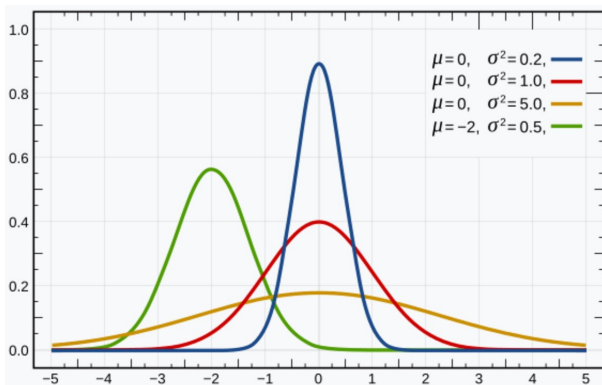


# $\chi^2$ Tests

Grinnell College

# Normal Distribution

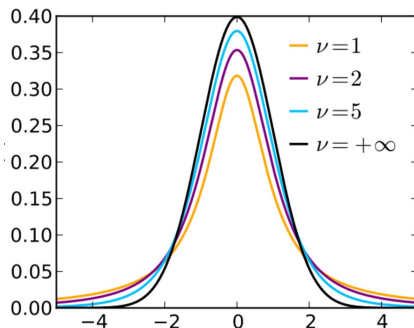


Distribution function:  $f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

This thing shows up *everywhere*

- ▶ population distributions
- ▶ sampling distributions for means and proportions
  - ▶ CLT

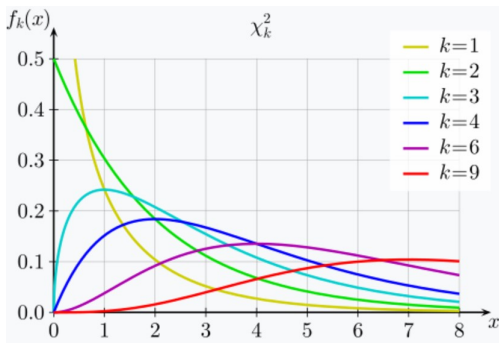
# t-Distribution



Distribution function:  $\frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi} \Gamma(\frac{\nu}{2})} \left(1 + \frac{x^2}{\nu}\right)^{-\left(\frac{\nu+1}{2}\right)}$ , parameter:  $\nu = \text{df}$

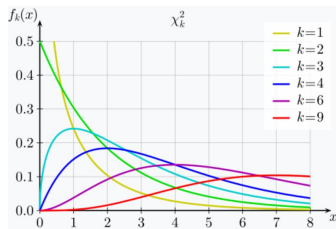
- ▶  $\mu = 0$ ,  $\sigma^2 = \frac{\nu}{\nu-2}$  for  $\nu > 2$ ,  $\infty$  for  $1 < \nu \leq 2$ , otherwise undefined
- ▶ shows up when we standardize a Normal variable but don't know  $\sigma$
- ▶ T test-statistic:  $T = \frac{\bar{X} - \mu}{s/\sqrt{n}} \sim t_{n-1}$

# $\chi^2$ -Distribution ("kai"-squared)



- ▶  $\mu = k, \sigma^2 = 2k$
- ▶ only parameter is  $k = \text{df} = \text{degrees of freedom}$

# $\chi^2$ -Distribution



Ok, but where does this distribution come from?

Suppose I have a whole bunch of Standard Normal variables  $X_1, X_2, \dots, X_k$  that are all independent (not influencing each other)

- ▶ If I square and add them, I get something that has a  $\chi^2$  distribution
- ▶  $df = k$ , where  $k$  is the number of Normals I used

$$\sum_{i=1}^k X_i^2 \sim \chi_k^2$$

# $\chi^2$ -Distribution

$$Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1) \rightarrow Z^2 \sim \chi_1^2$$

Could have done this:

$$Z_1 = \frac{\bar{x}_1 - \mu_1}{\sigma_1/\sqrt{n_1}} \sim N(0, 1)$$

$$Z_2 = \frac{\bar{x}_2 - \mu_2}{\sigma_2/\sqrt{n_2}} \sim N(0, 1)$$

$$Z_1^2 + Z_2^2 = \left(\frac{\bar{x}_1 - \mu_1}{\sigma_1/\sqrt{n_1}}\right)^2 + \left(\frac{\bar{x}_2 - \mu_2}{\sigma_2/\sqrt{n_2}}\right)^2 \sim \chi_2^2$$

- ▶ could have used  $\chi^2$  distribution for means and diff. in means
- ▶ but... Normal distribution is easier
- ▶ this **does** come in more handy for differences in proportions

# $\chi^2$ Hypothesis Tests

The most common forms of hypotheses we will be testing with the  $\chi^2$ -test:

$$H_0 : p_1 = p_2 \text{ OR } H_0 : p_1 - p_2 = 0$$

$$H_A : p_1 \neq p_2 \text{ OR } H_A : p_1 - p_2 \neq 0$$

If we have  $N > 2$  groups

$$H_0 : p_1 = \dots p_N \text{ OR } H_0 : \text{all proportions are the same}$$

$H_A$  : at least one proportion is different than the others

Pro's and Con's:

- ▶  $\chi^2$ -tests work for any number of groups
- ▶ only works for ' $\neq$ ' hypotheses

# $\chi^2$ Hypothesis Tests

A bit more general – we can test if all the proportions are equal to different specified values, but writing the general hypothesis forms get complicated.

- ▶ there is no agreed on notation for this I am aware of, usually specified qualitatively with context

$$H_0 : p_1 = p_{1(o)}, p_2 = p_{2(o)}, \dots, p_N = p_{N(o)}$$

**OR**  $H_0$  : all proportions equal the specified value

$H_A$  : at least one proportion is different than the specified value

## Problem

The 'final' exam is coming up. Suppose (hypothetically), that I am lazy and want to write an exam that is particularly easy to grade, so I make 1 question with 5 multiple-choice answers. But, also suppose (hypothetically) that I am mean and make the questions incredibly difficult.

Given that there are 25 students in this class, I can look through the results and make the following table to compare how many people got a question right (Observed) and compare that to how many I would expect to get the question right by guessing (Expected).

Answer	A	B	C	D	E
Expected	5	5	5	5	5
Observed	6	7	6	2	5

## Goodness of Fit

The  $\chi^2$  (chi squared or “kai” squared) **goodness of fit** test allows us to compare *expected* proportions in  $k$  groups against those we *observe*

$$\chi^2 = \sum_{i=1}^k \frac{(\text{Expected}_i - \text{Observed}_i)^2}{\text{Expected}_i}$$

Under the null hypothesis (equality of proportions to ‘expected’ proportions), for  $k$  groups, the  $\chi^2$  goodness of fit test statistic follows a  $\chi^2$  distribution with  $k - 1$  degrees of freedom

$$\chi^2 \sim \chi^2(k - 1)$$

- $df = k - 1$  comes from the fact that if we know  $k - 1$  proportions, we know the value for the last group since they all add to 1

Answer	A	B	C	D	E
Expected	5	5	5	5	5
Observed	6	7	6	2	4

$$\begin{aligned}
 \chi^2 &= \sum_{i=1}^k \frac{(\text{Expected}_i - \text{Observed}_i)^2}{\text{Expected}_i} \\
 &= \frac{(5 - 6)^2}{5} + \frac{(5 - 7)^2}{5} + \frac{(5 - 6)^2}{5} + \frac{(5 - 2)^2}{5} + \frac{(5 - 4)^2}{5} \\
 &= \frac{1 + 4 + 1 + 9 + 1}{5} = \frac{16}{5} = 3.2
 \end{aligned}$$

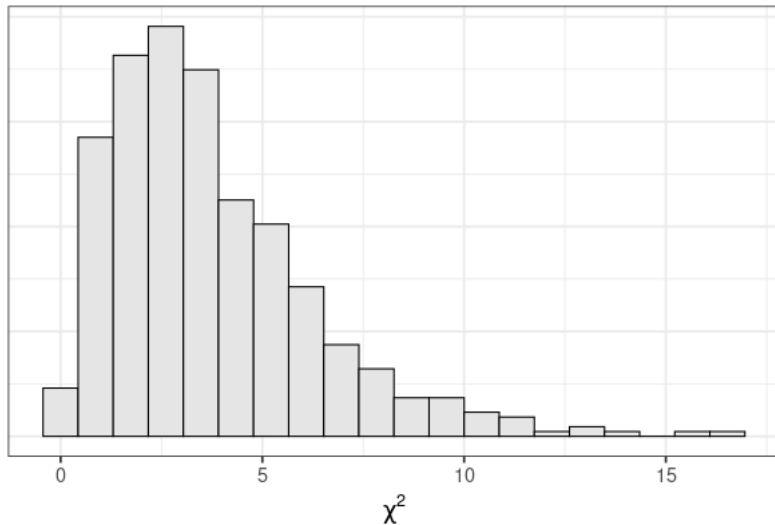
# Samples

	A	B	C	D	E
Sample 1	6	8	5	3	3
Sample 2	3	2	2	10	8
Sample 3	2	6	3	3	11
Sample 4	6	5	3	4	7
Sample 5	5	7	5	4	4
Sample 6	4	4	6	4	7
Sample 7	3	5	7	3	7
Sample 8	5	7	4	3	6
Sample 9	5	5	5	6	4
Sample 10	5	2	10	3	5

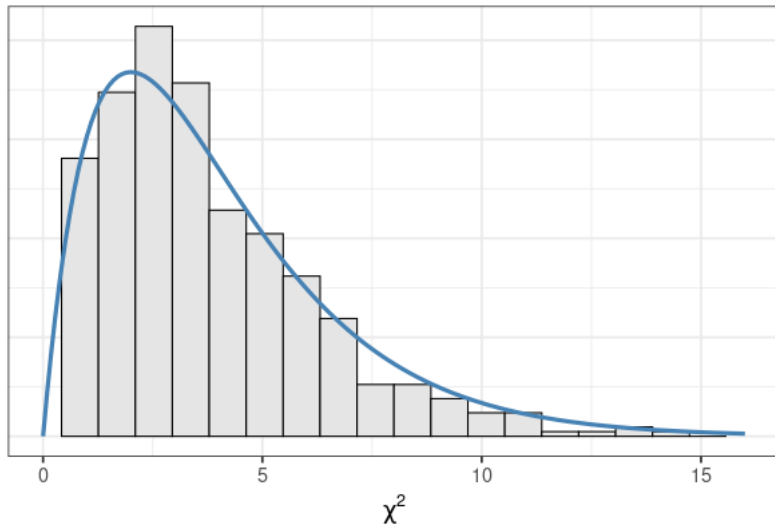
# Samples

	A	B	C	D	E	chi
Sample 1	6	8	5	3	3	4
Sample 2	3	2	2	10	8	11
Sample 3	2	6	3	3	11	11
Sample 4	6	5	3	4	7	2
Sample 5	5	7	5	4	4	1
Sample 6	4	4	6	4	7	2
Sample 7	3	5	7	3	7	3
Sample 8	5	7	4	3	6	2
Sample 9	5	5	5	6	4	0
Sample 10	5	2	10	3	5	8

## Histogram of $\chi^2$ Statistics for $df = 4$

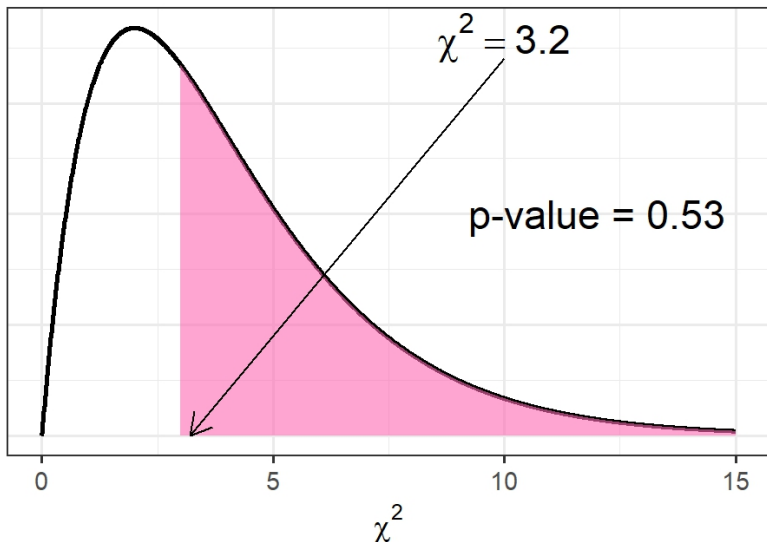


## Histogram of $\chi^2$ Statistics for $df = 4$

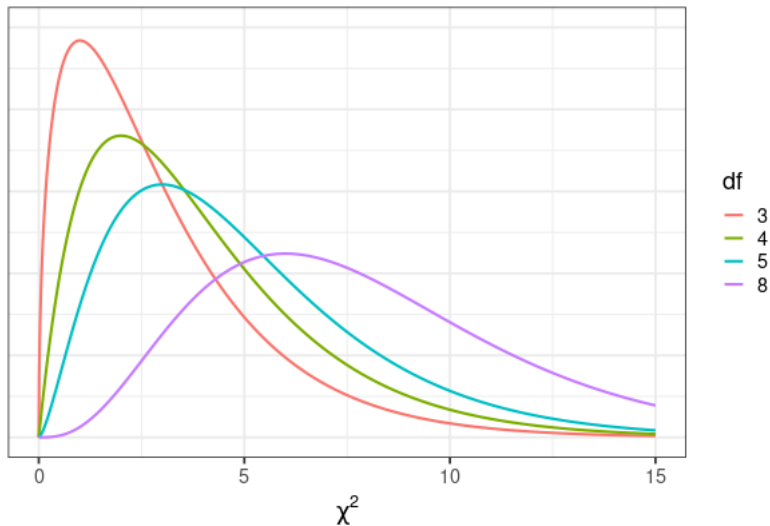


## *p*-value for exam questions

Chi-squared distribution with  $df = 4$



## Histogram of $\chi^2$ Statistics



## $p$ -value for $\chi^2$

$$\chi^2 = \sum_{i=1}^k \frac{(\text{Expected}_i - \text{Observed}_i)^2}{\text{Expected}_i}$$

A few things to note about this statistic:

- ▶ It's always positive (or equal to zero)
- ▶ The more our observed values deviate from our expected, the larger it gets

From this, we get two facts:

- ▶ Our  $p$ -value is computed as the area *to the right* of our test statistic
- ▶ Greater values of  $\chi^2$  indicate more evidence against the null hypothesis

# Jury Example

Prospective jurors are supposed to be randomly chosen from the eligible adults in a community. The American Civil Liberties Union (ACLU) studied the racial composition of the jury pools in 10 trials in Alameda County, California. Display below is the racial and ethnicity composition of the  $n = 1,453$  individuals included in the jury pools, along with the distribution of eligible jurors according to US Census data:

Race Ethnicity	White	Black	Hispanic	Asian	Other	Total
Jury Size	780	117	114	384	58	1453
Census Percentage	54%	18%	12%	15%	1%	100%

## Jury Example – R code

Race Ethnicity	White	Black	Hispanic	Asian	Other	Total
Jury Size	780	117	114	384	58	1453
Census Percentage	54%	18%	12%	15%	1%	100%

```
> obs <- c(780, 117, 114, 384, 58)
> percent <- c(.54, .18, .12, .15, .01)
> exp <- 1453*percent
> exp
[1] 784.62 261.54 174.36 217.95 14.53
> ((obs - exp)^2 / exp) %>% sum()
[1] 357.3625
>
> # P-value
> pchisq(357.3625, df=4, lower.tail = F)
[1] 4.510471e-76
```