

# CHI-SQUARE TESTS

---

# Chi-Square Tests

- A Chi-Square Test for Independence
- Chi-Square Goodness-of-Fit Tests

# Contingency Table Example

Left-Handed vs. Gender

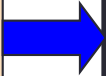
Dominant Hand: Left vs. Right

Gender: Male vs. Female

Sample results organized in a contingency table:

sample size =  $n = 300$ :

120 Females, 12  
were left handed  
180 Males, 24 were  
left handed



Gender	Hand Preference		
	Left	Right	
Female	12	108	120
Male	24	156	180
	36	264	300

# Chi-square Test on Independence

Consider a contingency table with  $r$  rows and  $c$  columns, test

- $H_0$ : The two categorical variables are independent, i.e., there is no relationship between them
- $H_1$ : The two categorical variables are dependent.

The Chi-square test statistic is

$$\chi^2 = \sum_{\text{all cells}} \frac{(f_{ij} - \hat{E}_{ij})^2}{\hat{E}_{ij}}$$

- where:
- $f_{ij}$  = observed cell frequency for  $i^{\text{th}}$  row and  $j^{\text{th}}$  column
- $r_i$  =  $i^{\text{th}}$  row total,  $c_j$  =  $j^{\text{th}}$  column total

$$\hat{E}_{ij} = \frac{r_i c_j}{n} = \text{expected cell frequency for } i^{\text{th}} \text{ row and } j^{\text{th}} \text{ column under independence}$$

(Assumed: each cell in the contingency table has expected frequency of at least 5)

# Chi-square Test on Independence

Consider a contingency table with  $r$  rows and  $c$  columns,

The test statistic

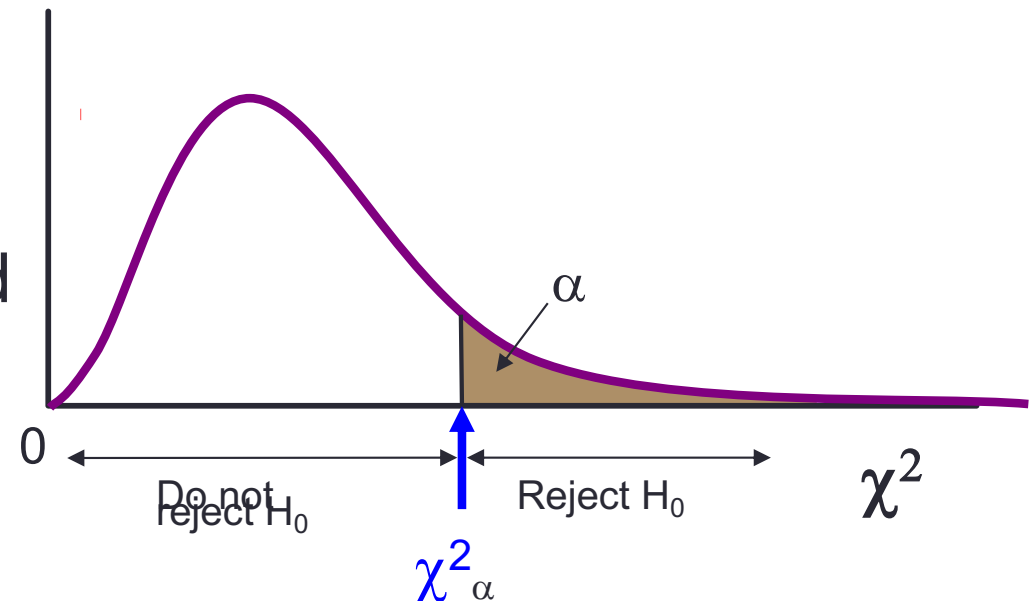
$$\chi^2 = \sum_{\text{all cells}} \frac{(f_{ij} - \hat{E}_{ij})^2}{\hat{E}_{ij}}$$

The null distribution is the chi-square distribution with  $(r-1)(c-1)$  degrees of freedom.

**Reject  $H_0$  if**

$\chi^2 > \chi_{\alpha}^2$  or if p-value  $< \alpha$

$\chi_{\alpha}^2$  and the p-value are based on  $(r-1)(c-1)$  degrees of freedom



# Observed vs. Expected Frequencies

Gender	Hand Preference		
	Left	Right	
Female	Observed = 12	Observed = 108	120
	Expected = 14.4	Expected = 105.6	
Male	Observed = 24	Observed = 156	180
	Expected = 21.6	Expected = 158.4	
	36	264	300

# The Chi-Square Test Statistic

Gender	Hand Preference		
	Left	Right	
Female	Observed = 12	Observed = 108	120
	Expected = 14.4	Expected = 105.6	
Male	Observed = 24	Observed = 156	180
	Expected = 21.6	Expected = 158.4	
	36	264	300

The test statistic is:

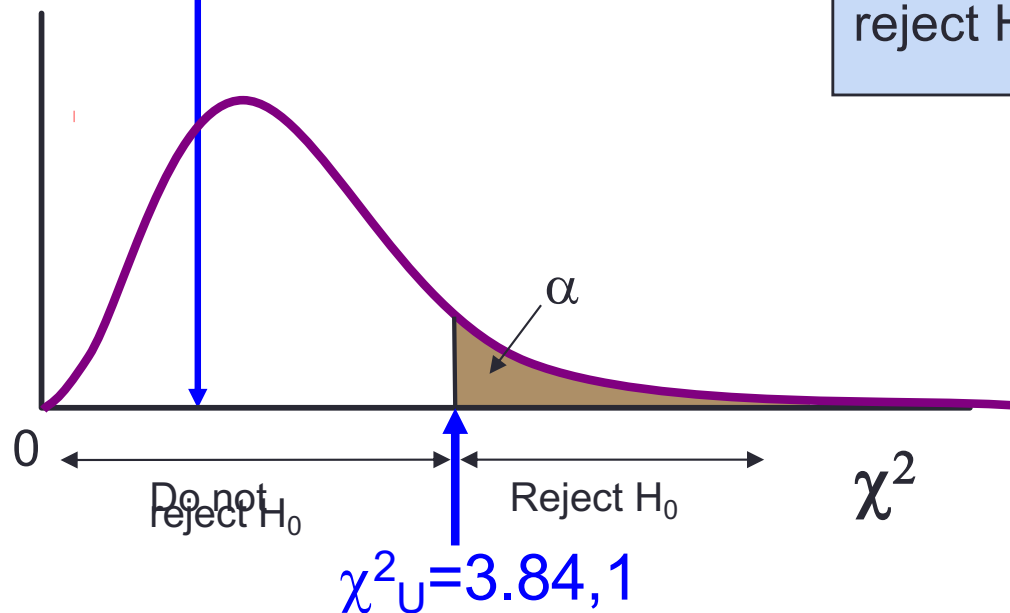
$$\begin{aligned}
 \chi^2 &= \sum_{\text{all cells}} \frac{(f_{ij} - \hat{E}_{ij})^2}{\hat{E}_{ij}} \\
 &= \frac{(12 - 14.4)^2}{14.4} + \frac{(108 - 105.6)^2}{105.6} + \frac{(24 - 21.6)^2}{21.6} + \frac{(156 - 158.4)^2}{158.4} = 0.7576
 \end{aligned}$$

# Decision Rule

The test statistic is  $\chi^2 = 0.7576$ ,  $\chi_\alpha^2$  with 1 d.f. = 3.841

## Decision Rule:

If  $\chi^2 > 3.841$ , reject  $H_0$ , otherwise, do not reject  $H_0$



Here,  
 $\chi^2 = 0.7576 < \chi_\alpha^2 = 3.841$ ,  
so we **do not reject  $H_0$**  and  
conclude that there is not sufficient  
evidence that the two proportions  
are different at  $\alpha = 0.05$

R code for  $\chi_\alpha^2$ : `qchisq(0.95,1)`



## Example: The Client Satisfaction Case

Client Satisfaction				
Fund	High	Low	Med	All
Bond	15	3	12	30
Stock	24	2	4	30
TaxDef	1	15	24	40
All	40	20	40	100

**H<sub>0</sub>: client satisfaction is independent of fund type**

**H<sub>a</sub>: client satisfaction depends upon fund type**

$$\begin{aligned}\chi^2 &= \sum_{\text{all cells}} \frac{(f_{ij} - \hat{E}_{ij})^2}{\hat{E}_{ij}}; \quad E_{BH} = \frac{r_B c_H}{n} = \frac{(30)(40)}{100} = 12, \dots, E_{TM} = \frac{r_T c_M}{n} = \frac{(40)(40)}{100} = 16 \\ &= \frac{(15-12)^2}{12} + \frac{(3-6)^2}{6} + \dots + \frac{(24-16)^2}{16} \\ &= 0.7500 + 1.5000 + \dots + 4.0000 = \mathbf{46.4375}\end{aligned}$$

$$\chi^2 = 46.4375 > 9.4877 = \chi^2_{.05}$$

Reject H<sub>0</sub>

$$p\text{-value} = P(\chi^2 > 46.4375) = 0.0000$$

## R function: `Chisq.test()`

```
M=as.table(rbind(c(15,3,12),c(24,2,4),c(1,15,24)))
dimnames(M)=list(Fund= c("Bond", "Stock",'TaxDef'),
                  Satisfaction=c("High","Low", "Med"))
```

```
Xsq <- chisq.test(M)) # Prints test summary
Xsq$observed # observed counts (same as M)
Xsq$expected # expected counts under the null
```

```
> Xsq$observed # observed counts (same as M)
```

```
      Satisfaction
Fund   High Low Med
Bond    15   3  12
Stock   24   2   4
TaxDef   1  15  24
```

```
> Xsq$expected # expected counts under the null
```

```
      Satisfaction
Fund   High Low Med
Bond    12   6  12
Stock   12   6  12
TaxDef   16   8  16
```

```
> Xsq
```

Pearson's Chi-squared test

data: M

X-squared = 46.4375, df = 4, p-value = 1.997e-09

# Chi-Square Goodness-of-Fit Test

- Does sample data conform to a hypothesized distribution?
- Examples:
  - Are technical support calls equal across all days of the week? (i.e., do calls follow a uniform distribution?)
  - Do measurements from a production process follow a normal distribution?

# Chi-Square Goodness-of-Fit Test

- Are technical support calls equal across all days of the week? (i.e., do calls follow a uniform distribution?)
  - Sample data for 10 days per day of week:

Sum of calls for this day:

Monday	290
Tuesday	250
Wednesday	238
Thursday	257
Friday	265
Saturday	230
Sunday	192
	<hr/>
	$\Sigma = 1722$

# Logic of Goodness-of-Fit Test

- If calls **are** uniformly distributed, the 1722 calls would be expected to be equally divided across the 7 days:

$$\frac{1722}{7} = 246 \quad \text{expected calls per day if uniform}$$

# Observed vs. Expected Frequencies

	Observed $f_o$	Expected $f_e$
Monday	290	246
Tuesday	250	246
Wednesday	238	246
Thursday	257	246
Friday	265	246
Saturday	230	246
Sunday	192	246
TOTAL	1722	1722

# Chi-Square Test Statistic

$H_0$ : The distribution of calls is uniform over days of the week

$H_1$ : The distribution of calls is not uniform

- The test statistic is

$$\chi^2 = \sum_k \frac{(f_o - f_e)^2}{f_e} \quad (\text{where } df = k - 1)$$

where:

$k$  = number of categories

$f_o$  = observed frequency

$f_e$  = expected frequency

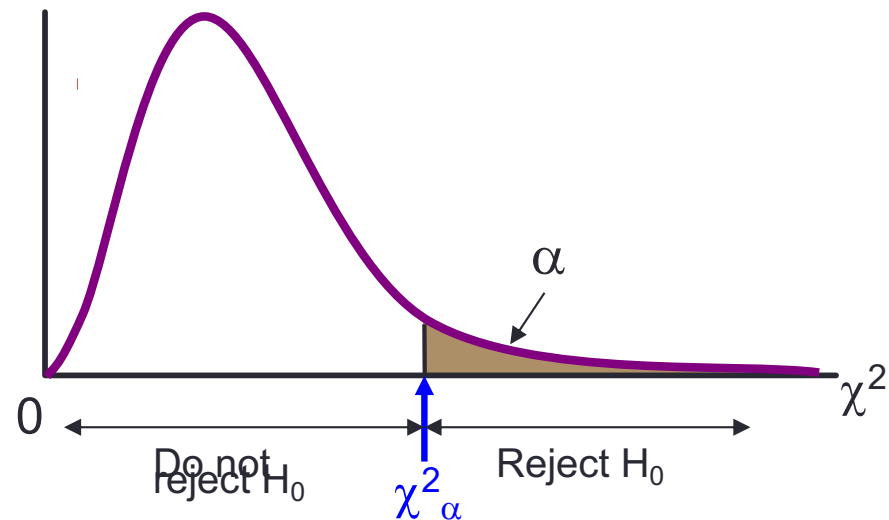
# The Rejection Region

$H_0$ : The distribution of calls is uniform over days of the week

$H_1$ : The distribution of calls is not uniform

$$\chi^2 = \sum_k \frac{(f_o - f_e)^2}{f_e}$$

- Reject  $H_0$  if  $\chi^2 > \chi^2_\alpha$





# Chi-Square Test Statistic

$H_0$ : The distribution of calls is uniform over days of the week

$H_1$ : The distribution of calls is not uniform

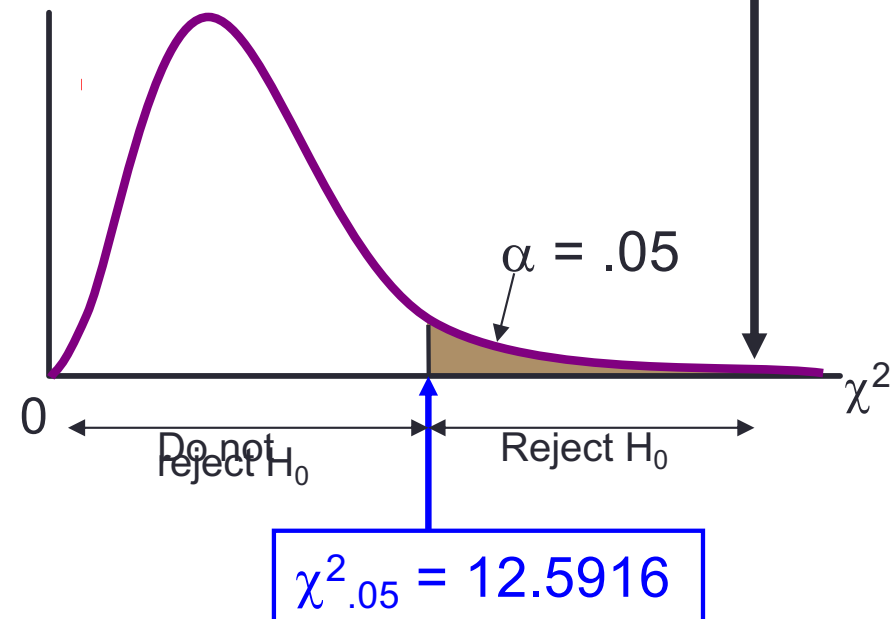
$$\chi^2 = \frac{(290 - 246)^2}{246} + \frac{(250 - 246)^2}{246} + \dots + \frac{(192 - 246)^2}{246} = 23.05$$

$k - 1 = 6$  ( $k = 7$  days of the week) so use 6 degrees of freedom:

$$\chi^2_{.05} = 12.5916$$

## Conclusion:

$\chi^2 = 23.05 > \chi^2_{\alpha} = 12.5916$   
so **reject  $H_0$**  and conclude that the distribution is not uniform



## R function: `Chisq.test()`

- `x<-c(290,250,238,257,265,230,192)`
- `prob=rep(1/length(x),length(x))`
- `chisq.test(x, p = prob)`

```
> chisq.test(x, p = prob)
```

```
Chi-squared test for given probabilities
```

```
data:  x
```

```
X-squared = 23.049, df = 6, p-value = 0.0007803
```

# Normal Distribution Example

- Do measurements from a production process follow a normal distribution with  $\mu = 50$  and  $\sigma = 15$ ?
- Process:
  - Get sample data
  - Group sample results into classes (cells)  
(Expected cell frequency must be at least 5 for each cell)
  - Compare actual cell frequencies with expected cell frequencies

# Normal Distribution Example

*(continued)*

- Sample data and values grouped into classes:

150 Sample Measurements		Class	Frequency
80		less than 30	10
65		30 but < 40	21
36		40 but < 50	33
66		50 but < 60	41
50		60 but < 70	26
38		70 but < 80	10
57		80 but < 90	7
77		90 or over	2
59			
...etc...			
		TOTAL	150

# Normal Distribution Example

*(continued)*

- What are the **expected frequencies** for these classes for a normal distribution with  $\mu = 50$  and  $\sigma = 15$ ?

Class	Frequency	Expected Frequency
less than 30	10	?
30 but < 40	21	
40 but < 50	33	
50 but < 60	41	
60 but < 70	26	
70 but < 80	10	
80 but < 90	7	
90 or over	2	
TOTAL	150	

# Expected Frequencies

Value	P(X < value)	Expected frequency
less than 30	0.09121	13.68
30 but < 40	0.16128	24.19
40 but < 50	0.24751	37.13
50 but < 60	0.24751	37.13
60 but < 70	0.16128	24.19
70 but < 80	0.06846	10.27
80 but < 90	0.01892	2.84
90 or over	0.00383	0.57
TOTAL	1.00000	150.00

Expected frequencies  
in a sample of size  
**n=150**, from a normal  
distribution with  
 $\mu=50$ ,  $\sigma=15$

**Example:**

$$\begin{aligned}
 P(x < 30) &= P\left(z < \frac{30 - 50}{15}\right) \\
 &= P(z < -1.3333) \\
 &= .0912
 \end{aligned}$$

$$(.0912)(150) = 13.68$$

Combine class groups so no class has  
expected frequency <1

# The Test Statistic

Class	Frequency (observed, $f_o$ )	Expected Frequency, $f_e$
less than 30	10	13.68
30 but < 40	21	24.19
40 but < 50	33	37.13
50 but < 60	41	37.13
60 but < 70	26	24.19
70 but < 80	10	10.27
80 or over	9	3.41
TOTAL	150	150.00

The test statistic is

$$\chi^2 = \sum_k \frac{(f_o - f_e)^2}{f_e}$$

- Reject  $H_0$  if

$$\chi^2 > \chi_\alpha^2$$

(with  $k - 1$  degrees of freedom)

# The Rejection Region

$H_0$ : The distribution of values is normal with  $\mu = 50$  and  $\sigma = 15$

$H_1$ : The distribution of calls does not have this distribution

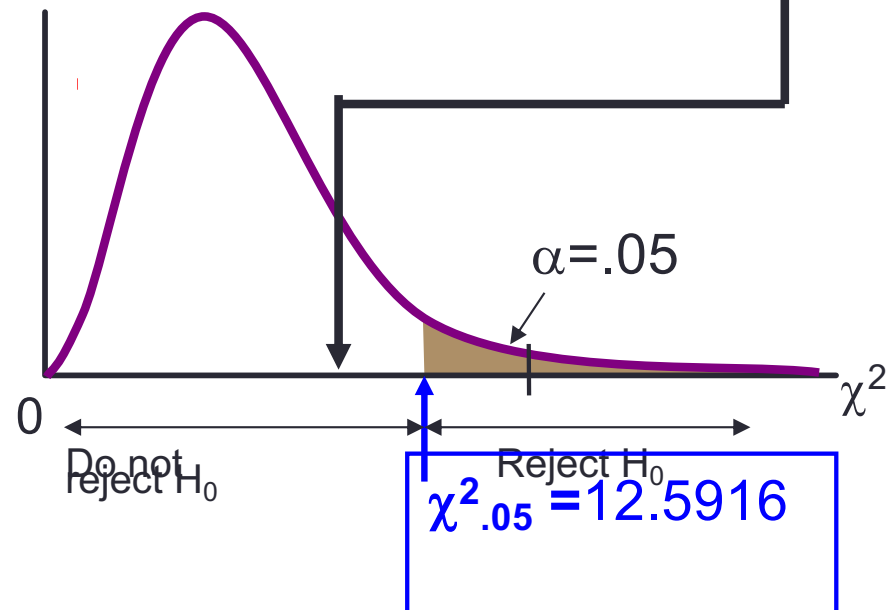
$$\chi^2 = \sum_k \frac{(f_o - f_e)^2}{f_e} = \frac{(10 - 13.68)^2}{13.68} + \dots + \frac{(9 - 3.41)^2}{3.41} = 11.580$$

7 classes so use  $7 - 1 = 6$  d.f.:

$$\chi^2_{.05} = 12.5916$$

## Conclusion:

$\chi^2 = 11.580 < \chi^2_{\alpha} = 11.0705$  so not **reject**  $H_0$ . There is no sufficient evidence that the data are not normal with  $\mu = 50$  and  $\sigma = 15$





# Exercise

- Write a program to perform a test: Whether the following measurements from a production process follow a **normal distribution** with  $\mu = 50$  and  $\sigma = 15$ ?

```
57 66 54 55 66 70 48 28 37 53 50 71 24 51 53 53 43 29 46 46
27 64 47 42 52 45 62 51 84 32 60 61 33 50 36 73 41 64 59 40
59 77 56 49 66 29 74 44 20 72 38 73 41 27 55 24 43 47 49 52
75 36 76 59 43 32 44 29 39 78 72 38 42 51 39 53 41 29 32 62
72 42 32 53 35 55 42 107 60 63 46 68 55 51 37 50 50 65 33 53
38 30 28 55 40 39 38 69 57 36 60 17 67 62 39 32 39 47 51 88
33 74 47 43 33 56 47 42 33 63 49 59 26 40 50 48 31 65 14 57
82 52 87 46 58 48 35 36 57 42
```

The data is assigned to 7 classes:

```
x<-c(57,66,54,55,...)
```

less than 30

30 but < 40

40 but < 50

50 but < 60

60 but < 70

70 but < 80

80 or over