>	1,494	2,886	2,075	6,455
	<i>Average</i>	4,796	9,853	2,121	16,770
	<i>Below Average</i>	4,085	4,027	1,054	9,166
Column Total		10,375	16,766	5,250	32,391

*Source:*  
GSS,  
1972-2008

# Bonus: Computing Gamma

Cells contribute discordant pairs if there are other cells with lower rank on one variable, higher rank on another

For the top-left:  $(1,494)(17,055)=25,480,170$

For the top-middle:  $(2,886)(2,121+1,054)=9,163,050$

For the middle-left:  $(4,796)(4,027+1,054)=24,368,476$

For the middle-middle:  $(9,853)(1,054)=10,385,062$

Summing, there are  $n_d =$   
 $25,480,170 + 9,163,050 +$   
 $24,368,476 + 10,385,062 =$   
 $69,396,758$  discordant pairs

		Family Income in Childhood			Row Total
		<i>Below Average</i>	<i>Average</i>	<i>Above Average</i>	
Family Income as an Adult	<i>Above Average</i>	1,494	2,886	2,075	6,455
	<i>Average</i>	4,796	9,853	2,121	16,770
	<i>Below Average</i>	4,085	4,027	1,054	9,166
Column Total		10,375	16,766	5,250	32,391

Source:  
GSS,  
1972-2008

# Bonus: Computing Gamma

Gamma can be expressed as:

$$G = \frac{n_c - n_d}{n_c + n_d}$$

So in our example:

$$G = \frac{130,314,698 - 69,396,758}{130,314,698 + 69,396,758}$$
$$= 0.305$$

Source: GSS, 1972-2008

		Family Income in Childhood			Row Total
		<i>Below Average</i>	<i>Average</i>	<i>Above Average</i>	
Family Income as an Adult	<i>Above Average</i>	1,494	2,886	2,075	6,455
	<i>Average</i>	4,796	9,853	2,121	16,770
	<i>Below Average</i>	4,085	4,027	1,054	9,166
Column Total		10,375	16,766	5,250	32,391

# Bonus: Computing Gamma Example

Is there an association between X and Y?

$\chi^2=11.011$  w/ 4 df (reject  $H_0$  w/  $\alpha=0.05$ )

Gamma = \_\_\_\_\_ X

		<i>Low</i>	<i>Medium</i>	<i>High</i>	Row Total
Y	<i>High</i>	12	3	2	17
	<i>Medium</i>	7	10	4	21
	<i>Low</i>	6	8	9	23
Column Total		25	21	15	N=61

# Bonus: Computing Gamma Example

Is there an association between X and Y?

$$\chi^2 = 11.011 \text{ w/ 4 df (reject } H_0 \text{ w/ } \alpha = 0.05)$$

$$n_c = (2)(7+10+6+8) + (3)(7+6) + (4)(6+8) + (10)(6) = 217$$

$$n_d = (12)(10+4+8+9) + (3)(4+9) + (7)(8+9) + (10)(9) = 620$$

$$G = \frac{n_c - n_d}{n_c + n_d} = \frac{217 - 620}{217 + 620}$$

Gamma = -0.481

		X			Row Total
		Low	Medium	High	
Y	High	12	3	2	17
	Medium	7	10	4	21
	Low	6	8	9	23
Column Total		25	21	15	N=61