

CS 2810 April 5 Day 20

chi square "goodness of fit"

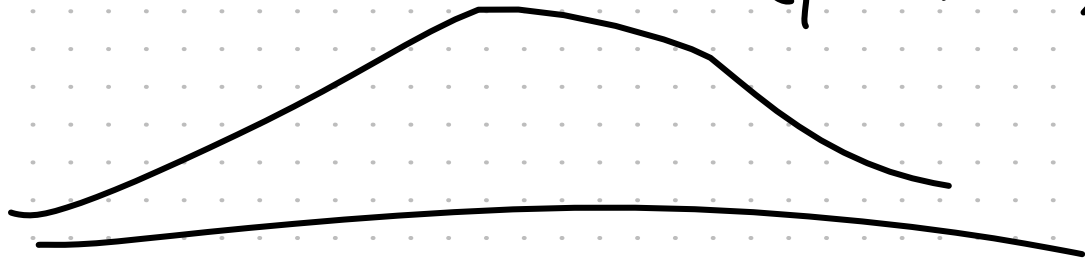
Admin:

- short class
- reschedule quiz4 to may 3? (with section 2)
 - Ell Hall AUD (8-10am)
 - please email me if you can't make it

STANDARD Normal Distribution

Z IS OFTEN USED AS NAME OF SAMPLE FROM THIS DISTRIBUTION

$$Z \sim N(\mu=0, \sigma^2=1)$$



IF $X \sim N(\mu, \sigma^2)$ THEN $\frac{X-\mu}{\sigma} \sim N(\mu=0, \sigma^2=1)$

"Z-SCORE"

Let X be a normally distributed random variable with mean 7 and variance 10. Identify the linear function of X so that it has a "standard" normal distribution (mean 0 and variance 1).

$$X \sim N(\mu=7, \sigma^2=10)$$

$$X_0 = 14$$

OUTCOME



$$Z = \frac{X - \mu}{\sigma} = \frac{X - 7}{\sqrt{10}} = \frac{1}{\sqrt{10}} X - \frac{7}{\sqrt{10}} \sim N(\mu=0, \sigma^2=1)$$

$$Z_0 = \frac{14 - 7}{\sqrt{10}} \approx \frac{7}{3.1} \approx 2.2$$

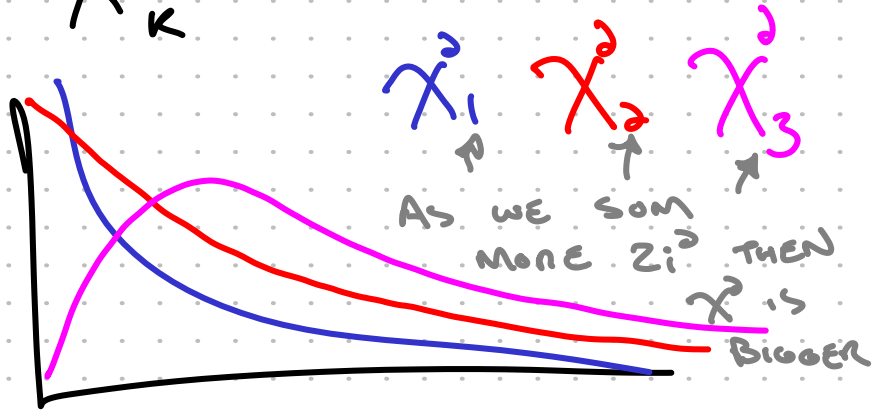
OBSERVATION IS 2 STD DEV ABOVE MEAN

CHI-SQUARE DISTRIBUTION

LET $Z_i \sim N(0, 1)$ BE IID STANDARD NORMAL

CHI-SQUARE IS
SUM OF K
STANDARD
NORMAL
SQUARED

$$\sum_{i=1}^K Z_i^2 \sim \chi^2_K$$



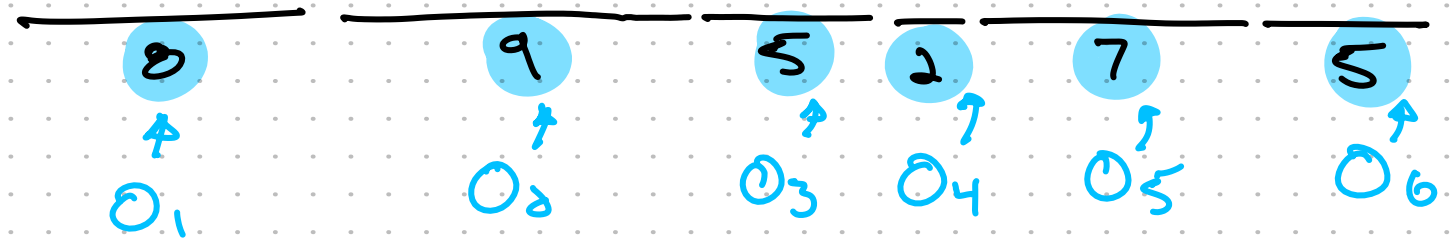
CHI-SQUARE "GOODNESS OF FIT" TEST

ARE THESE OUTCOMES OF A 6-SIDED DIE FAIR?

Do these $N=36$ outcomes come from a fair (uniform) 6-sided die?
(we sort outcomes below so they're easier to work with):

H_0 : PROB OF EACH OUTCOME IS $1/6$
 H_1 : NO ITS NOT

1, 1, 1, 1, 1, 1, 1, 1, 2, 2, 2, 2, 2, 2, 2, 2, 2, 3, 3, 3, 3, 3, 4, 4, 5, 5, 5, 5, 5, 5, 5, 5, 6, 6, 6, 6, 6



O_i How many times did we observe outcome i ?

FOR CHI-SQUARE TEST

H_0 : OBSERVATIONS COME FROM $P(x)$

H_1 : NO THEY DON'T

ASSUME

→ OBSERVATIONS INDEPENDENT

Assuming H_0 (Die is fair in this Ex) what
is expected count of each outcome?

$$E_i = N \cdot p_i = 36 \cdot \frac{1}{6} = 6$$

TOTAL OUTCOMES
(36 IN THIS EX)

PROB OF EACH
OUTCOME UNDER
 H_0

CONTINGENCY TABLE + χ^2 TEST STATISTIC

| | 1 | 2 | 3 | 4 | 5 | 6 |
|--------------------|---|---|---|---|---|---|
| OBSERVED (O_i) | 8 | 9 | 5 | 2 | 7 | 5 |
| EXPECTED (E_i) | 6 | 6 | 6 | 6 | 6 | 6 |

$$\chi^2 = \sum_i \frac{(O_i - E_i)^2}{E_i}$$

| | | | | | | |
|--|---------------------|---------------------|---------------------|---------------------|---------------------|---------------------|
| | $\frac{(8-6)^2}{6}$ | $\frac{(9-6)^2}{6}$ | $\frac{(5-6)^2}{6}$ | $\frac{(2-6)^2}{6}$ | $\frac{(7-6)^2}{6}$ | $\frac{(5-6)^2}{6}$ |
|--|---------------------|---------------------|---------------------|---------------------|---------------------|---------------------|

$$\chi^2 = \frac{4}{6} + \frac{9}{6} + \frac{1}{6} + \frac{16}{6} + \frac{1}{6} + \frac{1}{6} = \frac{32}{6}$$

ICA 1:

- Does there exist a minimum or maximum chi-squared statistic?
- Describe what kind of O_i and E_i achieve this min or max chi-squared statistic.
- Which values of the chi-squared stat are most typical of the null hypothesis? Justify your response with one or two sentences.

| | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---------------|---------------|---------------|----------------|---------------|---------------|
| OBSERVED (O_i) | 8 | 9 | 5 | 2 | 7 | 5 |
| EXPECTED (E_i) | 6 | 6 | 6 | 6 | 6 | 6 |
| $\chi^2 = \sum_i \frac{(O_i - E_i)^2}{E_i}$ | $\frac{4}{6}$ | $\frac{9}{6}$ | $\frac{1}{6}$ | $\frac{16}{6}$ | $\frac{1}{6}$ | $\frac{1}{6}$ |

$\chi^2 = 5.3$

min chi square is 0. it is achieved when the observed count of each outcome equals the expected count of each outcome.

chi-square = 0 is most typical of the null hypothesis
expected counts (E_i) assume the null hypothesis

MODELLING NULL HYPOTHESIS

Assuming the null hypothesis (die is fair) then the chi square statistic follows a chi square distribution with $k = \text{size of sample space} - 1$ degrees of freedom ($df=6-1=5$ in this example)

$$\S.3 = \chi^2 = \sum_i \frac{(E_i - O_i)^2}{E_i} \sim \chi^2_k$$

STATISTIC

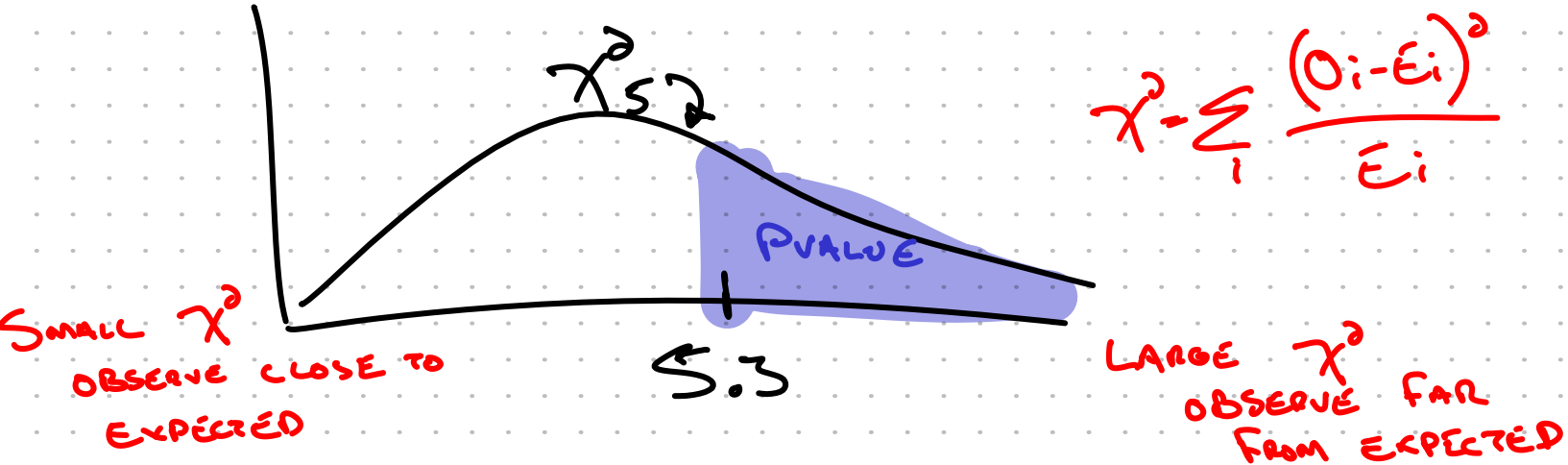
DISTRIBUTION

$k = \# \text{ CATEGORIES} - 1$
IN THIS EXAMPLE:
 $k = 6 - 1 = 5$

Computing P-value with Chi-Squared Goodness of Fit

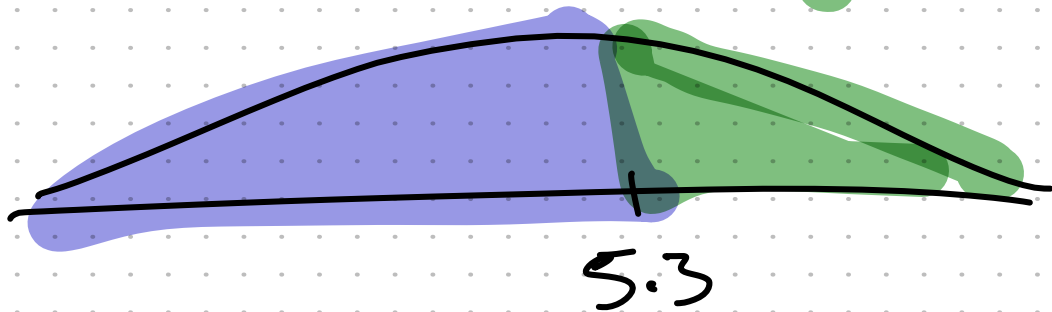
Remember: P-value is prob of all outcomes which are less consistent with null hypothesis

Assuming H_0 , $\chi^2 = 5.3$ is χ^2 DISTRIBUTED WITH $DF=5$



CDF(5.3)

P VALUE =
 $1 - \text{CDF}(5.3) = .38$



$$P_{\text{VAL}} > \alpha = .05$$

DO NOT REJECT H_0
MAKE NO CLAIMS

$$P_{\text{VAL}} < \alpha = .05$$

REJECT H_0 CLAIM H_1

"D.E IS NOT FAIR"

ICA 2

A "silly die" is supposed to roll higher outcomes more often than others:

| OUTCOME | 1 | 2 | 3 | 4 | 5 | 6 |
|---------|--------|--------|--------|--------|--------|--------|
| PROB | $1/21$ | $2/21$ | $3/21$ | $4/21$ | $5/21$ | $6/21$ |

If a die is observed to roll:

1, 1, 1, 2, 2, 2, 3, 3, 3, 3, 4, 4, 4, 4, 4, 4, 4, 4, 5, 5, 5, 5, 6, 6, 6, 6, 6, 6, 6, 6, 6, 6, 6

$\overline{\quad}$ $\overline{\quad}$ $\overline{\quad}$ $\overline{\quad}$ $\overline{\quad}$ $\overline{\quad}$
 3 3 4 7 4 12

perform a hypothesis test to (potentially) claim the die is not "silly".

H_0 : Die is "silly"
 H_1 : Die is NOT "silly"

~~5~~

| OUTCOME | 1 | 2 | 3 | 4 | 5 | 6 | |
|---|-------------------------|-------------------------|---|-------------------------|-------------------------|-------------------------|--|
| PROB | $\frac{1}{21}$ | $\frac{2}{21}$ | $\frac{3}{21}$ | $\frac{4}{21}$ | $\frac{5}{21}$ | $\frac{6}{21}$ | |
| EXPECTED E_i | $\frac{1}{21} \cdot 33$ | $\frac{2}{21} \cdot 33$ | $\frac{3}{21} \cdot 33$ | $\frac{4}{21} \cdot 33$ | $\frac{5}{21} \cdot 33$ | $\frac{6}{21} \cdot 33$ | |
| OBSERVED O_i | 3 | 3 | 4 | 7 | 4 | 12 | $3+3+4+7+4+12$ $= 33$ TOTAL OUTCOMES |
| $\chi^2 = \sum \frac{(E_i - O_i)^2}{E_i}$ | 1.29 | .006 | .108 | .081 | 1.89 | .701 | |
| | $\chi^2 = 4.08$ | | $P\text{-VAL} = 1 - \text{CHID.CDF}(4.08, \text{DF}=5) \approx .53$ DON'T REJECT H_0 | | | | |

CHI-SQUARE "BINNING"

χ^2 REQUIRES A FINITE SAMPLE SPACE

WE CAN "BIN" A DISTRIBUTION:



EXPECTED


$$N \cdot \int_{x_1}^{x_2} P(x) dx$$

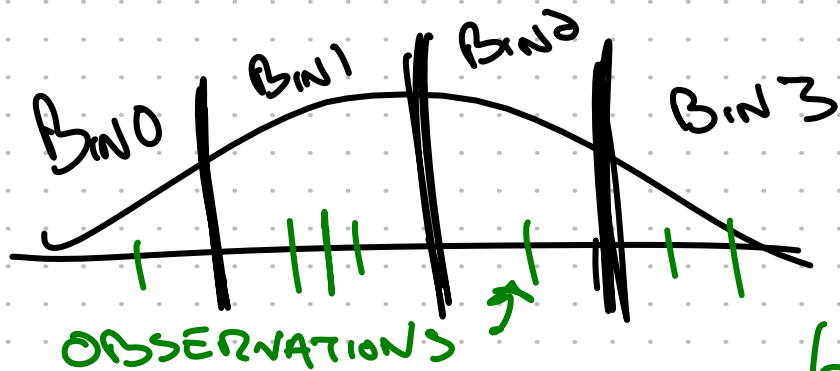
TOTAL OBSERVATIONS

PROB IN BIN ①

OBSERVED ARE SPLIT INTO BINS



CHOOSING BINS  IMPACTS ANALYSIS SENSITIVITY



7 TOTAL OBS

$$P(\text{Bin 1}) = \int_{-\infty}^{\infty} P(x) dx$$

→ CHOOSE BIN
EDGES SO BINS
HAVE EQUAL PROB

| | Bin 1 | Bin 2 | Bin 3 | Bin 4 |
|-------|-------|-------|-------|-------|
| O_i | 1 | 3 | 1 | 2 |
| E_i | | | | |
| | | | | |

$$Z = \frac{\bar{X} - \bar{Y}}{S}$$

Assuming $N_x = N_y$

$$Z \sim N(0, 1)$$

