

## Chapter 12: Inference on Categorical Data

$$\sum P(x) = 1$$

$$0 \leq P(x) \leq 1$$

Review:

$\bar{X}, x$

- Discrete Random Variables: Probability Distributions, Expected Value, Binomial Distribution

- 6.1 Discrete Random Variables and their Probability Distributions

- 6.1 Expected Value of Probability Distribution

$$E(X) = \mu = \sum [x \cdot P(x)] \quad \text{key "long run behavior"}$$

- 6.2 Binomial Probability Distribution

- 4 requirements: ① fixed trials  $n$  ② trials independent

$$\text{Expected Value: } E(X) = \mu = n \cdot p \quad (\sum E_i = 1)$$

Carrying over Dice

$x$	$P(x)$	$x$	$P(x)$
-1	0.5	1	1/6
-5	0.2	2	1/6
16	0.2	3	1/6
20	0.1	4	1/6
		5	1/6
		6	1/6

③ there are only two (BI) outcomes possible  
④ probability of success is constant.

notation  $p$  = prob. success  
 $q = 1 - p$  = prob. failure.

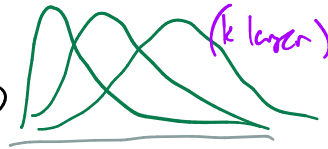
- Chi-Squared Distribution

- 9.3 Chi-Squared Distribution

- Shape: not symmetric, skewed (2)

- Depends on  $df = k$

- Can you find probability under  $\chi^2$ -distribution? What about critical values?



yes,  $\chi^2 \text{cdf}(a, b, df)$

Inv  $\chi^2$ ? No! Table!

### Section 12.1: Goodness-of-Fit Test

In this section, we study a procedure to test hypotheses about a probability distribution. For example, you might want to test if a dice is fair with each side having probability  $p = \frac{1}{6}$ .

Def A **goodness-of-fit test** is used to test the hypothesis that an observed frequency distribution fits (or conforms to) some claimed distribution.

The M&M company claims that the distribution of plain M&M candies in a bag is 23% blue, 23% orange, 15% green, 12% red, 15% yellow, and 12% brown. Even though this is their claim, do you think this represents the true proportions of color distribution in all of the M&M bags? **How would we check?**

claimed prop.  $\rightarrow$  sample data

$H_0: p_1 = 0.23, p_2 = 0.23, p_3 = 0.15, p_4 = 0.12, p_5 = 0.15, p_6 = 0.12$

$H_A: \text{At least one of the proportions differs from the claimed proportion.}$

(Blue) (Orange) (Green) (Red) (Yellow) (Brown)

### What to Compare and How to Compare It

Expected Counts ( $E_i$ )	Observed Counts ( $O_i$ )
The number in each category we would expect to see if $H_0$ is true.	Observe how many in your <u>sample</u> are in each category.
* Two ways of calculating the expected counts:	* This information will be given.
1. If the expected counts are EQUAL, then $E = \frac{n}{k}$	
<ul style="list-style-type: none"> <li>where <math>n</math> = total observed</li> <li>where <math>k</math> = # of categories</li> </ul>	
2. If the expected counts are not equal, then calculate using $E_i = \mu_i = n_i \cdot p_i$ where $i = 1, 2 \dots k$	

\*If the observation and experiment counts are "close", then Fail to Reject  $H_0$

(keep  $H_0$ )

LOGIC of HT

If the observation and experiment counts "far apart" (BIG DIFFERENCE), then Reject  $H_0$

EX 1: Finding the expected counts.

(a) A single die is rolled 45 times with the following results. Assuming that the die is fair and all outcomes are equally likely, find the expected frequency  $E$  for each empty cell.

$\hookrightarrow E$  are all same!

Outcome	1	2	3	4	5	6
Observed $O_i$	13	6	12	9	3	2
Expected $E_i$	7.5	7.5	7.5	7.5	7.5	7.5

$$\mu = E = \frac{n}{k} = \frac{\text{total observations}}{\# \text{ categories}} = \frac{45}{6} = 7.5$$

other formula:  $E = n \cdot p = 45 \left( \frac{1}{6} \right) = 7.5$   
 $\text{fair}$

(b) Jon works as an usher at a theatre. The theatre has 1000 seats that are accessed through five entrances. Each guest should use the entrance that's marked on their ticket. Entrances A and B should each have 30% of the guests using these entrances. Entrance C should have 20% of the guests using its entrance. Entrances D and E should each have 10% of the guests using these entrances. Find the expected frequency for each  $E$  for each entrance.

Entrance	A	B	C	D	E
Observed $O_i$	398	202	205	87	108
Expected $E_i$	300	300	200	100	100

Note  
 $\sum E_i = n$

$(i=1,2)$   $E_i = n_i \cdot p_i$   
 $= 1000 \cdot 0.3$   
 $= 300$

$(i=3)$   
 $E_i = n_i \cdot p_i$   
 $= 1000 \cdot (0.20)$   
 $= 200$

$E_i = n_i \cdot p_i$  ( $i=3,4$ )  
 $= 1000 \cdot (0.10)$   
 $= 100$

### Steps for Hypothesis Test for Goodness-of-Fit

#### What to Find...

- Number of categories,  $k$
- Expected Counts,  $E_i$

#### Check Requirements

- The data has to be randomly selected.
- The sample data consist of frequency counts for each of the different categories. (none missing)
- For each category, the expected frequency is at least 5.

#### Step 1: Hypotheses

$H_0: p_1 = p_2 = \dots = p_k$  (EQUAL OUTCOMES - DICE)  
 $H_A: \text{At least one of the probabilities is different from the others}$   
 $p_i \neq \#$

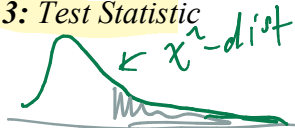
OR

$H_0: p_1 = \# \ p_2 = \#, \dots, p_k = \#$  (Potentially different - M & M's)  
 $H_A: \text{At least one of the probabilities is different the claimed distribution}$

#### Step 2: Level of Significance

$\alpha = P(\text{Type I Error})$   
 if not given, assume 0.05

#### Step 3: Test Statistic



$$\chi_0^2 = \sum \frac{(O-E)^2}{E}$$

ALWAYS RIGHT-TAILED TEST!

Note: To compute the test statistic you will need to use lists on your calculator!

Calc  $L1 = O_i$   $L2 = E_i$   $L3 = \frac{(L1 - L2)^2}{L2}$  Test statistic  $\chi_0^2 = \text{sum}(L3)$

Stat > Edit > enter: L1 = O, L2 = E, L3 = (L1-L2)^2/L2, then  $\chi^2_0 = \text{sum}(L3)$

$k = \# \text{ of categories}$

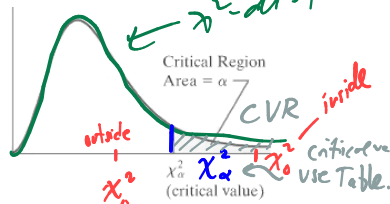
### Steps for Hypothesis Test for Goodness-of-Fit (Cont.)

**Step 4:** Find a Critical Value or P-Value to check either using the Critical Value Method or P-Value Method.

#### CRITICAL REGION METHOD

\* Table VIII

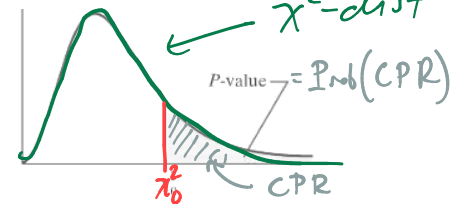
$$df = k - 1$$



- $\left\{ \begin{array}{l} \text{Reject } H_0 \text{ if } \chi^2 \text{ lies in the critical region} \\ \text{Fail to Reject } H_0 \text{ if } \chi^2 \text{ doesn't lie in the critical region} \end{array} \right.$

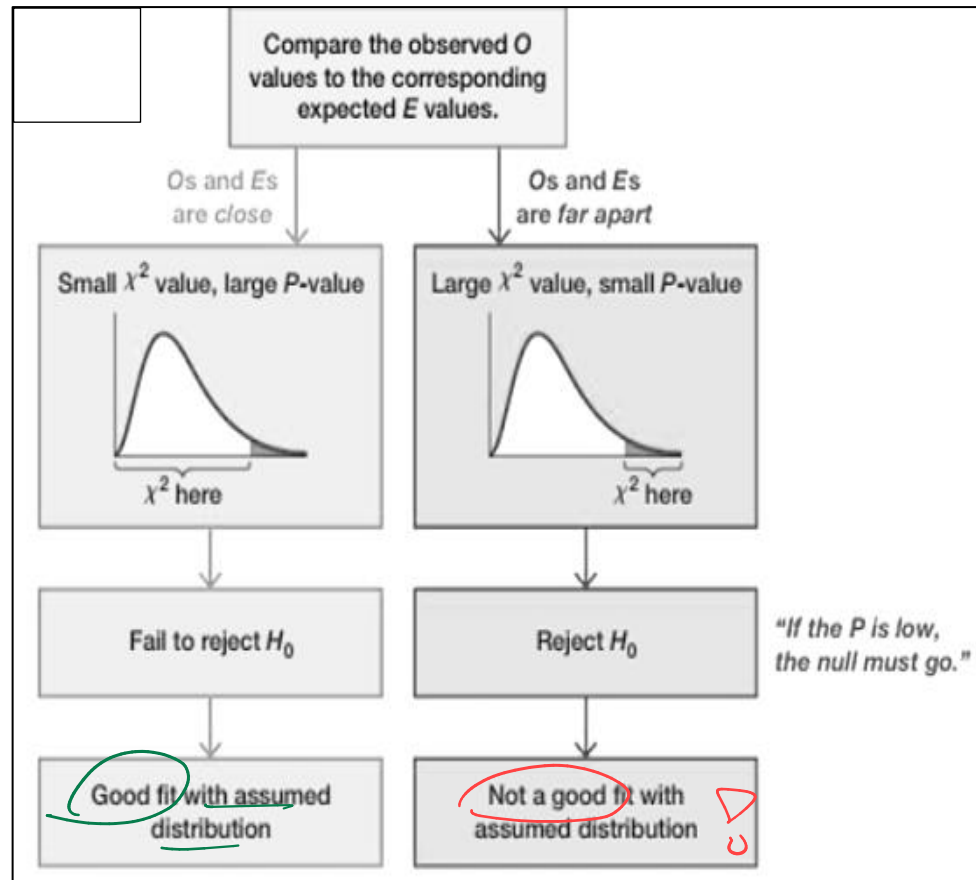
#### P-VALUE METHOD

$$df = k - 1$$



- $\left\{ \begin{array}{l} \text{Reject } H_0 \text{ if } P\text{-value} \leq \alpha \\ \text{Fail to Reject } H_0 \text{ if } P\text{-value} > \alpha \end{array} \right.$

**Step 5:** Make a decision and draw a conclusion.



### GRAPHING CALCULATOR (TI-83 OR 84)

Instructions:

STAT  $\Rightarrow$  TESTS  $\Rightarrow$  D:  $\chi^2$  GOF -Test

$k = 6 = \# \text{ of colors}$

$$L3 = (L1 - L2)^2$$

Ex 2: The M&M company claims that the distribution of plain M&M candies in a bag is 23% blue, 23% orange, 15% green, 12% red, 15% yellow, and 12% brown. Suppose we took a simple random sample of 400 M&Ms from the populations of all M&Ms. The results are shown below:

COLOR	Blue $p_1$	Orange $p_2$	Green $p_3$	Red $p_4$	Yellow $p_5$	Brown $p_6$
FREQUENCY	53	66	38	96	88	59
EXPECTED	92	92	60	48	60	48

$400 \cdot 0.23 \quad 400 \cdot 0.23 \quad 400 \cdot 0.15 \quad 400 \cdot 0.12 \quad 400 \cdot 0.15 \quad 400 \cdot 0.12$

Find  $E_i = n p_i$

Is the proportion of each color different than the claim of the M&M's manufacturer?

Null and Alternative Hypothesis

$$\begin{cases} H_0: p_1 = 0.23, p_2 = 0.23, p_3 = 0.15, p_4 = 0.12, p_5 = 0.15, p_6 = 0.12 \\ H_A: \text{At least one proportion differs from the claimed proportions.} \end{cases}$$

Test Statistic

$$\chi^2 = \sum \frac{(O - E)^2}{E} = \text{sum}(L3) = 95.53$$

$$L1 = O$$

$$L2 = E$$

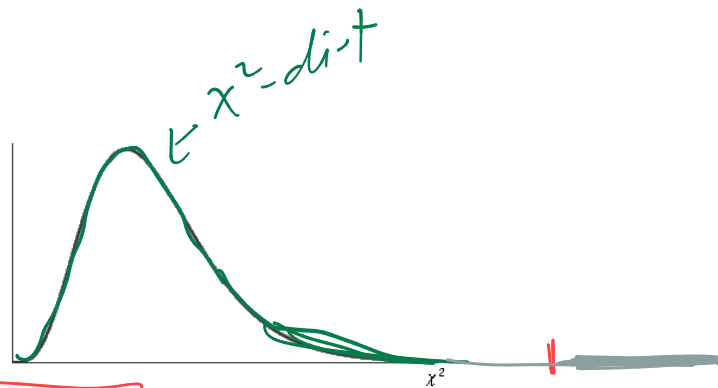
$$L3 = (L1 - L2)^2 / L2$$

P-value/Critical Region

$$P\text{-Val} = P(\chi^2 > 95.53)$$

$$= \chi^2 \text{cdf}(\text{low}, \text{high}, \text{df})$$

$$= \chi^2 \text{cdf}(95.53, 1E99, 5) = 4.62 \times 10^{-19} = 0.000000000000000000462 = 0+$$



Decision about Null Hypothesis

$$\alpha = 0.05$$

$$P = 0.000000000000000000462$$

$$P < \alpha \rightarrow P \text{ low, Null go!}$$

Reject  $H_0$

Conclusion

"There is enough statistical evidence to support the claim that at least one proportion of M&M does differs from the manufacturer's claimed proportion!"

(Calc)

Ex 3: A company sells their products exclusively by mail. The company's management wants to find out if the number of orders received at the company's office on each of the five days of the week is the same. The company took a random sample of 400 orders received during a four-week period. The following table lists the frequency distribution for these orders by the day of the week.

$k = 5$  (# categories)

	Monday	Tuesday	Wednesday	Thursday	Friday
Number of Orders	92	71	65	83	89
Expected Number	$400 \cdot (1/5)$ 80	80	80	80	80

$E = n \cdot p$   
 $400(1/5)$

Test the claim that the orders are evenly distributed over the five days of the week. Use  $\alpha = .025$

LOS

Null and Alternative Hypothesis

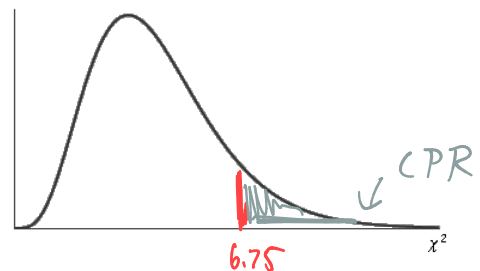
$$\begin{cases} H_0: p_1 = p_2 = p_3 = p_4 = p_5 \text{ "the proportion of orders is the same"} \\ H_A: \text{At least one is different.} \end{cases}$$

Test Statistic

$$\chi^2 = 6.75$$

P-value/Critical Region

$$P = 0.150$$



Decision about Null Hypothesis

$$\alpha = 0.025$$

$P > \alpha \rightarrow P \text{ high, Null!}$

Fail to Reject  $H_0$

$$P = 0.150$$

Conclusion

"There isn't enough statistical evidence to support the claim that at least one day of the week receives a different # of orders."