#### Math 11 Calculus-Based Introductory Probability and Statistics

Eddie Aamari S.E.W. Assistant Professor

eaamari@ucsd.edu math.ucsd.edu/~eaamari/ AP&M 5880A

Today:

•  $\chi^2$  tests

## The $\chi^2$ test

The  $\chi^2$  distribution is used to compare counts in a table (to a list of expected values). We will study two cases where you use the  $\chi^2$ distribution:

- Goodness-of-fit (one-way table)
- Test of Homogeneity/Independence

(two-way table)

## The $\chi^2$ test

The  $\chi^2$  distribution is used to compare counts in a table (to a list of expected values). We will study two cases where you use the  $\chi^2$  distribution:

- Goodness-of-fit (one-way table)
- Test of Homogeneity/Independence

(one-way table) (two-way table)

Two goodness-of-fit questions:

- 1. If we look at the birth months of NHL players, do they resemble what we might see in the larger U.S. population?
- 2. If we breed a bunch of peas, do we really get the results expected from Mendel's theory of genetics?

#### The Hockey Conundrum

¢ Month	‡ Birth in US	% <sup>‡</sup> Birth in US	‡ NHL Players	<pre>% NHL <sup> ‡ </sup> Players </pre>
January	319297	8.0	773	10.3
February	299235	7.5	708	9.5
March	335786	8.4	737	9.9
April	308809	7.7	708	9.5
May	334437	8.3	685	9.2
June	336251	8.4	600	8.0
July	347934	8.7	619	8.3
August	362798	9.1	533	7.1
September	350711	8.7	570	7.6
October	347354	8.7	557	7.4
November	330832	8.3	495	6.6
December	335111	8.4	494	6.6

#### The Hockey Conundrum

÷ Month	‡ Birth in US	% <sup>‡</sup> Birth in US	‡ NHL Players	∲ NHL Players
January	319297	8.0	773	10.3
February	299235	7.5	708	9.5
March	335786	8.4	737	9.9
April	308809	7.7	708	9.5
May	334437	8.3	685	9.2
June	336251	8.4	600	8.0
July	347934	8.7	619	8.3
August	362798	9.1	533	7.1
September	350711	8.7	570	7.6
October	347354	8.7	557	7.4
November	330832	8.3	495	6.6
December	335111	8.4	494	6.6

A  $\chi^2$  test is basically comparing all the proportions above (although we end up using counts instead).

In this sense, it generalizes the study of proportions.

#### It's All About the Counts

Observed counts are the values you see in your data.

(These will be integers)

Expected counts are the values you expect if you apply your theory to the data set size you are exploring.

(These can be fractions)

\$ Month	Observed Counts	Expected Counts
January	773	595.7314
February	708	558.3006
March	737	626.4960
April	708	576.1634
May	685	623.9790
June	600	627.3635
July	619	649.1612
August	533	676.8939
September	570	654.3424
October	557	648.0791
November	495	617.2530
December	494	625.2366

 $\leftarrow$  In the U.S., 8% of birth are in January. So, of the 7479 NHL players, we expect

 $7479 \times 0.08 \simeq 595$ 

players born in January.

#### It's All About the Counts

Observed counts are the values you see in your data.

(These will be integers)

Expected counts are the values you expect if you apply your theory to the data set size you are exploring.

(These can be fractions)

\$ Month	Observeð Counts	Expected Counts
January	773	595.7314
February	708	558.3006
March	737	626.4960
April	708	576.1634
May	685	623.9790
June	600	627.3635
July	619	649.1612
August	533	676.8939
September	570	654.3424
October	557	648.0791
November	495	617.2530
December	494	625.2366

 $\leftarrow$  In the U.S., 8% of birth are in January. So, of the 7479 NHL players, we expect

 $7479 \times 0.08 \simeq 595$ 

players born in January.

In general, to find an expected count, multiply your data size by the theorized percentage of that cell.

We set up hypotheses as we would for any test:

- $H_0\colon$  The birth pattern of NHL players (by month) matches the larger U.S. population
- $H_A :$  There is something different in the birth patterns of NHL players and Americans

We set up hypotheses as we would for any test:

- $H_0\colon$  The birth pattern of NHL players (by month) matches the larger U.S. population
- ${\cal H}_A :$  There is something different in the birth patterns of NHL players and Americans

We need a statistic that captures how strange our data is under  $H_0$ .

$$\chi^2 = \frac{(O_1 - E_1)^2}{E_1} + \frac{(O_2 - E_2)^2}{E_2} + \dots + \frac{(O_k - E_k)^2}{E_k}$$

 $O_i$ : Observed count for category *i*  $E_i$ : Expected count for category *i* 

We set up hypotheses as we would for any test:

- $H_0\colon$  The birth pattern of NHL players (by month) matches the larger U.S. population
- ${\cal H}_A :$  There is something different in the birth patterns of NHL players and Americans

We need a statistic that captures how strange our data is under  $H_0$ .

$$\chi^2 = \frac{(O_1 - E_1)^2}{E_1} + \frac{(O_2 - E_2)^2}{E_2} + \dots + \frac{(O_k - E_k)^2}{E_k}$$

 $O_i$ : Observed count for category i $E_i$ : Expected count for category i

If  $H_0$  is true, what value should  $\chi^2$  be near?

We set up hypotheses as we would for any test:

- $H_0\colon$  The birth pattern of NHL players (by month) matches the larger U.S. population
- ${\cal H}_A :$  There is something different in the birth patterns of NHL players and Americans

We need a statistic that captures how strange our data is under  $H_0$ .

$$\chi^2 = \frac{(O_1 - E_1)^2}{E_1} + \frac{(O_2 - E_2)^2}{E_2} + \dots + \frac{(O_k - E_k)^2}{E_k}$$

 $O_i$ : Observed count for category *i*  $E_i$ : Expected count for category *i* 

If  $H_0$  is true, what value should  $\chi^2$  be near?

Under the null, we expect to have 
$$O_i \simeq E_i$$
,  
so  $\chi^2$  should be small (close to zero).

Anth +	Observed Counts	Expected Counts
January	773	595.7314
February	708	558.3006
March	737	626.4960
April	708	576.1634
May	685	623.9790
June	600	627.3635
July	619	649.1612
August	533	676.8939
September	570	654.3424
October	557	648.0791
November	495	617.2530
December	494	625.2366
	7479	7479.0000



$$\chi^2 = \frac{(773 - 595.73)^2}{595.73} + \frac{(708 - 558.30)^2}{558.30} + \ldots + \frac{(494 - 625.23)^2}{625.23} \simeq 257.13$$

The bigger  $\chi^2$  is, the more evidence we have against  $H_0$ .

#### Using $\chi^2$ to Find a *p*-Value

Getting a value for  $\chi^2$  is like getting a z-statistic or t-statistic. Now we need to know what sampling distribution it lives on so we can find an area.

#### Using $\chi^2$ to Find a *p*-Value

Getting a value for  $\chi^2$  is like getting a z-statistic or t-statistic. Now we need to know what sampling distribution it lives on so we can find an area.

When you analyse a null hypothesis across k categories (k = 12 months), the curve is  $\chi^2_{k-1}$ , the Chi-squared distribution with k-1 degrees of freedom.



For our hockey example, we want to know if a value 257.13 is strange to see on  $\chi^2_{12-1} = \chi^2_{11}$ .



For our hockey example, we want to know if a value 257.13 is strange to see on  $\chi^2_{12-1} = \chi^2_{11}$ .



We find the area under the curve to the right of 257.13. Here, the area is so small, the *p*-value is basically 0 ( $p < 2.2 \cdot 10^{-16}$ )

Since p < 0.05, we reject the null hypothesis. We are quite (!) sure something is strange about the birth month counts of hockey players. Summary of the  $\chi^2$  Goodness-of-Fit Test

- 1. You have a collection of counts for some phenomenon (birth of NHL players) divided by categories (months). You wish to compare these counts to those predicted by some theory.
- 2. Calculate the expected counts from your theory (fractional counts allowed).

3. Calculate 
$$\chi^2 = \sum_{\text{category}} \frac{(O_{\text{category}} - E_{\text{category}})^2}{E_{\text{category}}}$$

- 4. Look up this value on the curve  $\chi^2_{k-1}$ , where k is the number of categories.
- 5. Use the *p*-value to decide about  $H_0$ .

#### Mendel, Genetics, and Plunnett Squares



R/r is an allele that codes roundness

Y/y is an allele that codes for yellowness

R and Y are dominant r and y are recessive

In the cross, we get many genotypes (genetic combos), which results in just 4 phenotypes (visual traits)

Theory predicts the phenotypes occur in the ratios 9:3:3:1.





We must find the expected counts. A 9:3:3:1 ratio means  $\frac{9}{9+3+3+1} = \frac{9}{16}$  of the peas will be of the most prominent phenotype.



We must find the expected counts.

A 9:3:3:1 ratio means  $\frac{9}{9+3+3+1} = \frac{9}{16}$  of the peas will be of the most prominent phenotype.

We find the expected counts. (for instance,  $\frac{9}{16} \times (62 + 24 + 9 + 5) = 56.25$  for the first line)



We must find the expected counts.

A 9:3:3:1 ratio means  $\frac{9}{9+3+3+1} = \frac{9}{16}$  of the peas will be of the most prominent phenotype.

We find the expected counts. (for instance,  $\frac{9}{16} \times (62 + 24 + 9 + 5) = 56.25$  for the first line)

And we calculate the test statistic.

$$\chi^2 = \frac{(62 - 56.25)^2}{56.25} + \frac{(24 - 18.75)^2}{18.75} + \frac{(9 - 18.75)^2}{18.75} + \frac{(5 - 6.25)^2}{6.25} \simeq 7.378$$

We go look up the area to the right of the value  $\chi^2 = 7.378$  under the curve of the distribution  $\chi^2_{4-1} = \chi^2_3$ .



We go look up the area to the right of the value  $\chi^2 = 7.378$  under the curve of the distribution  $\chi^2_{4-1} = \chi^2_3$ .



Figure B.2: Areas in the chi-square table always refer to the right tail.

Upper	tail	0.3	0.2	0.1	0.05	0.02	0.01	0.005	0.001
df	1	1.07	1.64	2.71	3.84	5.41	6.63	7.88	10.83
	<b>2</b>	2.41	3.22	4.61	5.99	7.82	9.21	10.60	13.82
	3	3.66	4.64	6.25	7.81	9.84	11.34	12.84	16.27
	4	4.88	5.99	7.78	9.49	11.67	13.28	14.86	18.47
	<b>5</b>	6.06	7.29	9.24	11.07	13.39	15.09	16.75	20.52
	6	7.23	8.56	10.64	12.59	15.03	16.81	18.55	22.46
	7	8.38	9.80	12.02	14.07	16.62	18.48	20.28	24.32

We go look up the area to the right of the value  $\chi^2 = 7.378$  under the curve of the distribution  $\chi^2_{4-1} = \chi^2_3$ .



From the table, we see that  $\chi^2 = 7.378$  falls between the values 6.25 (associated to 0.1) and 7.81 (associated to 0.05).

We go look up the area to the right of the value  $\chi^2 = 7.378$  under the curve of the distribution  $\chi^2_{4-1} = \chi^2_3$ .



From the table, we see that  $\chi^2 = 7.378$  falls between the values 6.25 (associated to 0.1) and 7.81 (associated to 0.05).

Hence, the *p*-value satisfies 0.05 .

# Performing a $\chi^2$ Test Minitab



 $H_0$ : Our data follow Mendel's theory  $H_A$ : Our data contradict Mendel's theory.

With a *p*-value of 6.1%, we do not reject  $H_0$ . Our data, while strange, may be strange because of random variations.

 $H_0$ : Our data follow Mendel's theory  $H_A$ : Our data contradict Mendel's theory.

With a *p*-value of 6.1%, we do not reject  $H_0$ . Our data, while strange, may be strange because of random variations.



In essence, a  $\chi^2$  goodness-of-fit test tells you if the observed data differ from the data expected based on some theory by an amount that exceeds the type of variation we expect from random effects.

Said differently:  $\chi^2$  tells you if the two histograms are "close enough" to each other (null) or not (alternative).

#### But Wait! When Can We Use This Framework?

To use a  $\chi^2$  framework to test for goodness-of-fit, we must meet three conditions:

- You start with a one-dimensional  $(k \times 1)$  table of **observed** counts. You wish to compare these counts to those predicted by some theory.
- The **counts** in the cells of the table must be <u>independent</u> of one another. Randomly sampling the people that <u>comprise</u> these counts usually gives you this.
- The **expected count** for each cell must be at least 5. (Note: we don't require that the observed counts be at least 5, just the expected counts)

#### Detecting Game Cheating With $\chi^2$ Test

A gamer uses  $\chi^2$  goodness-of-fit to provide evidence that something is fishy in another gamer's uploaded video on Super Smash Bros:

#### Link

	Video	Data	p-value	Significant
<ul> <li>1 clip that starts with a 0 in the seconds' tens digit</li> <li>1 with a 1</li> <li>1 with a 2</li> </ul>	A Silly Combo Video	1, 7, 3, 8, 3, 2	0.08	No
	I Killed Mufasa	13, 9, 7, 6, 6, 9	0.55	No
	Silence	10, 6, 8, 10, 9, 9	0.95	No
<ul> <li>2 with a 3</li> <li>2 with a 4</li> </ul>	The Game is not Over	14, 10, 11, 17, 6, 15	0.28	No
<ul> <li>12 (yes, twelve) with a 5</li> </ul>	Version 2.0	4, 9, 4, 3, 9, 9	0.25	No
	510 Evolution: Darrell	7, 4, 2, 3, 7, 8	0.34	No
	600 Hours	1, 1, 1, 2, 2, 12	0.00006	Yes

#### What Else is $\chi^2$ Used For?

Goodness-of-Fist Test: **1 population** (NHL players, peas) split across a **categorical variable** (birth month, phenotype). You have a 1-dimensional table of counts.

Test for Independence: **2 or more populations** split across a **categorical variable**.

You have a 2-dimensional table of counts.

	No tumor	One tumor	Two or more tumors	Total
Control	74	5	1	80
Low dose	63	12	5	80
High dose	45	15	10	70
Total	182	32	16	230

Here, we have

- 3 different population of rats
- A categorical variable "Tumor development"

#### What Happens When Rats Smoke?

Pyrobenzene is a major component of cigarette smoke.

Researchers injected rats with different levels of pyrobenzene and looked for tumor development.

Is there an association between tumor development and pyrobenzene dosage?

	No tumor	One tumor	Two or more tumors	Total
Control	74	5	1	80
Low dose	63	12	5	80
High dose	45	15	10	70
Total	182	32	16	230

#### What Happens When Rats Smoke?

Pyrobenzene is a major component of cigarette smoke.

Researchers injected rats with different levels of pyrobenzene and looked for tumor development.

Is there an association between tumor development and pyrobenzene dosage?

	No tumor	One tumor	Two or more tumors	Total
Control	74	5	1	80
Low dose	63	12	5	80
High dose	45	15	10	70
Total	182	32	16	230

Clearly, the is an association in our sample. The question is: Is there an association in the larger population? Might our results be due to chance?

 $H_0$ :There is no association between pyrobenzene<br/>and tumor development in rats. $H_A$ :There is an association between pyrobenzene<br/>and tumor development in rats.

If there is <u>NO</u> association, then whatever **results we see on av**erage for a given tumor count should apply consistently to the different pyrobenzene levels.

	No tumor	One tumor	Two or more tumors	Total
Control	74	5	1	80
Low dose	63	12	5	80
High dose	45	15	10	70
Total	182	32	16	230

If there is <u>NO</u> association, then whatever **results we see on av**erage for a given tumor count should apply consistently to the different pyrobenzene levels.

	No tumor	One tumor	Two or more tumors	Total
Control	74 (63.3)	5	1	80
Low dose	63	12	5	80
High dose	45	15	10	70
Total	182	32	16	230

– There are 80 control rats. We expect no tumor to develop in 182/230 rats. Thus, we expect  $80 \times (182/230)$  control rats to have no tumor.

If there is <u>NO</u> association, then whatever **results we see on av**erage for a given tumor count should apply consistently to the different pyrobenzene levels.

	No tumor	One tumor	Two or more tumors	Total
Control	74 (63.3)	5	1	80
Low dose	63	12	5	80
High dose	45	15 (9.74)	10	70
Total	182	32	16	230

– There are 80 control rats. We expect no tumor to develop in 182/230 rats. Thus, we expect  $80 \times (182/230)$  control rats to have no tumor.

– There are 70 high-dose rats. We expect one tumor to develop in 32/230 rats. Thus, we expect  $70 \times (32/230)$  control rats to have one tumor.

If there is <u>NO</u> association, then whatever **results we see on average for a given tumor count** should apply consistently to the **different pyrobenzene levels**.

	No tumor	One tumor	Two or more tumors	Total
Control	74 (63.3)	5	1	80
Low dose	63	12	5	80
High dose	45	15 (9.74)	10	70
Total	182	32	16	230

– There are 80 control rats. We expect no tumor to develop in 182/230 rats. Thus, we expect  $80 \times (182/230)$  control rats to have no tumor.

– There are 70 high-dose rats. We expect one tumor to develop in 32/230 rats. Thus, we expect  $70 \times (32/230)$  control rats to have one tumor.

In general, the expected count for a cell is

 $\text{Expected}_{\text{cell}} = \frac{(\text{Row total}) \times (\text{Column total})}{\text{Table total}}.$ 

#### The $\chi^2$ Value for Testing Independence

 $\text{Expected}_{\text{cell}} = \frac{(\text{Row total}) \times (\text{Column total})}{\text{Table total}}.$ 

As before,

$$\chi^{2} = \sum_{\text{all cells}} \frac{(\text{Observed}_{cell} - \text{Expected}_{cell})^{2}}{\text{Expected}_{cell}}$$

This  $\chi^2$  value lives on the  $\chi^2_{df}$  distribution with df = (r-1)(c-1), where the two-way table has size  $r \times c$ .

	No tumor	One tumor Two or more tumors		Total
Control	74	5	1	80
Low dose	63	12	5	80
High dose	45	15	10	70
Total	102	32	16	230
		,	Our table is	$3 \times 3$

**Important:** Ignore the total column/row to determine r and c! (They're displayed for readability, but they're not part of data)

#### Doing This in Minitab



Names	No tumor	1 tumor	2+ tumors	
Control	74	5	1	
Low Dose	63	12	5	
High Dose	45	15	10	

#### Doing This in Minitab

Chi-Square Test for Ass	ociation ×
	Summarized data in a two-way table
	Columns containing the table:
	'No tumor'-'2+ tumors'
	~
	Labels for the table (optional)
	Rows: Names (column with row labels)
	Columns: (name for column category)
Select	Statistics Options
Help	QK Cancel

The probability of getting our observed data (top values in cells) in a world with no association (bottom values in cells) is 0.001. Our data would be a 1 in a 1000 event! This is so strange, we reject the null hypothesis.

#### Chi-Square Test for Association: Names, Worksheet columns



Without Minitab, you compute the value of the test statistic  $\chi^2 = 19.25$ , and you go look up the area to the right of this value under the curve  $\chi^2_{df}$ , where df = (3-1)(3-1) = 4.

Without Minitab, you compute the value of the test statistic  $\chi^2 = 19.25$ , and you go look up the area to the right of this value under the curve  $\chi^2_{df}$ , where df = (3-1)(3-1) = 4.



Without Minitab, you compute the value of the test statistic  $\chi^2 = 19.25$ , and you go look up the area to the right of this value under the curve  $\chi^2_{df}$ , where df = (3-1)(3-1) = 4.



Figure B.2: Areas in the chi-square table always refer to the right tail.

Upper	tail	0.3	0.2	0.1	0.05	0.02	0.01	0.005	0.001
df	1	1.07	1.64	2.71	3.84	5.41	6.63	7.88	10.83
	2	2.41	3.22	4.61	5.99	7.82	9.21	10.60	13.82
	3	3.66	4.64	6.25	7.81	9.84	11.34	12.84	16.27
	4	4.88	5.99	7.78	9.49	11.67	13.28	14.86	18.47
	5	6.06	7.29	9.24	11.07	13.39	15.09	16.75	20.52
	6	7.23	8.56	10.64	12.59	15.03	16.81	18.55	22.46
	7	8.38	9.80	12.02	14.07	16.62	18.48	20.28	24.32

Without Minitab, you compute the value of the test statistic  $\chi^2 = 19.25$ , and you go look up the area to the right of this value under the curve  $\chi^2_{df}$ , where df = (3-1)(3-1) = 4.



From the table, we see that  $\chi^2 = 19.25$  falls after the value 18.47 (associated to 0.001).

Without Minitab, you compute the value of the test statistic  $\chi^2 = 19.25$ , and you go look up the area to the right of this value under the curve  $\chi^2_{df}$ , where df = (3-1)(3-1) = 4.



From the table, we see that  $\chi^2 = 19.25$  falls after the value 18.47 (associated to 0.001).

Hence, the *p*-value satisfies p < 0.001.

In particular, p < 0.05, so we reject the null. It appears that there is an association between pyrobenzene and tumor development in rats.

#### What Conditions Must Be Met to Use These Tools?

To use a  $\chi^2$  to determine if there is an association between the rows and columns of a two-way table, we must meet the same conditions as in the one-way table goodness-of-fit test:

- We must have a table of **observed counts**.
- The data making up the cells must be **independent**.
- The **expected cell counts** must be 5 or more.

#### Helpful Observation:

If you have two **quantitative** variables, you can measure the strenght of an association using a correlation coefficient.

If you have two **qualitative** (categorical) variables, you can use a  $\chi^2$  test for the significance of an association.

#### What is This Test Really Doing?

The  $\chi^2$  test for association is looking to see how these distributions look relative to one another.



If pyrobenzene level and tumor count have no association, these three distributions should look nearly identical.

The  $\chi^2$  test is basically measuring how different these three distributions look.