

CS 2810 April 5 Day 20

chi square "goodness of fit"

Admin:

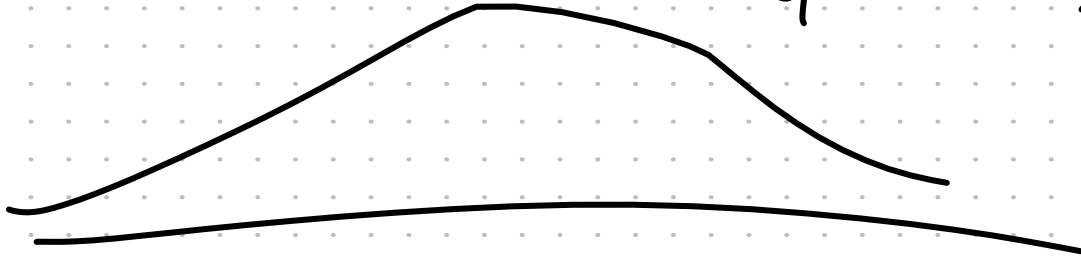
SWEET CLASS

SEC 3 FINALS SCHED

# 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).

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

$$X_7 = 11$$

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

DISTRIBUTED AS

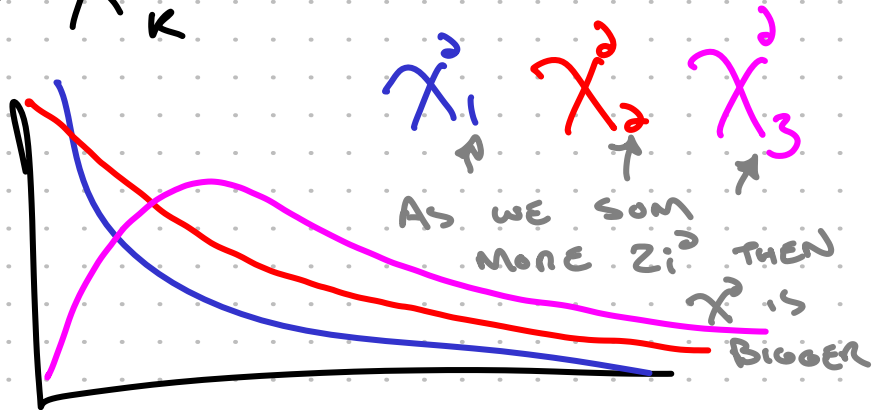
$$\frac{11 - 7}{\sqrt{10}} = \frac{4}{\sqrt{10}} \approx 1.2$$

# CHI-SQUARE DISTRIBUTION

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

CHI-SQUARE IS  
SUM OF  $K$   
STANDARD  
NORMAL  
SQUARED

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



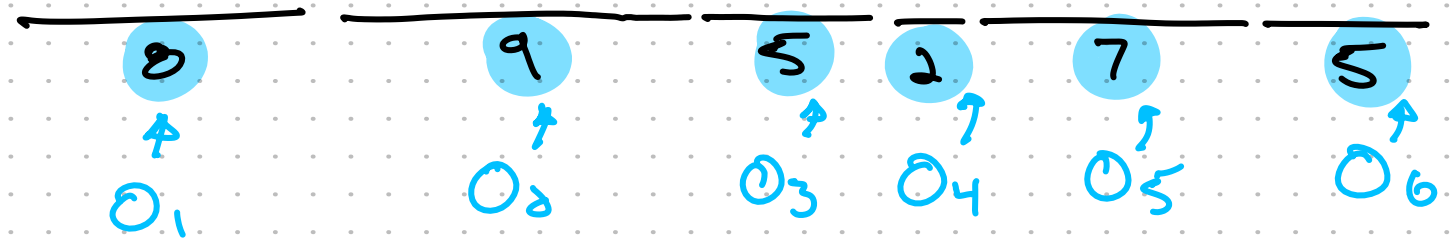
# CHI-SQUARE "GOODNESS OF FIT" TEST

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  $i$ ?

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

$$E_i = N \cdot p_i = 36 \cdot \frac{1}{6} = 6 \quad \text{for all outcomes}$$

TOTAL OUTCOMES  
(36 IN THIS EX)

PROB OF EACH  
OUTCOME UNDER  
 $H_0$

# CONTINGENCY TABLE + $\chi^2$ TEST STATISTIC

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

### 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.

chi-square = 0 happens only when each outcome is observed as many times as its expected

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

$\chi^2 = 5.3$



chi-square = 0 happens only when each outcome is observed as many times as its expected  
the chi-square stat closer to zero are more consistent with 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)

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

STATISTIC

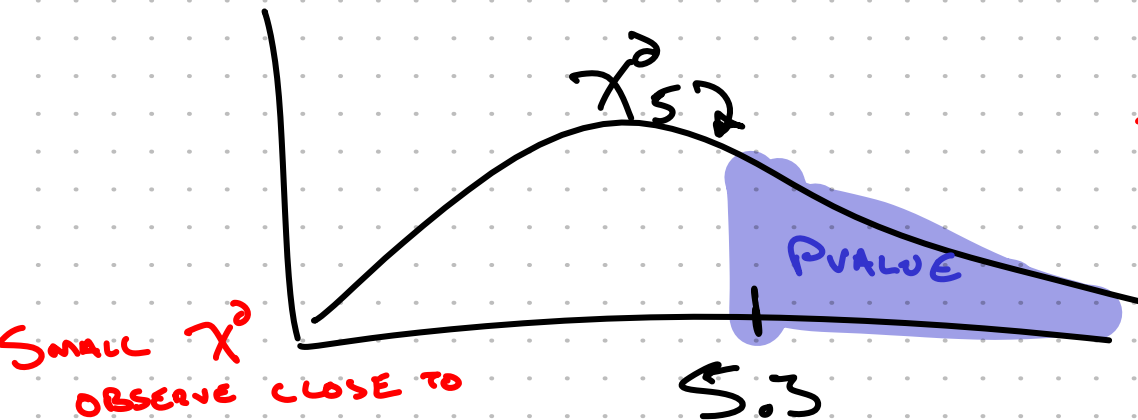
DISTRIBUTION

$k = \# \text{ CATEGORIES} - 1$   
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  $\chi^2$  DISTRIBUTED WITH  $DF=5$



SMALL  $\chi^2$   
OBSERVE CLOSE TO  
EXPECTED  
more consistent with  $H_0$

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

LARGE  $\chi^2$   
OBSERVE FAR  
FROM EXPECTED  
less consistent with  $H_0$

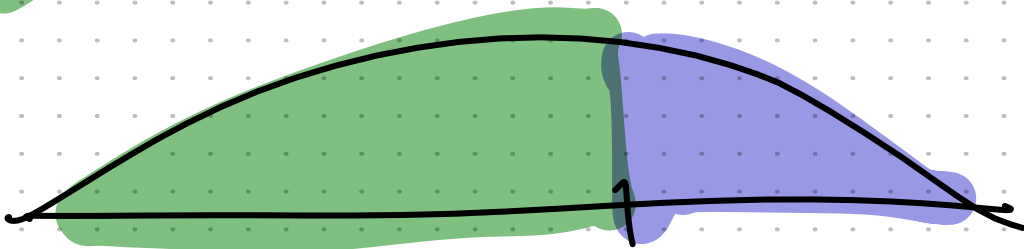
Finalize Chi2 example

NOT PVAL  $\rightarrow$

CDF(5.3,  $k=5$ )

PVALUE  $\downarrow$

$1 - \text{CDF}(5.3, k=5)$



PVAL = .38

5.3

$P_{VAL} > \alpha = .05$  ← **STRONG**

DON'T REJECT  $H_0$

... DIE COULD BE FAIR

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

REJECT  $H_0$ , CLAIM  $H_1$

... DIE 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, 6, 6

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 NOT SILLY

OUTCOME	1	2	3	4	5	6
PROB	$1/21$	$2/21$	$3/21$	$4/21$	$5/21$	$6/21$
$E_i$ EXPECTED # OUTCOMES	$\frac{1}{21} \cdot 33$	$\frac{2}{21} \cdot 33$	$\frac{3}{21} \cdot 33$	$\frac{4}{21} \cdot 33$	$\frac{5}{21} \cdot 33$	$\frac{6}{21} \cdot 33$
$O_i$ OBSERVED # OUTCOMES	3	3	4	7	4	12
$\chi^2 = \sum_i \frac{(E_i - O_i)^2}{E_i}$	1.29	.006	.108	.081	1.89	.701

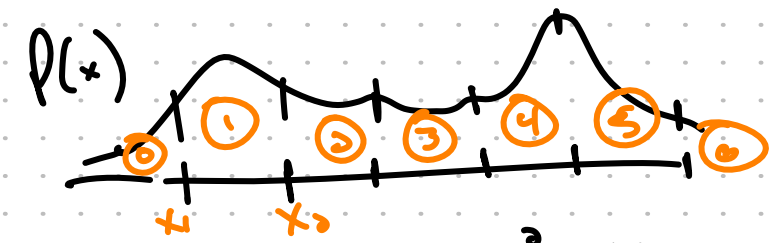
$3 + 3 + 4 + 7 + 4 + 12$   
 $= 33$  TOTAL TRIALS

$\chi^2 = 4.08$      $P = 1 - \text{CHI2.CDF}(4.08, DF=5) = .538$   
 DON'T REJECT  $H_0$

# CHI-SQUARE "BINNING"

$\chi^2$  REQUIRES A FINITE SAMPLE SPACE

WE CAN "BIN" A DISTRIBUTION:



EXPECTED

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

TOTAL OBSERVATIONS  $\uparrow$  PROB IN BIN  $\uparrow$  ①

OBSERVED ARE SPLIT INTO BINS



CHOOSING BINS  $\triangle$  IMPACTS ANALYSIS SENSITIVITY