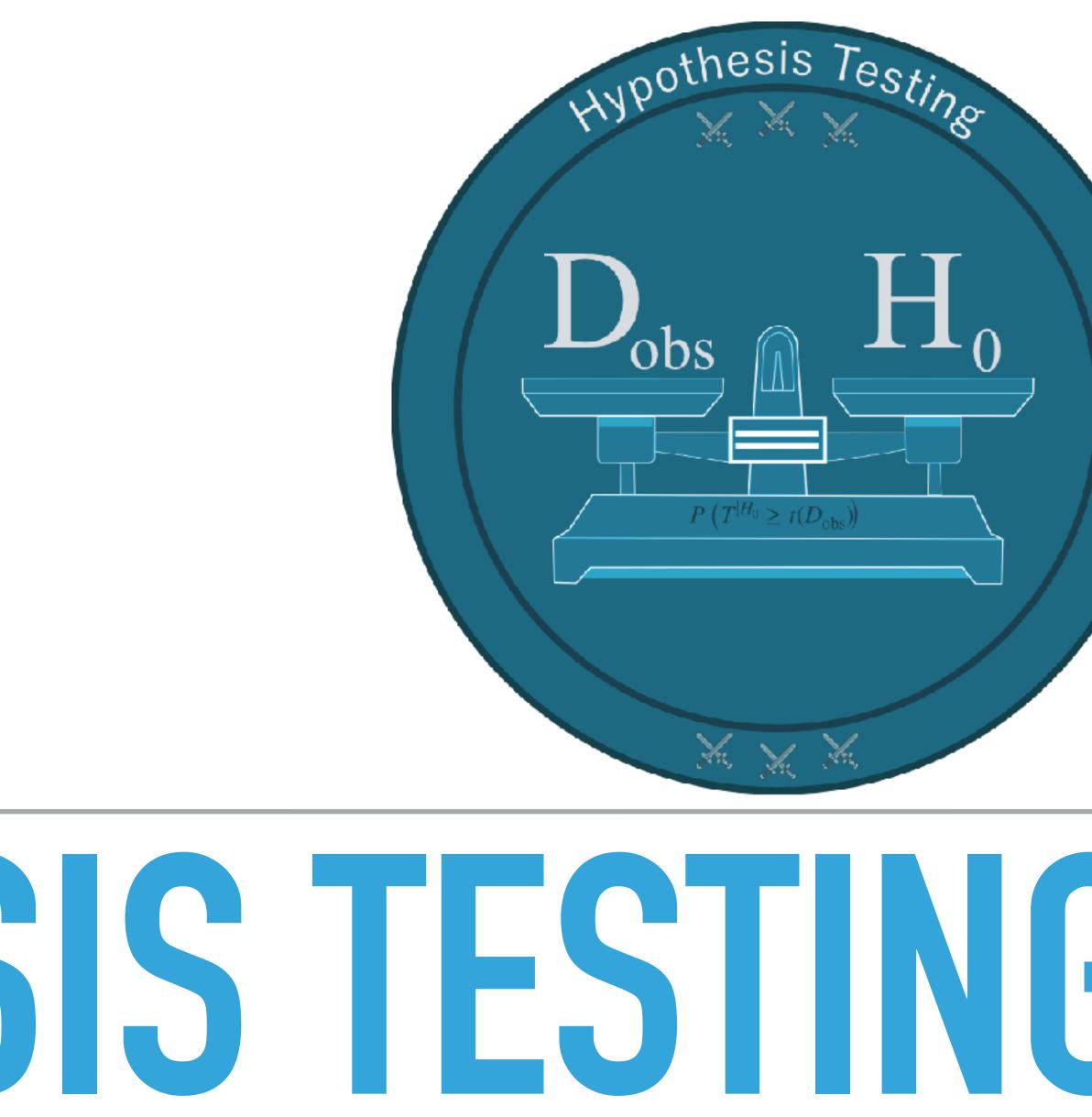


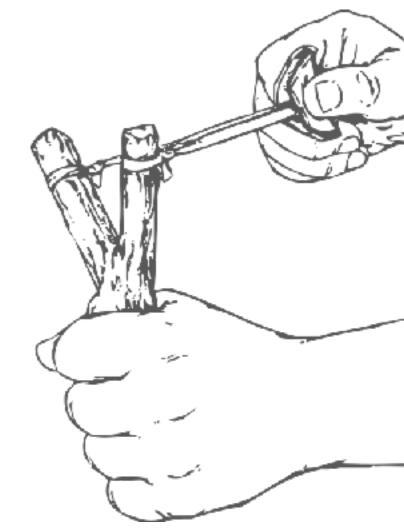
INTRODUCTION TO DATA ANALYSIS

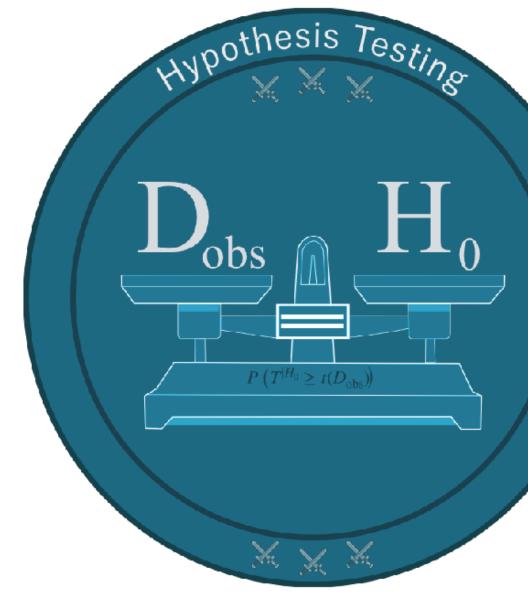




LEARNING GOALS

- become able to interpret & apply some statistical tests
 - Pearson's χ^2 -tests of independence
 - > z-test
 - one-sample *t*-test
 - two-sample t-test
 - one-way ANOVA
- understand differences and commonalities of different approaches to frequentist testing
 - **Fisher**
 - Neyman/Pearson
 - modern hybrid NHST



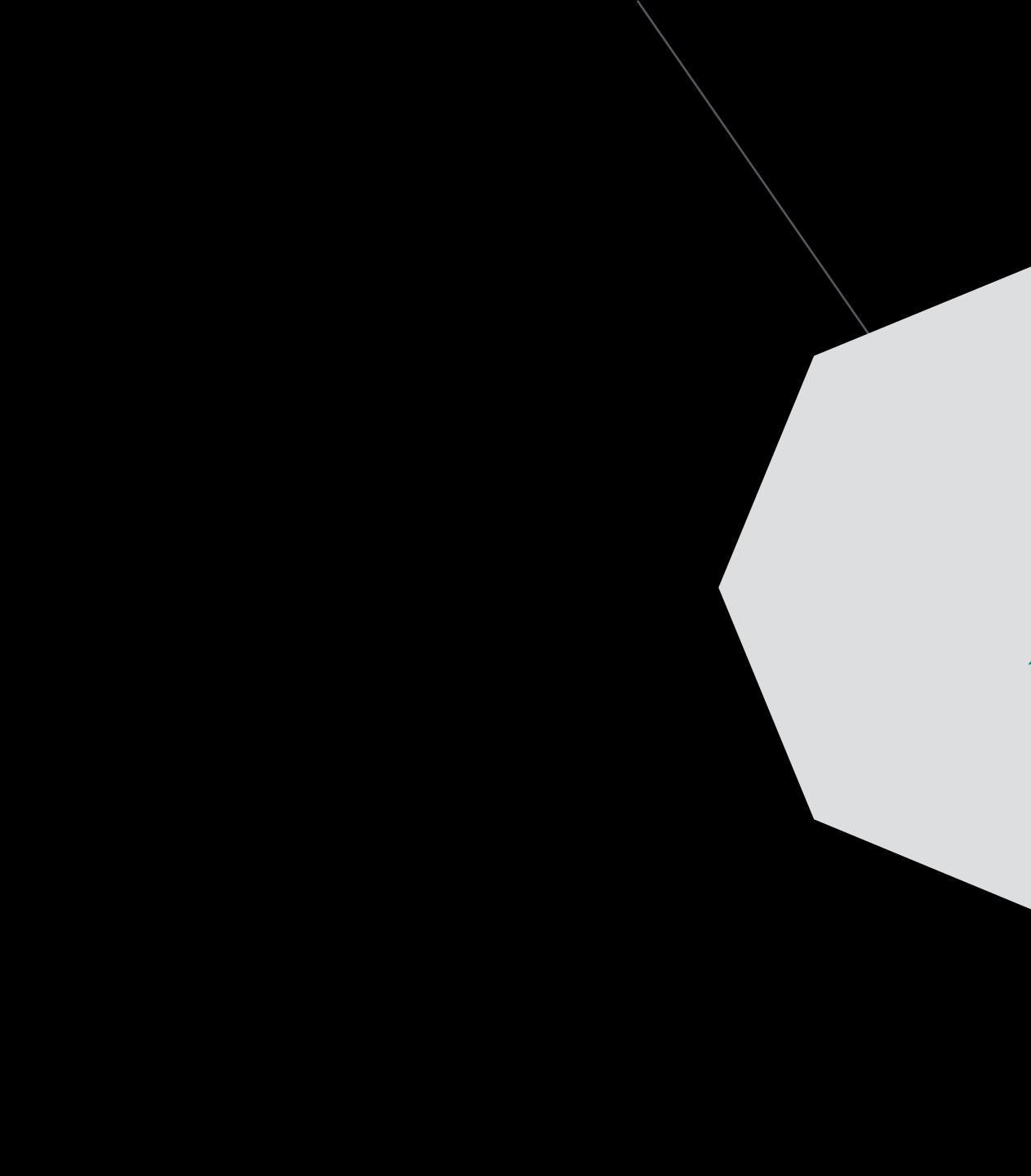










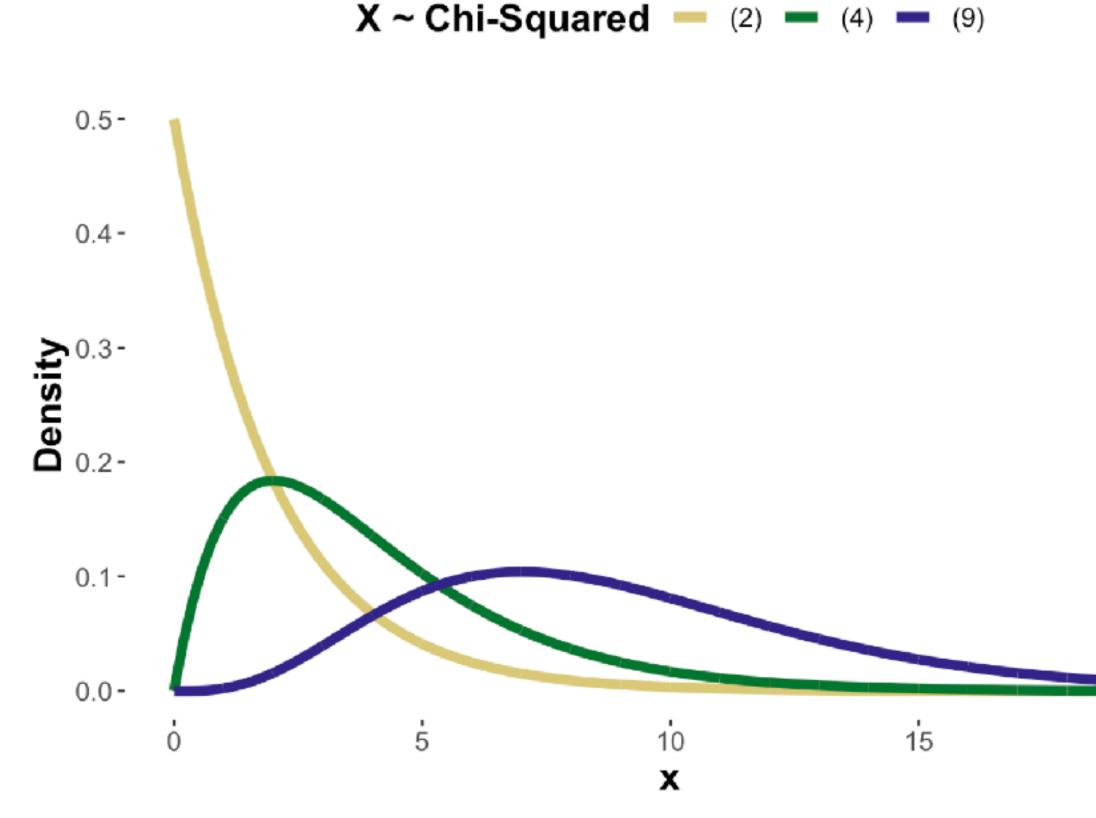


Pearson's 2-test goodness of fit

PEARSON χ^2 -**TESTS**

- tests for categorical data (with more than two categories)
- two flavors:
 - test of goodness of fit
 - test of independence
- sampling distribution is a χ^2 -distribution

standard normal random variables: X₁,...X_n derived RV: Y = X₁² + ... + X_n² it follows (by construction) that: y ~ χ²-distribution(n)



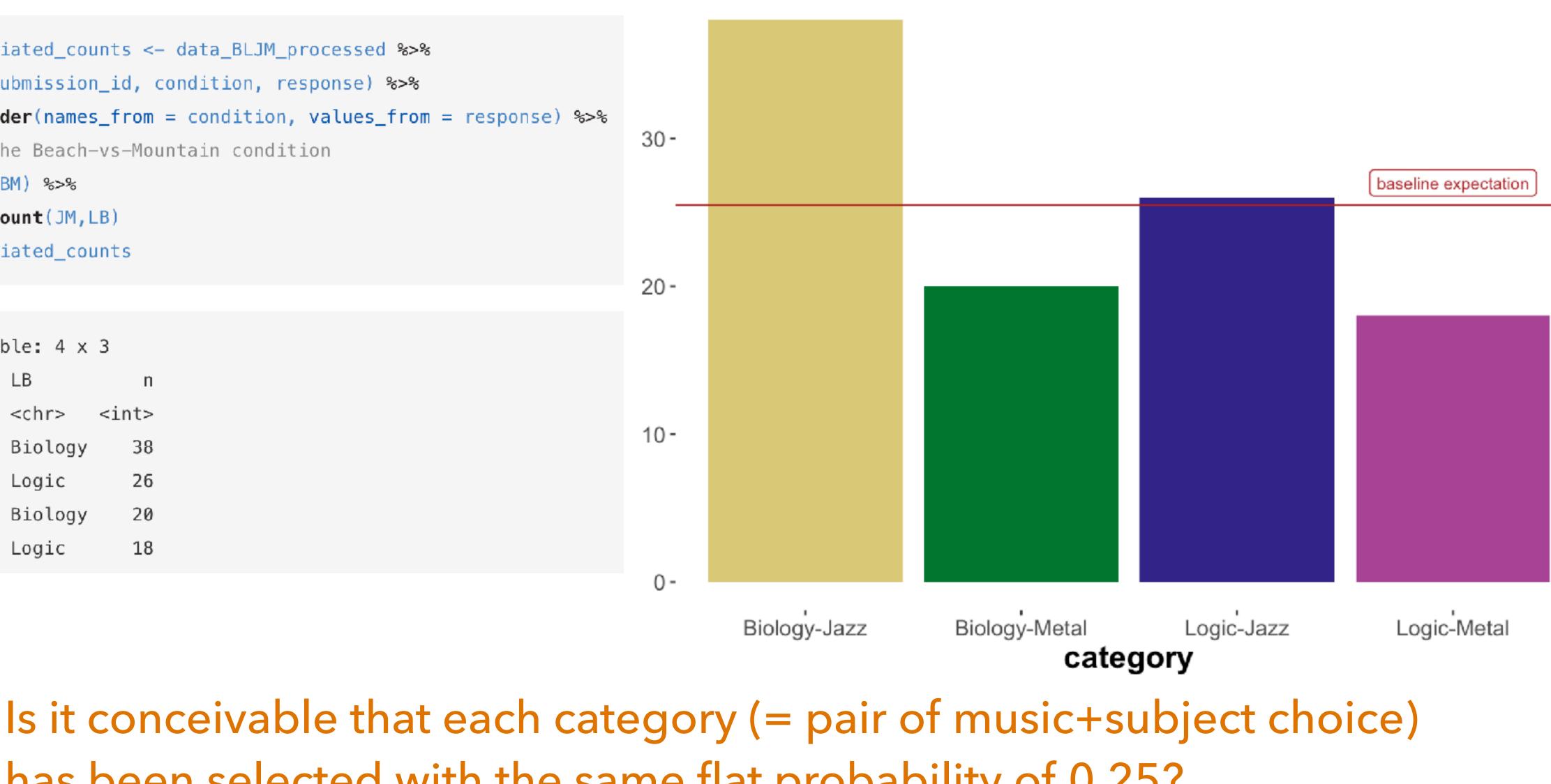


```
BLJM_associated_counts <- data_BLJM_processed %>%
 select(submission_id, condition, response) %>%
 pivot_wider(names_from = condition, values_from = response) %>%
 # drop the Beach-vs-Mountain condition
 select(-BM) %>%
 dplyr::count(JM,LB)
BLJM_associated_counts
```

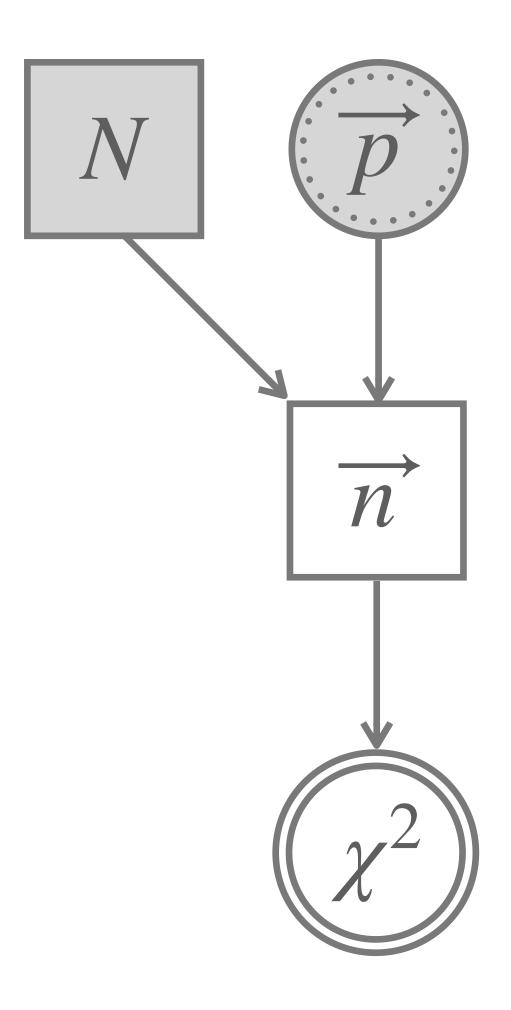
##	#	A tib	ole: 4 x	3
##		JM	LB	n
##		<chr></chr>	<chr></chr>	<int></int>
##	1	Jazz	Biology	38
##	2	Jazz	Logic	26
##	3	Metal	Biology	20
##	4	Metal	Logic	18

has been selected with the same flat probability of 0.25?

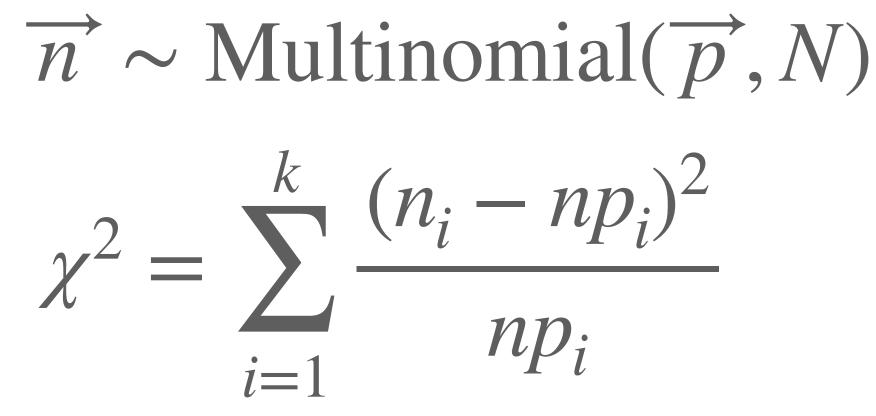


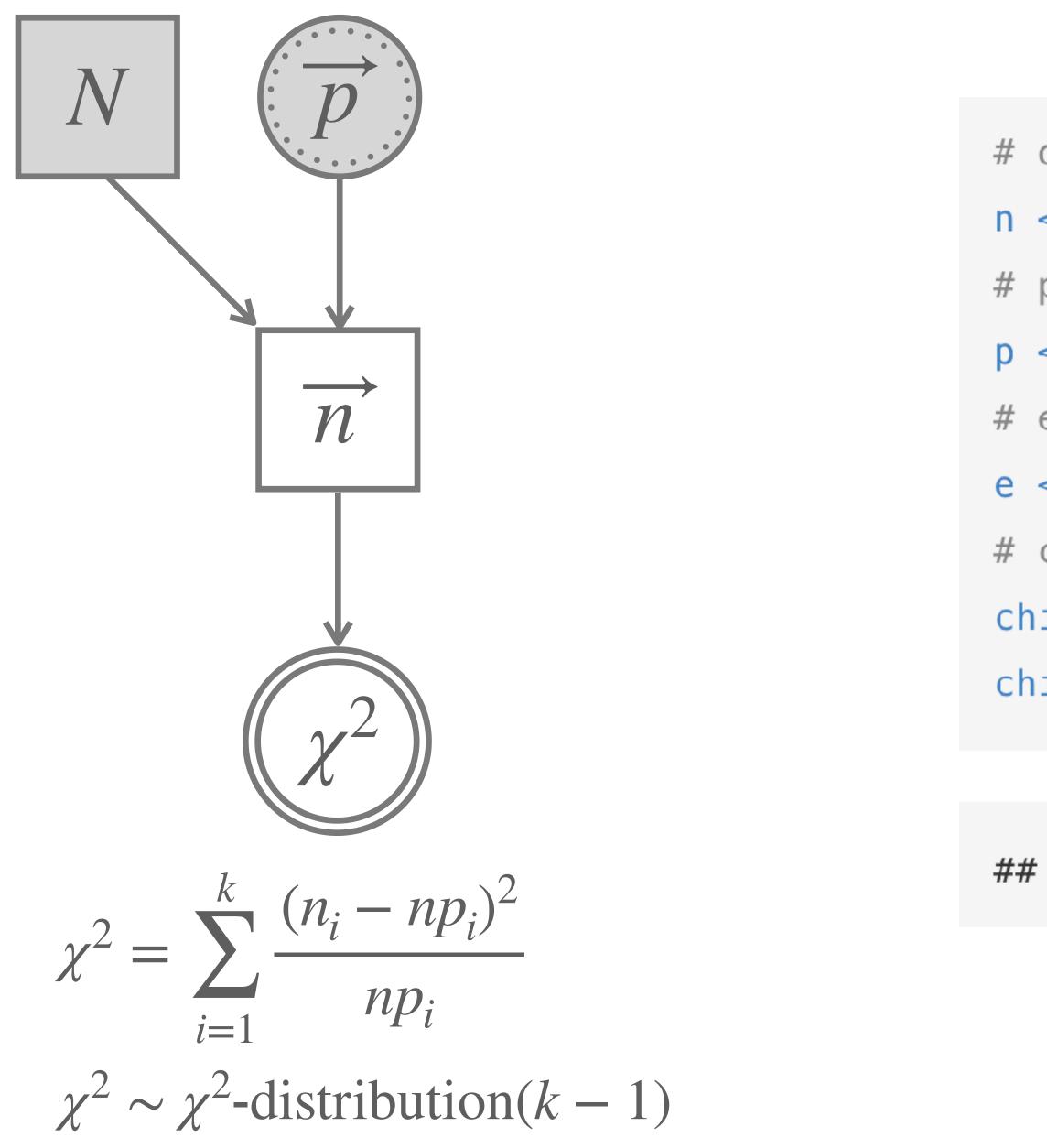


FREQUENTIST MODEL FOR PEARSON'S χ^2 **-TEST** [GOODNESS OF FIT]



Sampling distribution: $\chi^2 \sim \chi^2$ -distribution(k-1)



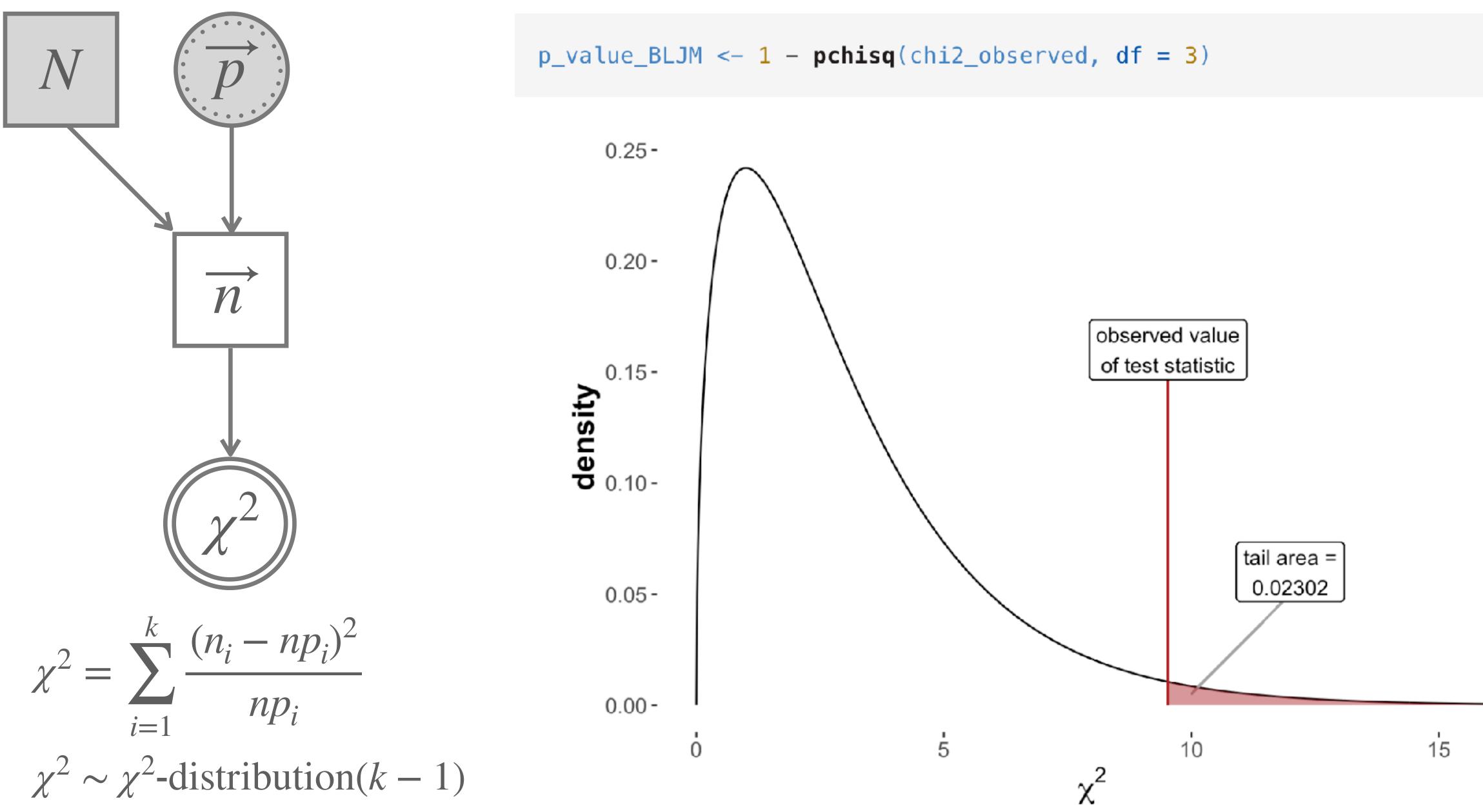




```
# observed counts
n <- counts_BLJM_choice_pairs_vector
# proprortion predicted
p <- rep(1/4,4)
# expected number in each cell
e <- sum(n)*p
# chi-squared for observed data
chi2_observed <- sum((n-e)^2 *1/e)
chi2_observed
```

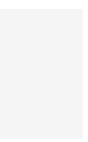
[1] 9.529412

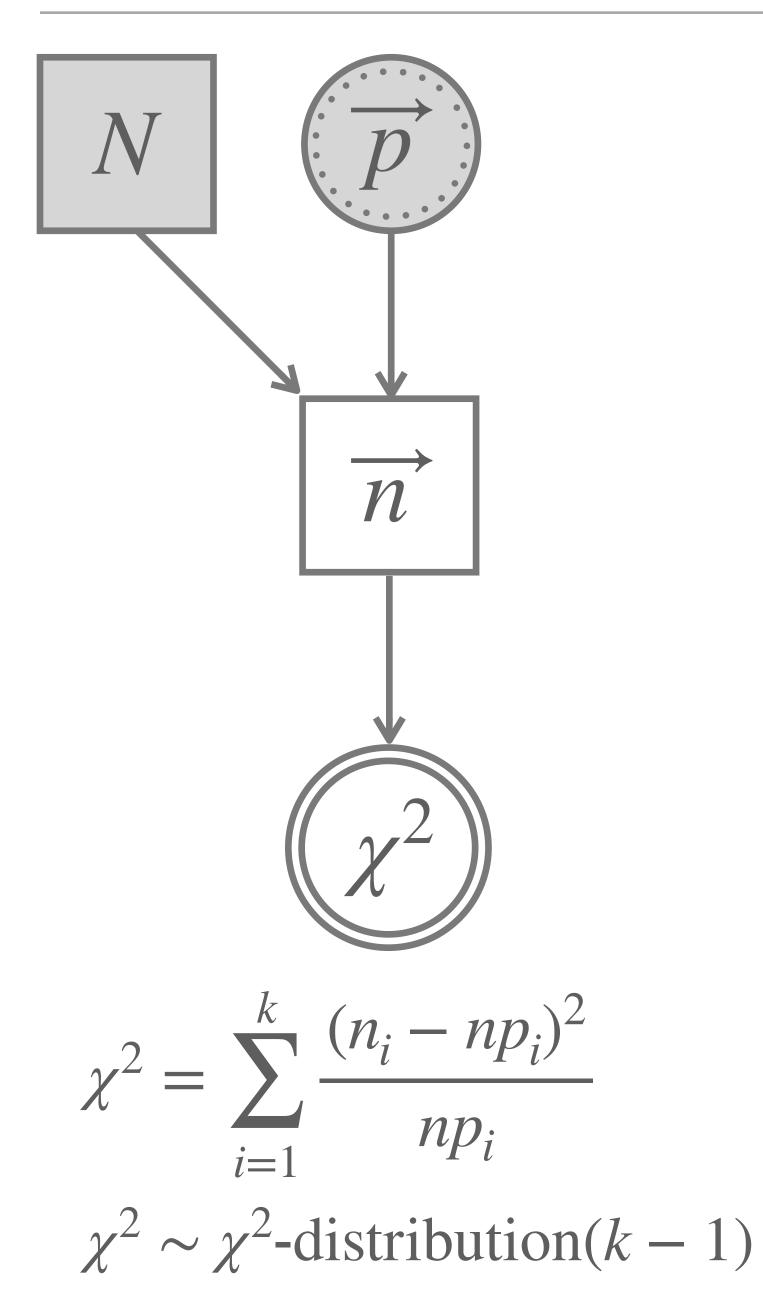












##

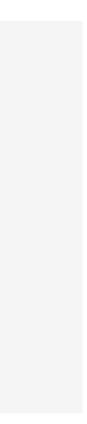


counts_BLJM_choice_pairs_vector <- BLJM_associated_counts %>% pull(n) chisq.test(counts_BLJM_choice_pairs_vector)

Chi-squared test for given probabilities

```
## data: counts_BLJM_choice_pairs_vector
## X-squared = 9.5294, df = 3, p-value = 0.02302
```





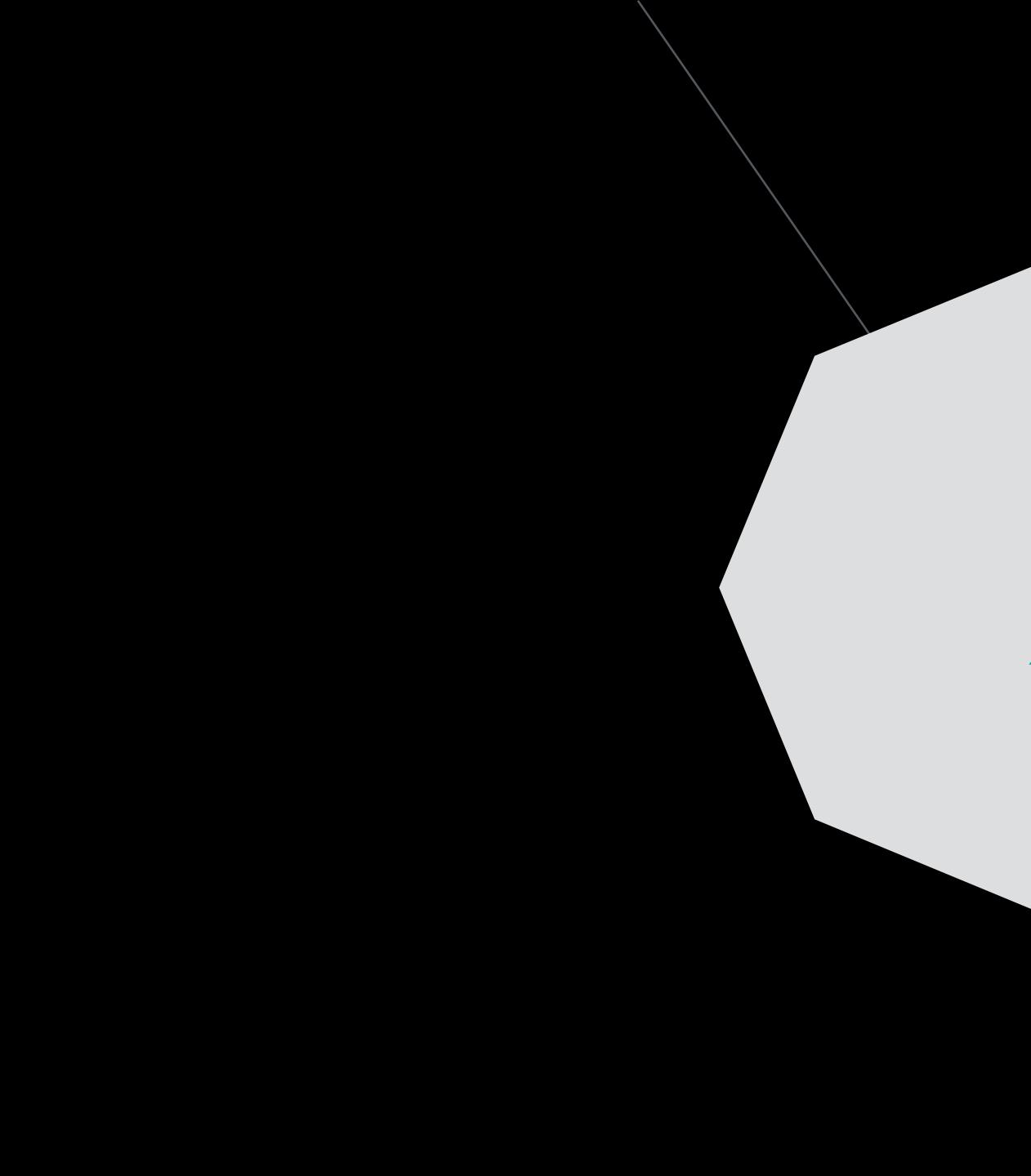
How to interpret / report the result:

Observed counts deviated significantly from what is expected if each category (here: pair of music+subject choice) was equally likely (χ^2 -test, with $\chi^2pprox 9.53$, df=3 and ppprox 0.023).



What about the lecturer's conjecture that (colorfully speaking) logic + metal = 😂?





Pearson's 2-test independence

STOCHASTIC INDEPENDENCE

events A and B are stochastically independent iff intuitively: learning one does not change beliefs about the other; • formally: $P(A \mid B) = P(A)$

• notice that $P(A \mid B) = P(A)$ entails that $P(B \mid A) = P(B)$ (see web-book)

STOCHASTIC INDEPENDENCE

Proposition 7.1 (Probability of conjunction of stochastically independent events) For any pair of events A and B with non-zero probability:

$$P(A \cap B) = P(A) \ P(B)$$

$$egin{aligned} P(A \cap B) &= P(A \mid B) \ P(B) \ &= P(A) \ P(B) \end{aligned}$$

Table 7.2: Joint probability table for a fli of 0.8 towards heads and where each

	heads	tails	Σ rows
black	0.8 imes 0.3=0.24	0.2 imes 0.3=0.06	0.3
white	0.8 imes 0.7=0.56	0.2 imes 0.7=0.14	0.7
Σ columns	0.8	0.2	1.0

- [if A and B are stoch. independent]
- *Proof.* By assumption of independence, it holds that $P(A \mid B) = P(A)$. But then:
 - [def. of conditional probability] [by ass. of independence]

lip-and-draw scenario where the coin has a bias	
of the two urns hold 3 black and 7 white balls.	

PEARSON'S χ^2 -**TEST** [INDEPENDENCE]

BLJM_table <- BLJM_associated_counts %>% select(-category) %>% pivot_wider(names_from = LB, values_from = n) BLJM_table

##	#	A tib	ole: 2 x	3
##		JM	Biology	L
##		<chr></chr>	<int></int>	<
##	1	Jazz	38	
##	2	Metal	20	

Is it conceivable that the outcome in each cell is given by independent choices of row and column options?

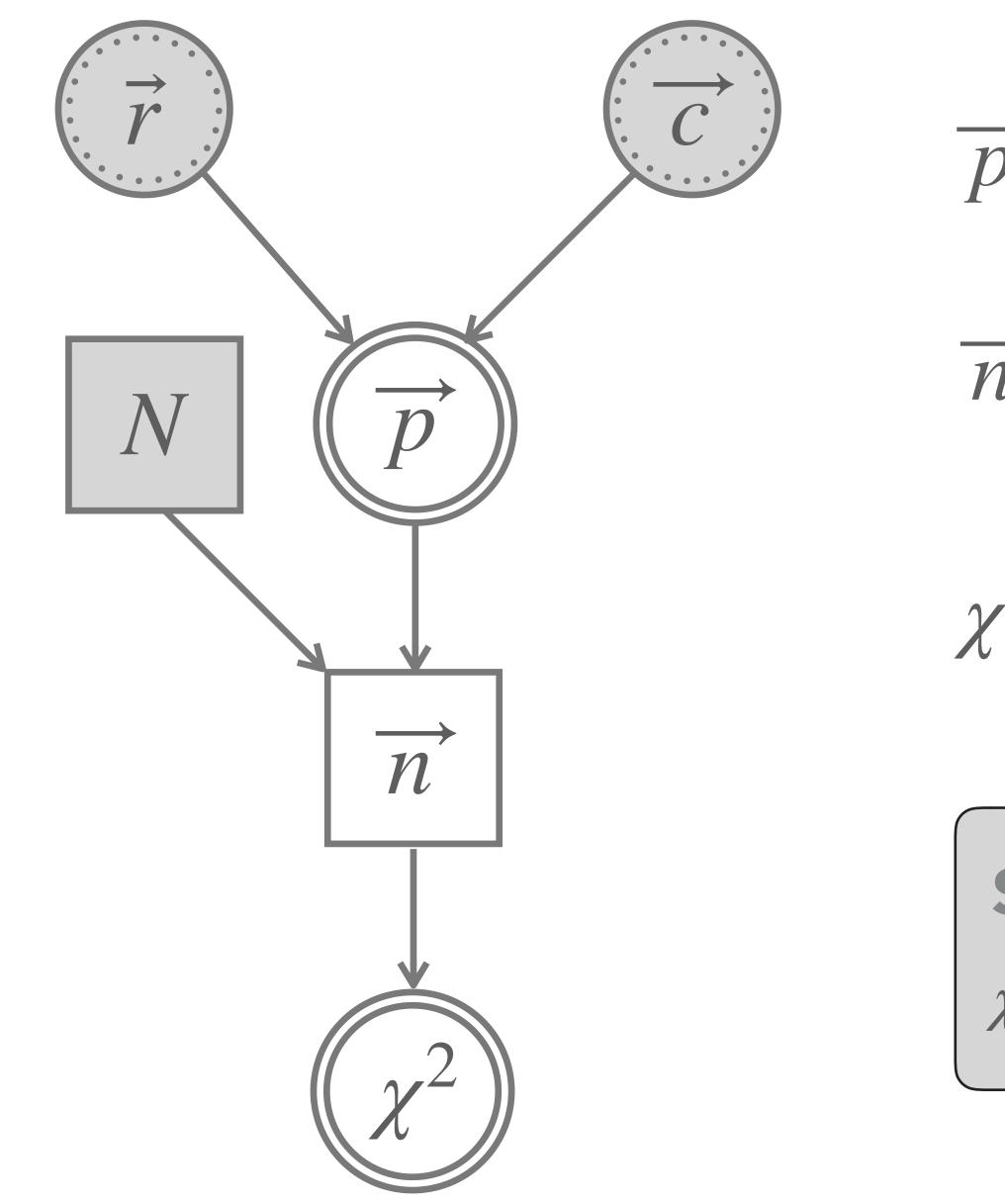
Hence: is the probability of a choice of cell the product of the probability of row- and column choices?



.ogic int>

26

18



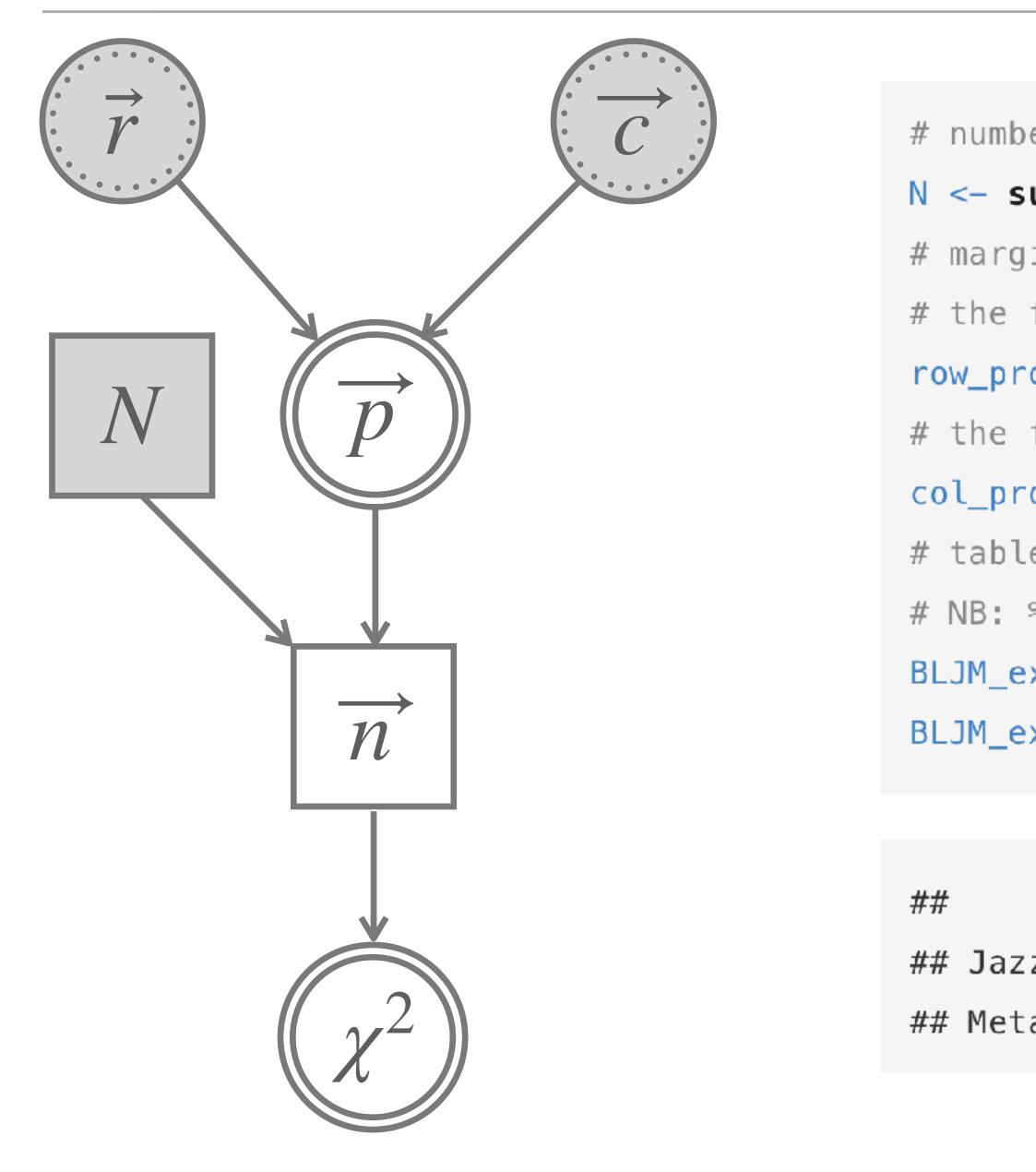
 \overrightarrow{p} = vec. of outer product $\overrightarrow{r} \& \overrightarrow{c}$

 $\overrightarrow{n} \sim \text{Multinomial}(\overrightarrow{p}, N)$

$$L^{2} = \sum_{i=1}^{k} \frac{(n_{i} - np_{i})^{2}}{np_{i}}$$

Sampling distribution:

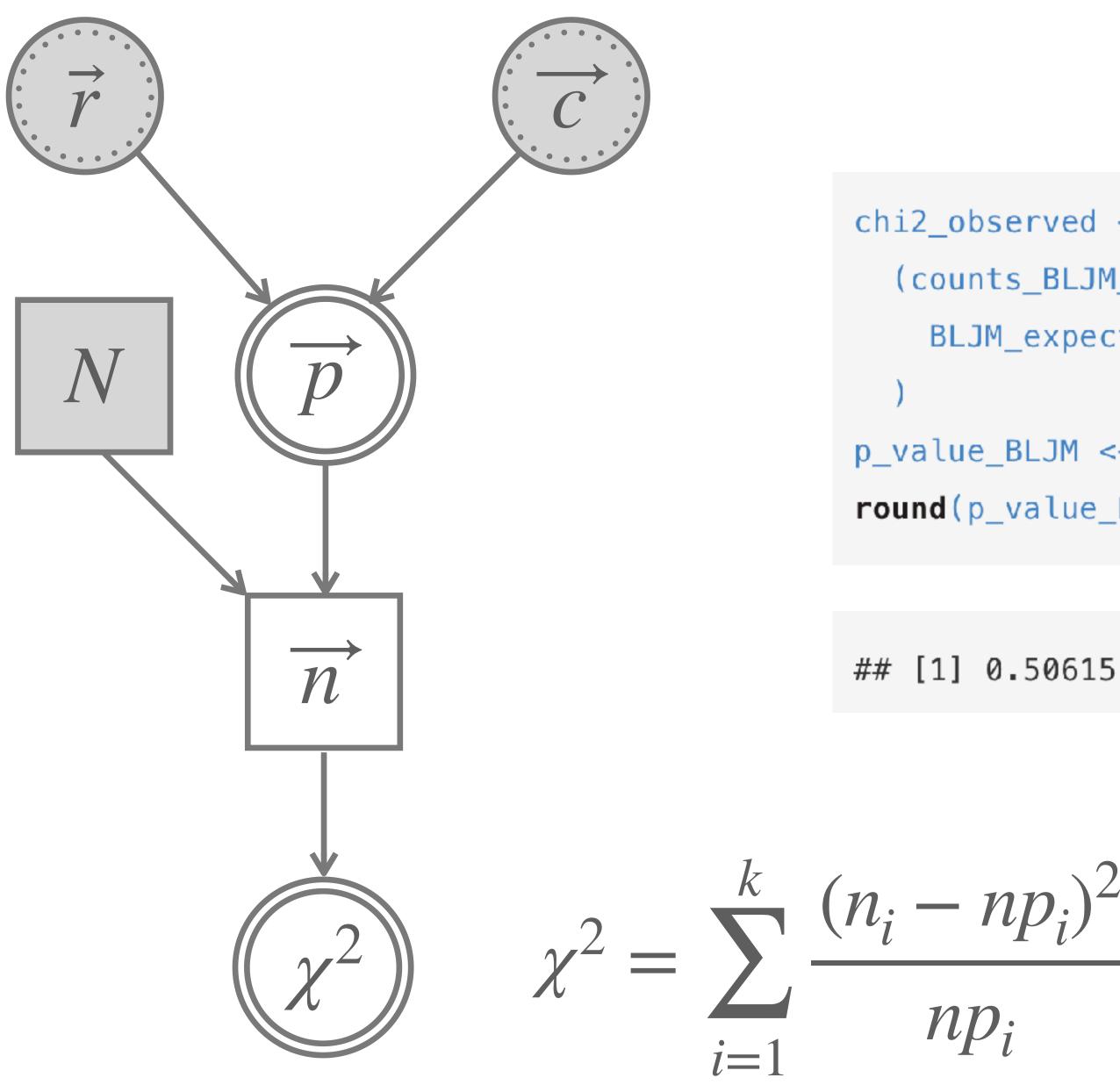
$$\chi^2 \sim \chi^2$$
-distribution $((k_r - 1) \cdot (k_c - 1))$



number of observations in total N <- sum(counts_BLJM_choice_pairs_matrix) # marginal proportions observed in the data # the following is the vector r in the model graph row_prob <- counts_BLJM_choice_pairs_matrix %>% rowSums() / N # the following is the vector c in the model graph col_prob <- counts_BLJM_choice_pairs_matrix %>% colSums() / N # table of expected observation under independence assumption # NB: %o% is the outer product of vectors BLJM_expectation_matrix <- (row_prob %o% col_prob) * N BLJM_expectation_matrix

Biology Logic
Jazz 36.39216 27.60784
Metal 21.60784 16.39216





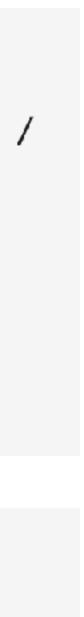


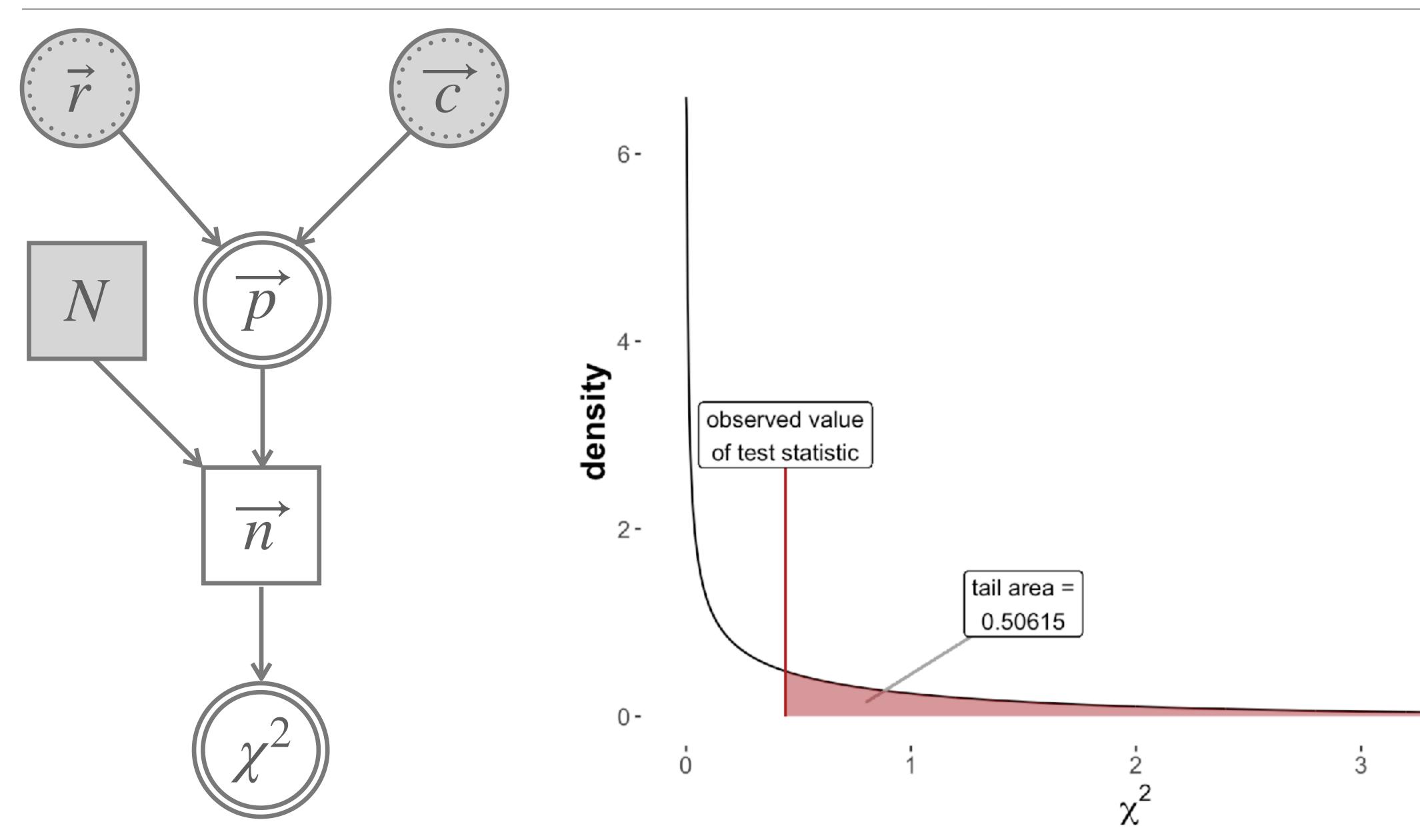
chi2_observed <- sum(</pre>

(counts_BLJM_choice_pairs_matrix - BLJM_expectation_matrix)^2 / BLJM_expectation_matrix

```
p_value_BLJM <- 1-pchisq(q = chi2_observed, df = 1)</pre>
round(p_value_BLJM,5)
```

$$(np_i)^2$$

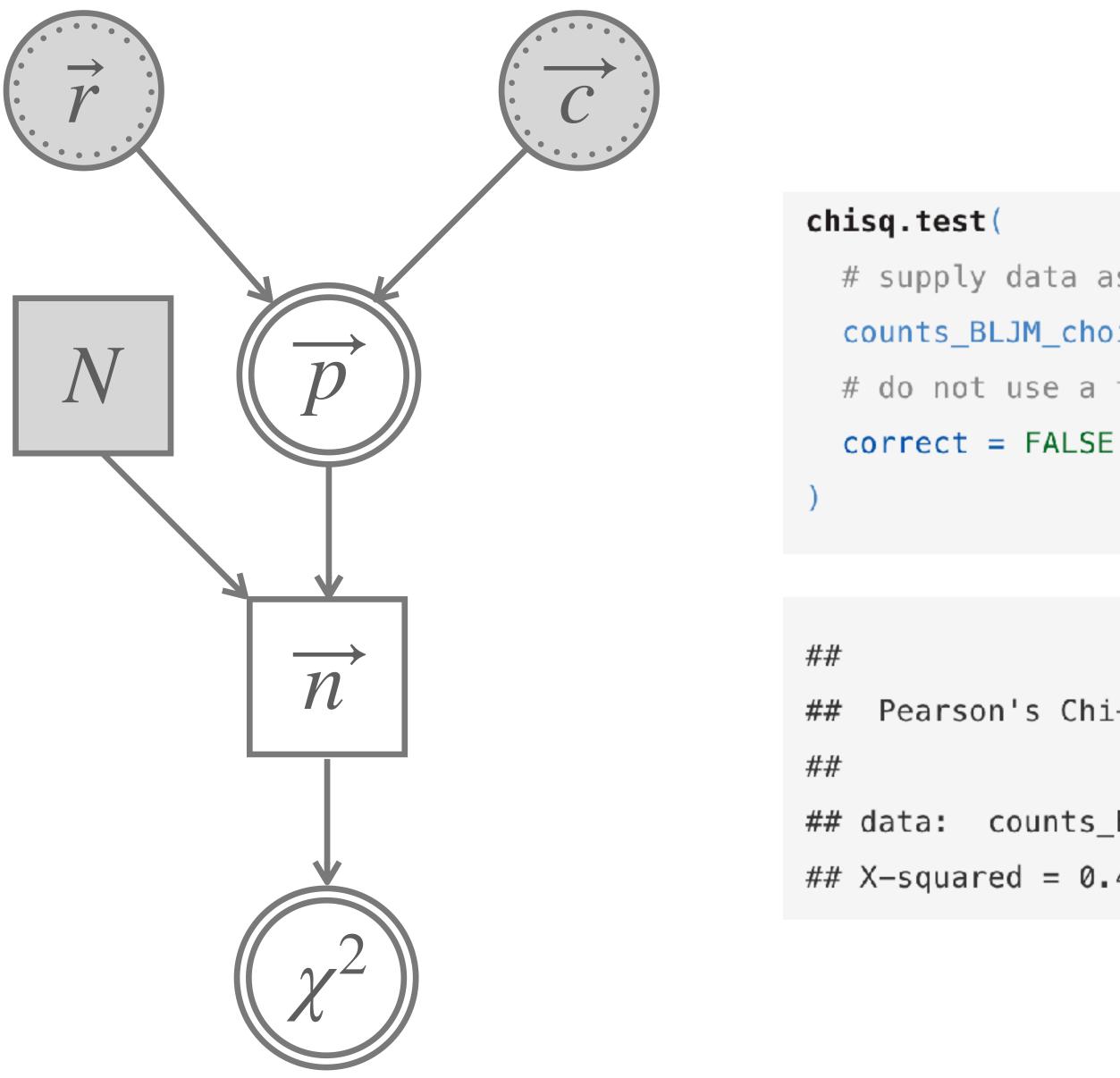














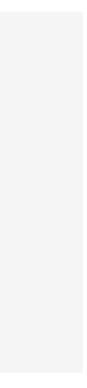
supply data as a matrix, not as a vector, for test of independence counts_BLJM_choice_pairs_matrix,

do not use a the default correction (because we didn't introduce it)

Pearson's Chi-squared test

data: counts_BLJM_choice_pairs_matrix ## X-squared = 0.44202, df = 1, p-value = 0.5061





How to interpret / report the result:

A χ^2 -test of independence did not yield a significant test result (χ^2 -test, with $\chi^2 pprox 0.44$, df = 1 and p pprox 0.5). Therefore, we cannot claim to have found any evidence for the research hypothesis of dependence.







SCENARIO FOR A *Z***-TEST** [ONE-SAMPLE]

- metric variable \vec{x} with samples from normal distribution
 - unknown μ
 - known σ [usually unrealistic!]

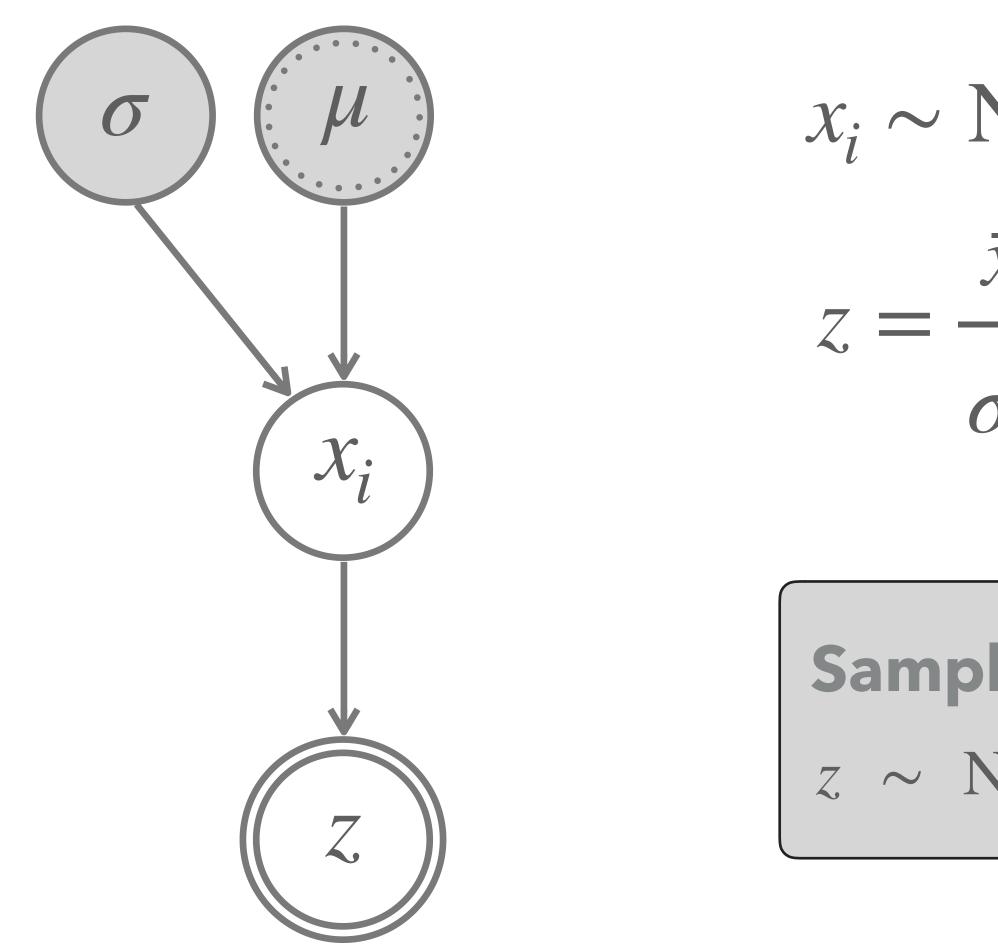
```
IQ_data <- c(87, 91, 93, 97, 100, 101, 103, 104,
             104, 105, 105, 106, 108, 110, 111,
             112, 114, 115, 119, 121)
mean(IQ_data)
```

[1] 105.3

Is it plausible to maintain that this data was generated by a normal distribution with mean 100 (if we assume that the standard deviation is known to be 15)?



FREQUENTIST MODEL FOR A *Z***-TEST** [ONE-SAMPLE]



$x_i \sim \text{Normal}(\mu, \sigma)$

$$\bar{x} - \mu$$

$$5/\sqrt{N}$$

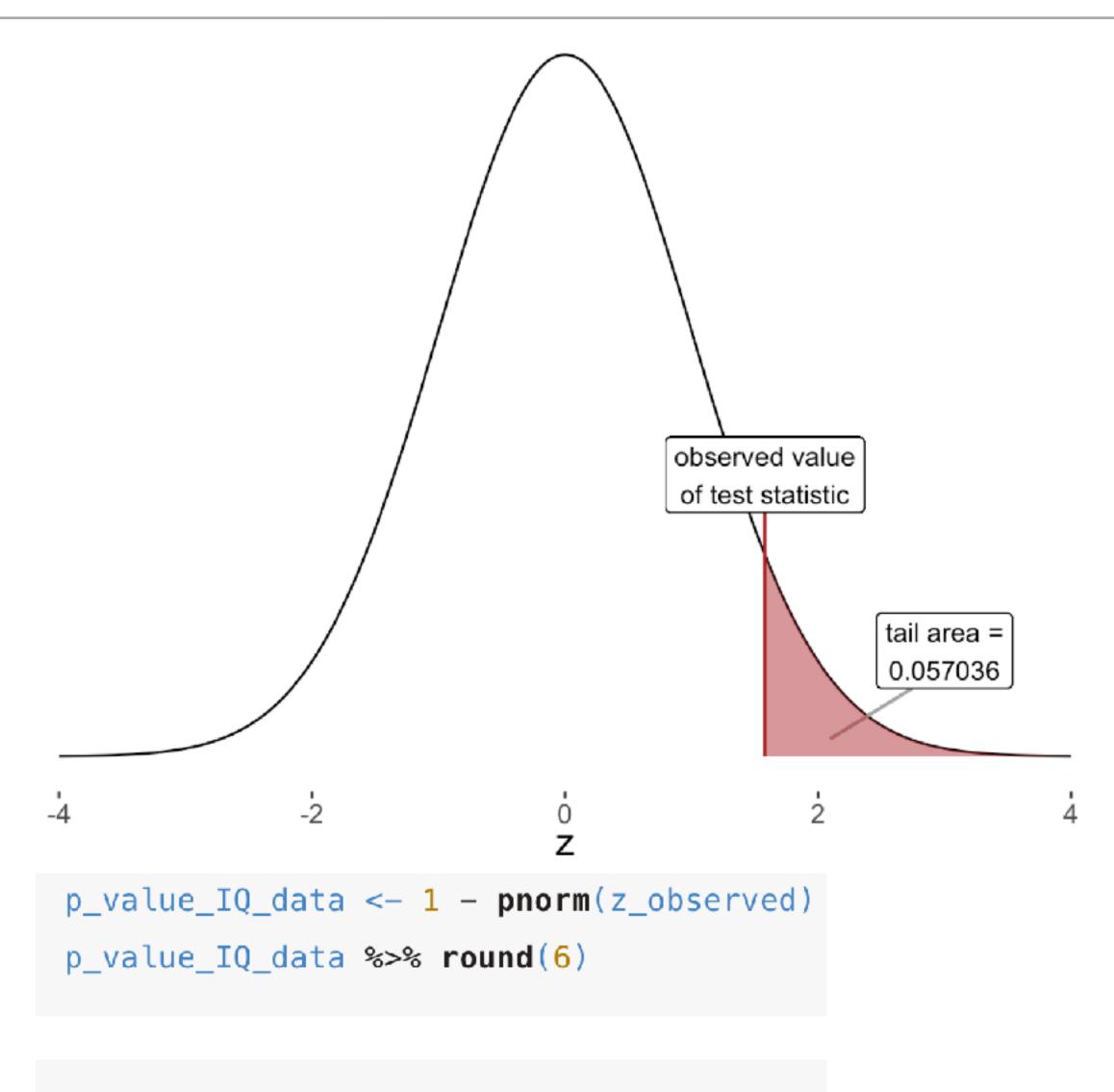
Sampling distribution: $z \sim \text{Normal}(0,1)$

FREQUENTIST Z-TEST [APPLICATION]

 $x_{i} \sim \text{Normal}(\mu, \sigma)$ $z = \frac{\bar{x} - \mu}{\sigma/\sqrt{N}}$ $z \sim \text{Normal}(0, 1)$

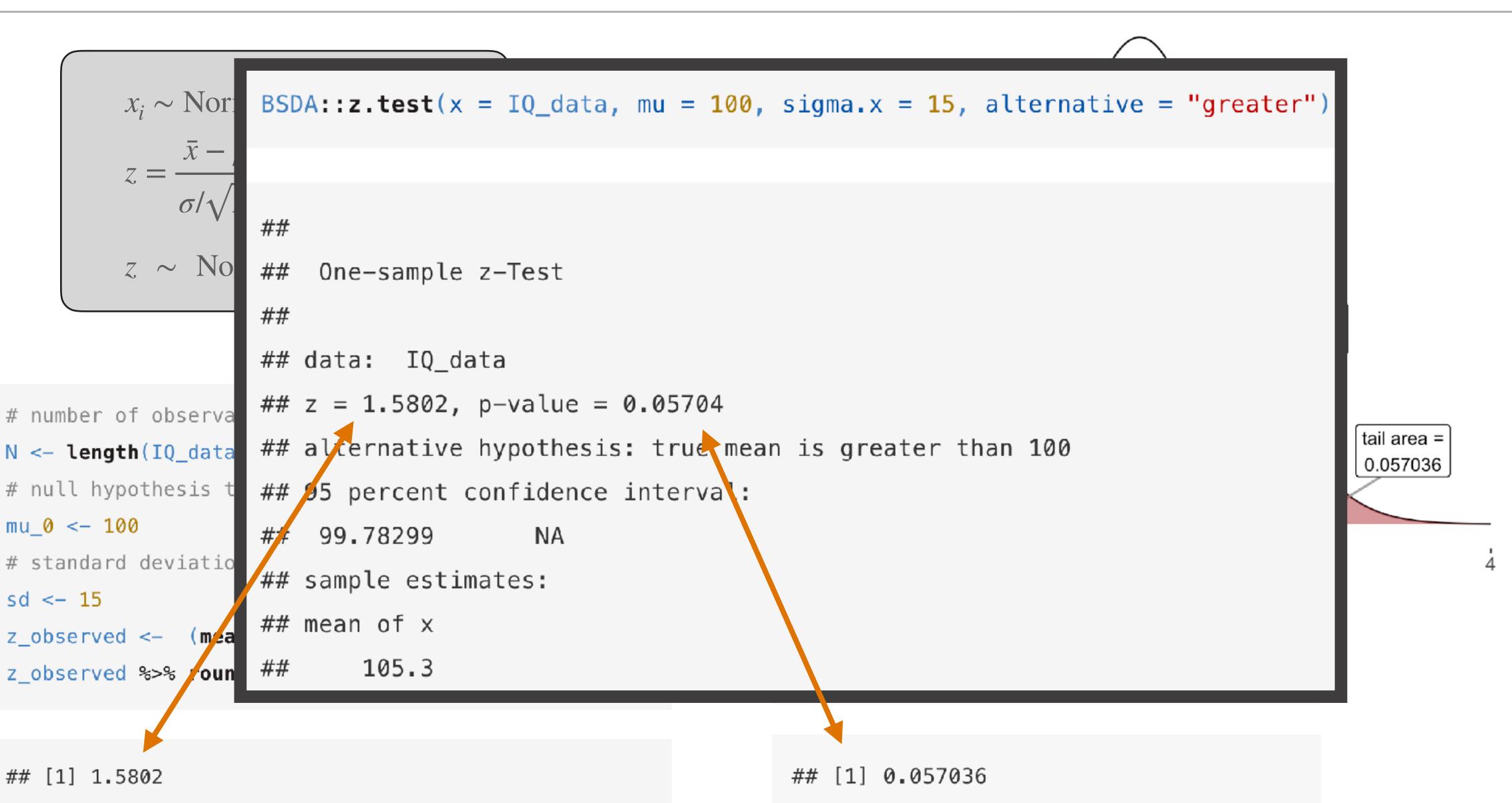
```
# number of observations
N <- length(IQ_data)
# null hypothesis to test
mu_0 <- 100
# standard deviation (known/assumed as true)
sd <- 15
z_observed <- (mean(IQ_data) - mu_0) / (sd / sqrt(N))
z_observed %>% round(4)
```

[1] 1.5802



[1] 0.057036

FREQUENTIST Z-TEST [APPLICATION]

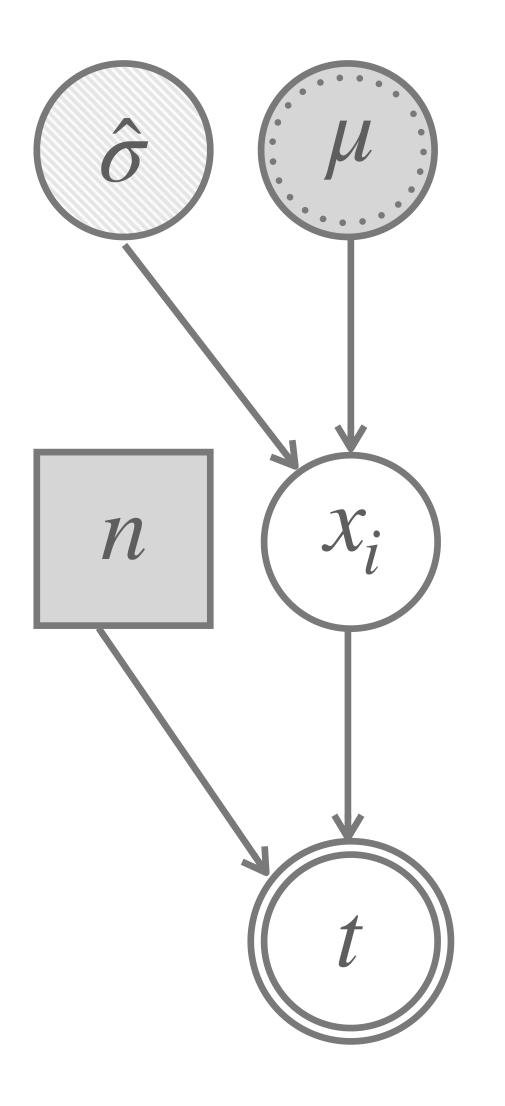


[1] 1.5802



one-sample t-test

FREQUENTIST T-TEST MODEL [ONE-SAMPLE]





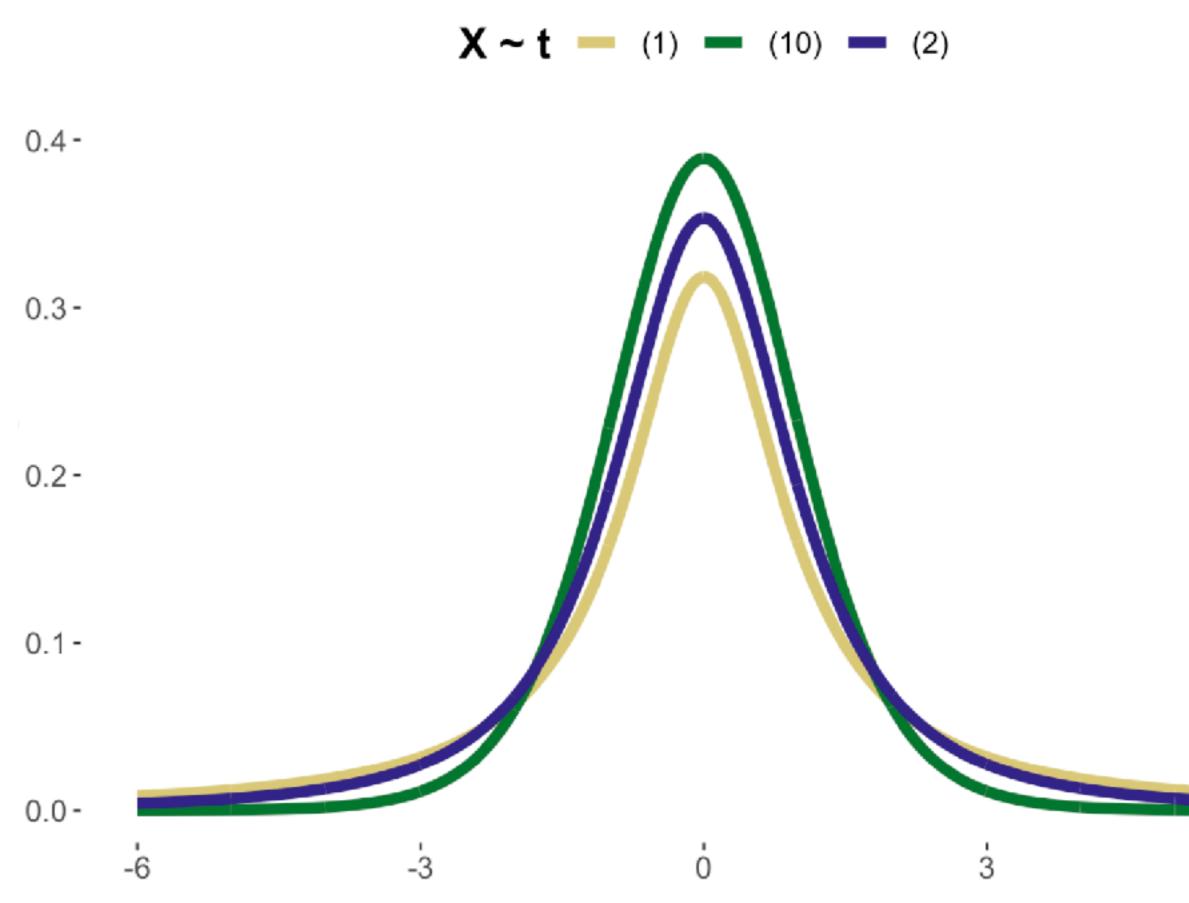
$x_i \sim \text{Normal}(\mu, \sigma)$

$$\left| \frac{1}{n-1} \sum_{i=1}^{n} (x_i - \mu_{\vec{x}})^2 \right|$$

$$-\mu_0$$

Sampling distribution: $t \sim \text{Student-t}(\nu = n - 1)$

two random variables: $x \sim \text{Normal}(0,1)$ $y \sim \chi^2$ -distribution(*n*) derived RV: $Z = \frac{X}{\sqrt{Y/n}}$ it follows (by construction) that: $z \sim \text{Student-t}(\nu = n - 1)$





FREQUENTIST T-TEST [APPLICATION]

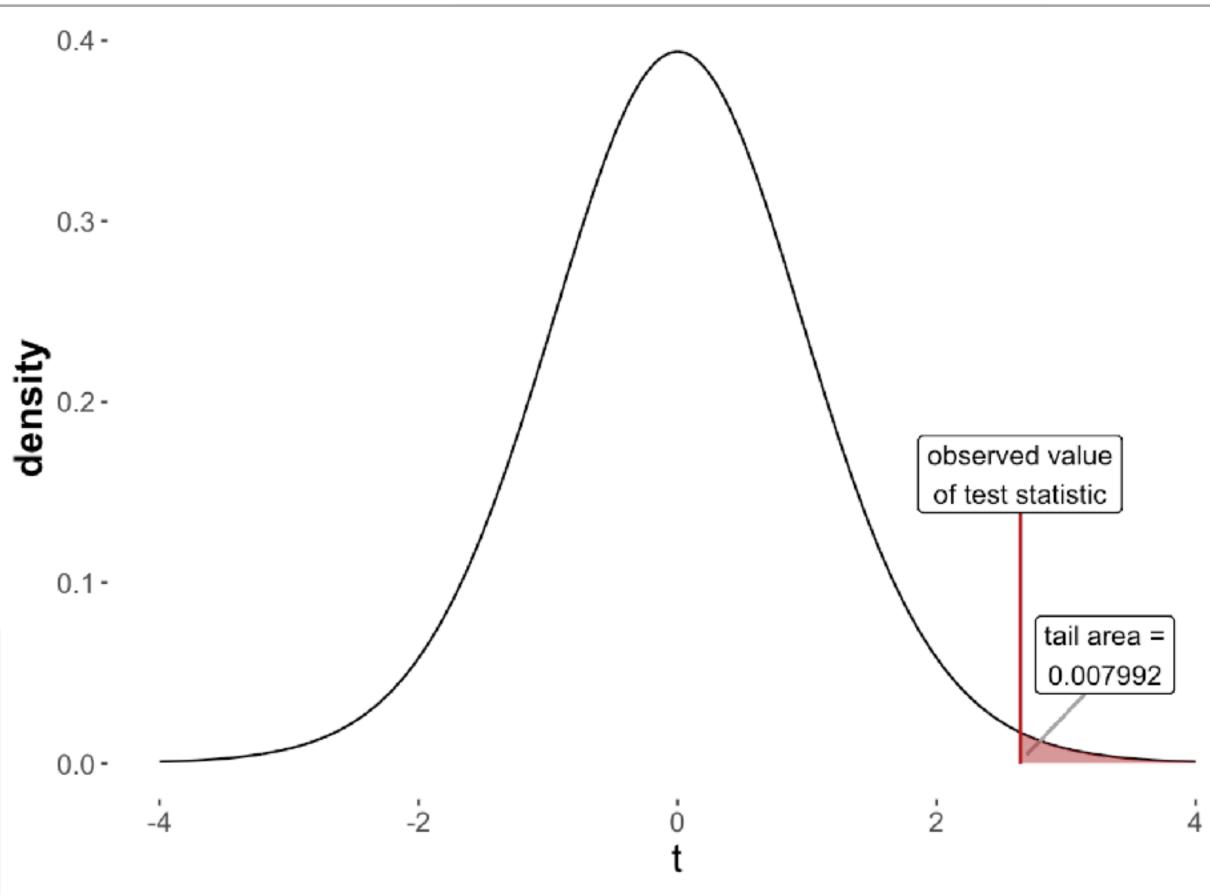
$$x_{i} \sim \text{Normal}(\mu, \sigma)$$

$$\hat{\sigma} = \sqrt{\frac{1}{n-1} \sum_{i=1}^{n} (x_{i} - \mu_{\overrightarrow{x}})^{2}}$$

$$t = \frac{\overline{x} - \mu_{0}}{\hat{\sigma}/\sqrt{n}}$$

$$t \sim \text{Student-t}(\nu = n - 1)$$

[1] 2.6446



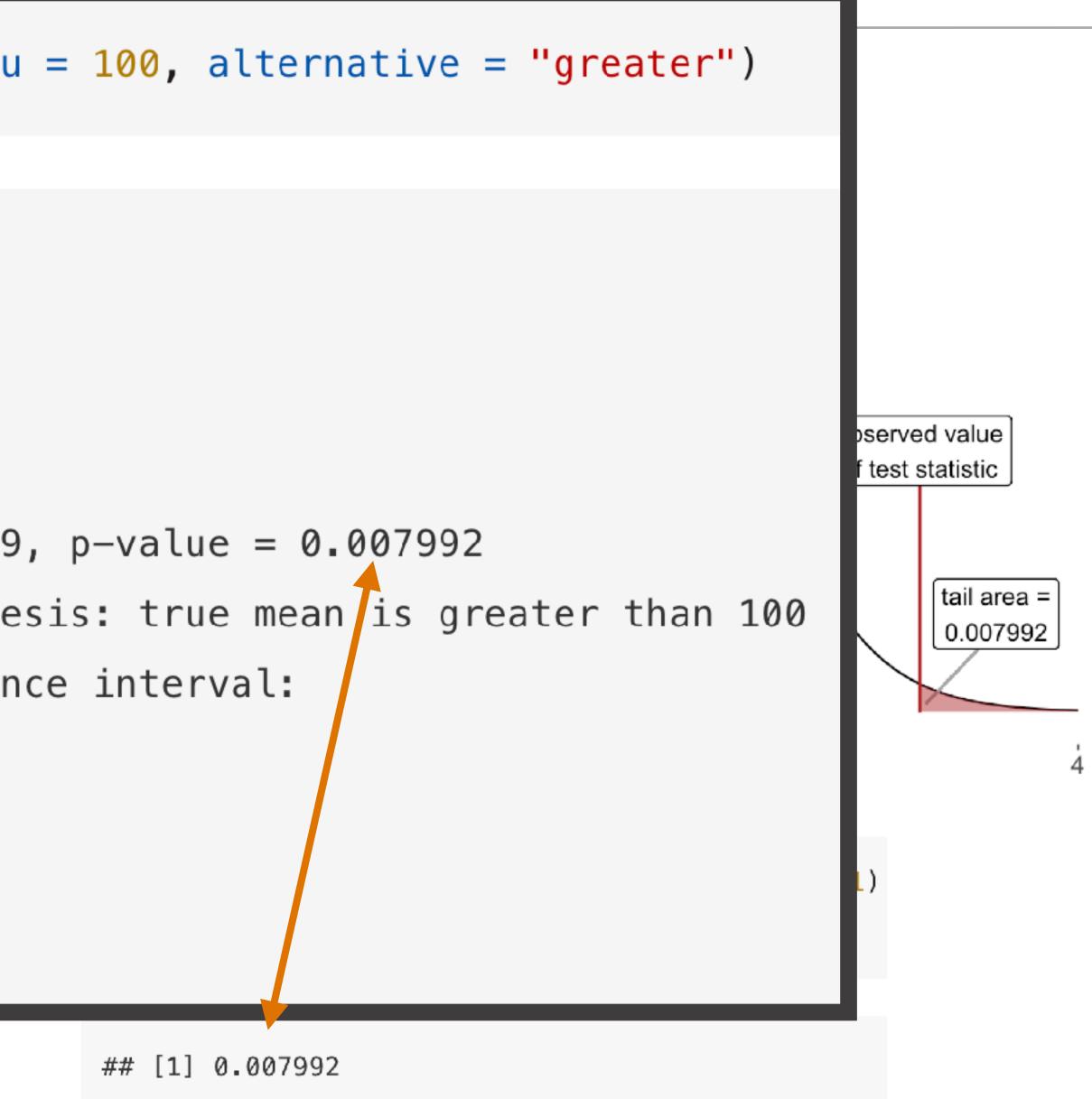
p_value_t_test_IQ <- 1 - pt(t_observed, df = N-1)
p_value_t_test_IQ %>% round(6)

[1] 0.007992

FREQUENTIST T-TEST [APPLICATION]

	<pre>t.test(x = IQ_data, mu</pre>
$x_i \sim Norma$	
$\hat{\sigma} = \sqrt{\frac{1}{n-1}}$	##
$t = \frac{\bar{x} - \mu_0}{\hat{\sigma}/\sqrt{n}}$	## One Sample t-test ##
$t \sim \text{Student}$	## data: IQ_data
	## t = 2.6446, df = 19
N <- length(IQ_data)	<pre>## alternative hypothe</pre>
<pre># fix the null hypothesis</pre>	<pre>## 95 percent confiden</pre>
<pre>mean_0 <- 100 # unlike in a z-test we use</pre>	## 101.8347 Inf
sigma_hat <- sd (IQ_data)	## sample estimates:
t_observed <- (mean (JQ_data)	## mean of x
t_observed %>% round (4)	## 105.3

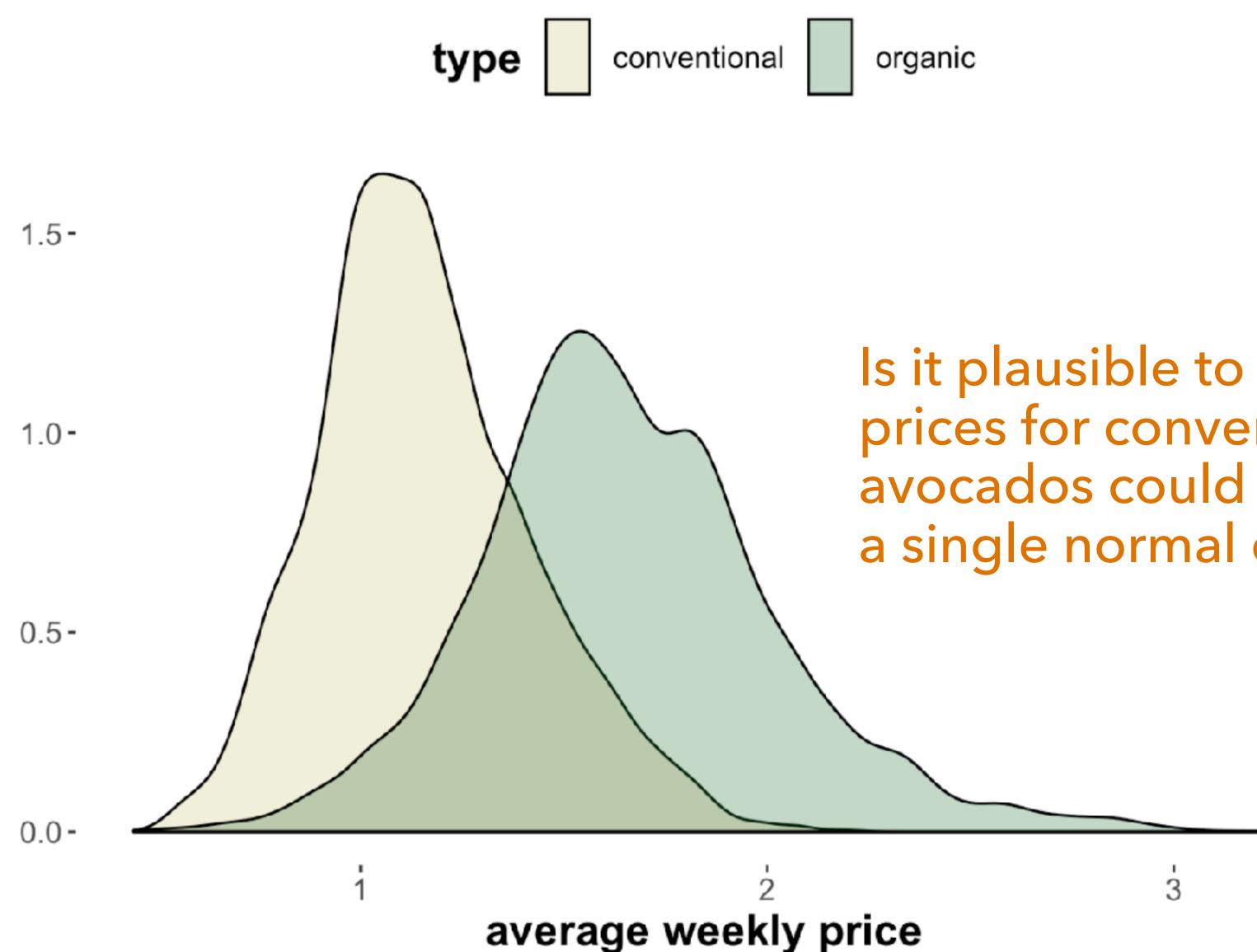
[1] 2.6446

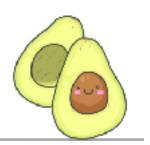




two-sample (unpaired cata, equal variance & unequal sample size)

COMPARING TWO GROUPS OF METRIC MEASURES



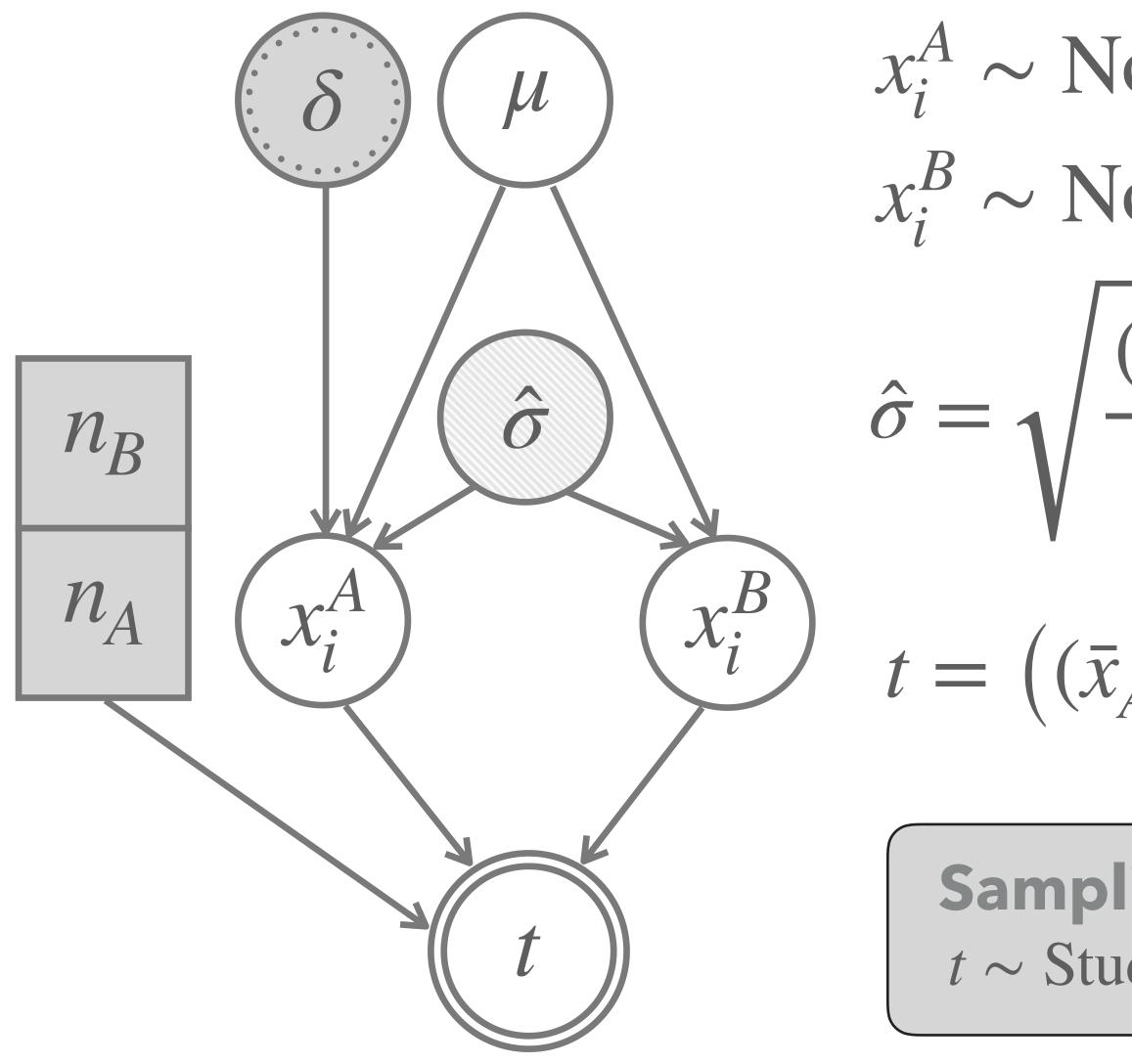


Is it plausible to assume that the observed prices for conventional and organic avocados could have been generated by a single normal distribution?





FREQUENTIST T-TEST MODEL [TWO-SAMPLE, UNPAIRED, EQUAL VARIANCE, UNEQUAL SAMPLE SIZES]



$$A^{A} \sim \text{Normal}(\mu + \delta, \sigma)$$

$$B^{B} \sim \text{Normal}(\mu, \sigma)$$

$$= \sqrt{\frac{(n_{A} - 1)\hat{\sigma}_{A}^{2} + (n_{B} - 1)\hat{\sigma}_{B}^{2}}{n_{A} + n_{B} - 2}} \left(\frac{1}{n_{A}} + \frac{1}{n_{B}}\right)$$

$$= \left((\bar{x}_{A} - \bar{x}_{B}) - \delta\right) \cdot \frac{1}{\hat{\sigma}}$$
Sampling distribution:

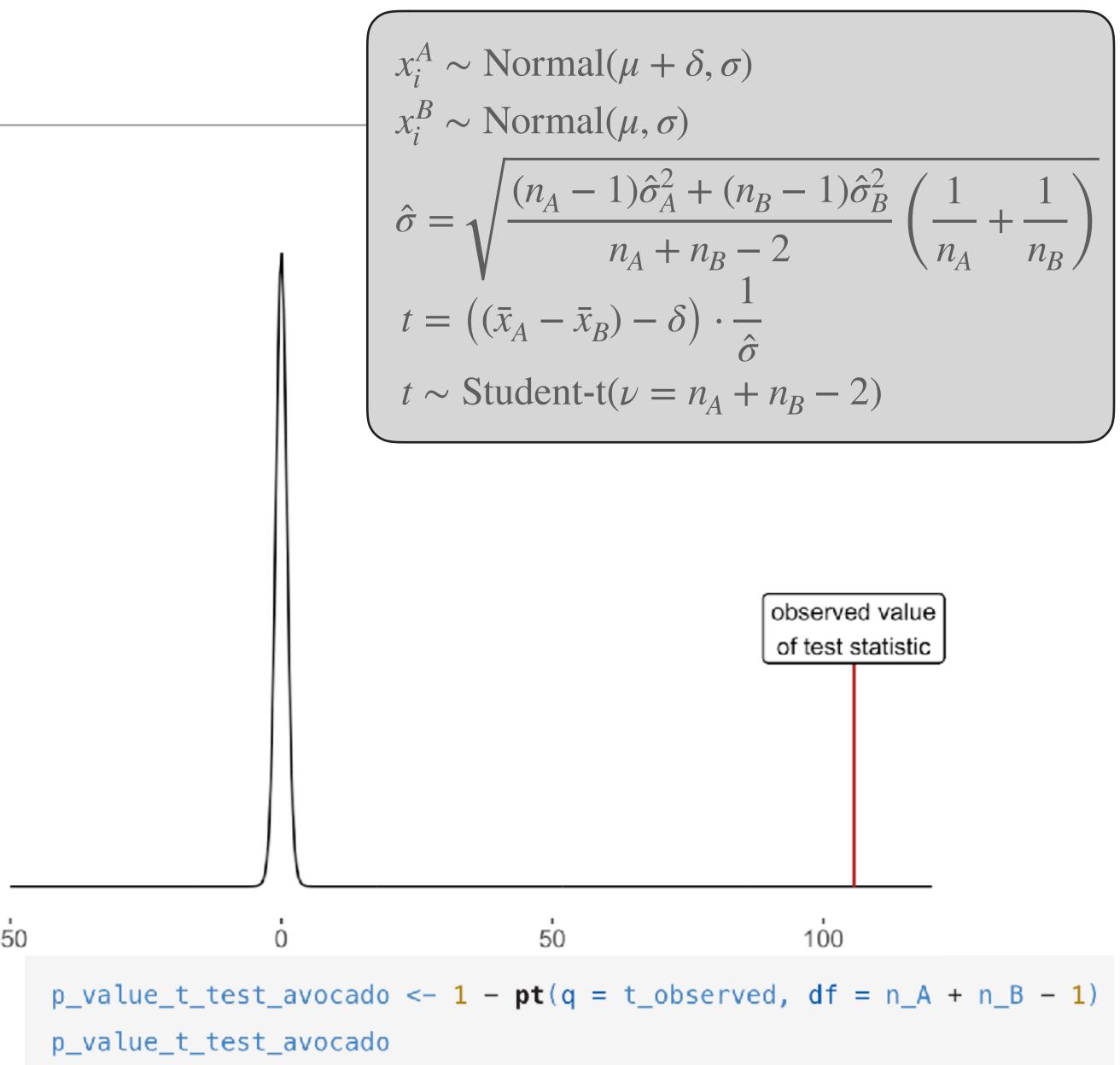
$$t \sim \text{Student-t}(\nu = n_{A} + n_{B} - 2)$$



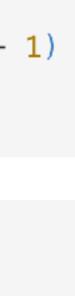


TWO-SAMPLE T-TEST EXAMPLE

# fix the null hypothesis: no difference between groups	
delta_0 <- 0	
# data (group A)	0.4 -
x_A <- avocado_data %>%	
<pre>filter(type == "organic") %>% pull(average_price)</pre>	
# data (group B)	
x_B <- avocado_data %>%	0.3 -
<pre>filter(type == "conventional") %>% pull(average_price)</pre>	
<pre># sample mean for organic (group A)</pre>	
<pre>mu_A <- mean(x_A)</pre>	0.2 -
<pre># sample mean for conventional (group B)</pre>	0.2
<pre>mu_B <- mean(x_B)</pre>	
<pre># numbers of observations</pre>	
<pre>n_A <- length(x_A)</pre>	0.1-
<pre>n_B <- length(x_B)</pre>	
# variance estimate	
sigma_AB <- sqrt(
$(((n_A -1) * sd(x_A)^2 + (n_B -1) * sd(x_B)^2) /$	0.0-
$(n_A + n_B - 2)) * (1/n_A + 1/n_B)$	-{
)	
t_observed <- (mu_A - mu_B - delta_0) / sigma_AB	
t_observed	



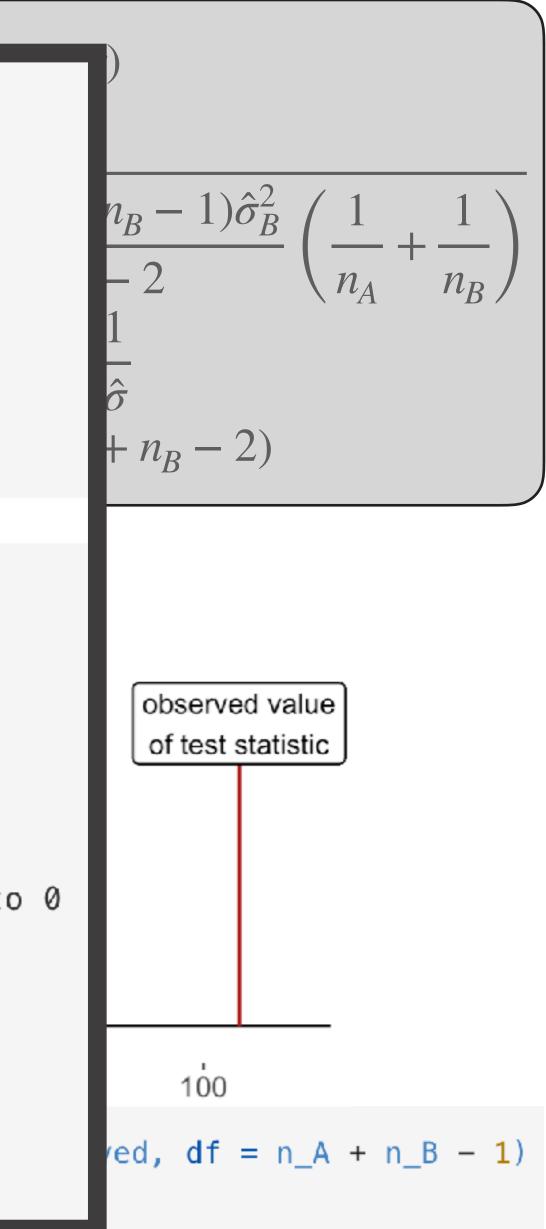
[1] 0

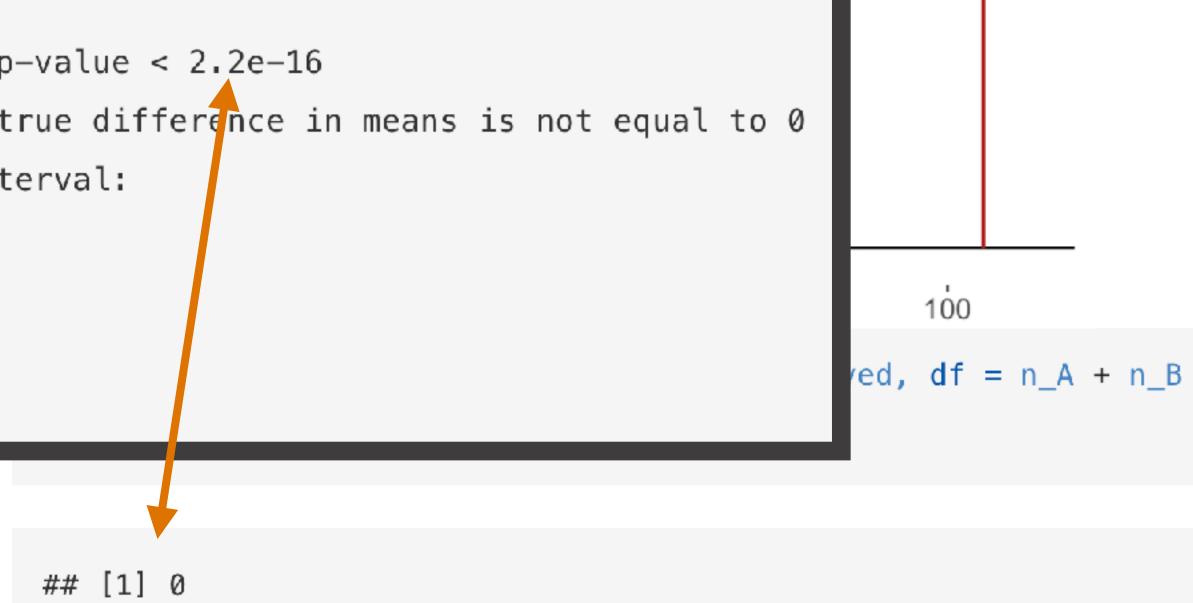


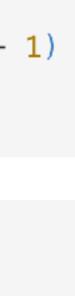
TWO-SAMPLE T-TE	t.test(
	<pre>x = x_A, # first</pre>
<pre># fix the null hypothesis: no di delta_0 <- 0</pre>	$y = x_B$, # sec v
# data (group A)	<pre>paired = FALSE, # measu</pre>
x_A <- avocado_data %>%	<pre>var.equal = TRUE, # we as</pre>
<pre>filter(type == "organic") %>%</pre>	mu = 0 # NH is
# data (group B)	
x_B <- avocado_data %>%	·
<pre>filter(type == "conventional")</pre>	
<pre># sample mean for organic (group</pre>	
mu_A <- mean(x_A)	##
<pre># sample mean for conventional (</pre>	## Two Sample t-test
mu_B <- mean(x_B)	##
<pre># numbers of observations</pre>	## data: x_A and x_B
<pre>n_A <- length(x_A)</pre>	## t = 105.59, df = 18247, p
<pre>n_B <- length(x_B)</pre>	<pre>## alternative hypothesis: t</pre>
<pre># variance estimate</pre>	
sigma_AB <- sqrt(<pre>## 95 percent confidence int</pre>
$(((n_A - 1) * sd(x_A)^2 + (n_P))$	## 0.4867522 0.5051658
$(n_A + n_B - 2)) * (1/p_A$	<pre>## sample estimates:</pre>
)	## mean of x mean of y
t_observed <- (mu_A - mu_B - del	## 1.653999 1.158040
t_observed	

[1] 105.5878

t vector of data measurements
vector of data measurements
urements are to be treated as unpaired
ssum equal variance in both groups
s delta = 0 (name 'mu' is misleading!)

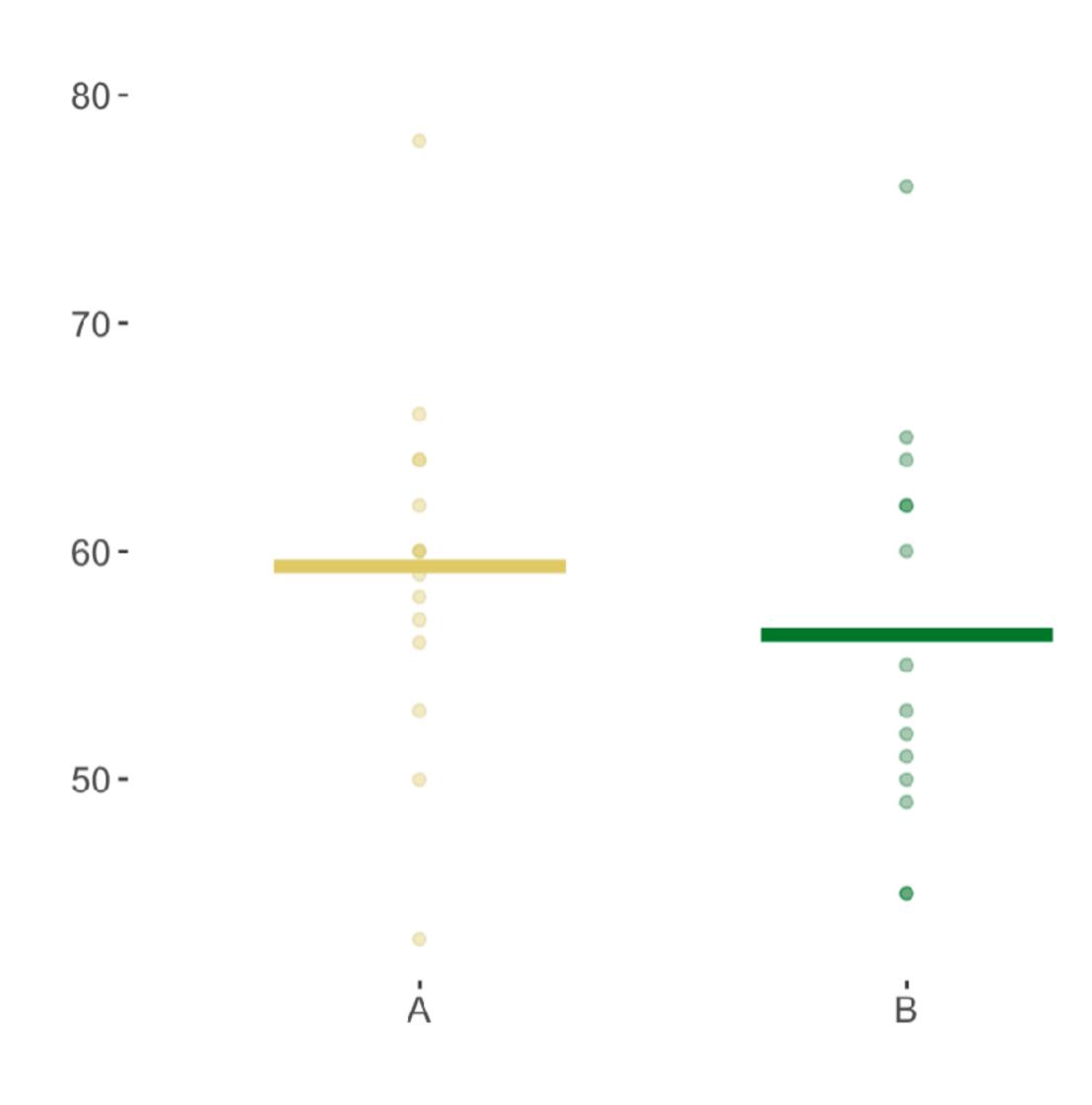








COMPARING K \geq 2 **GROUPS OF METRIC MEASURES**

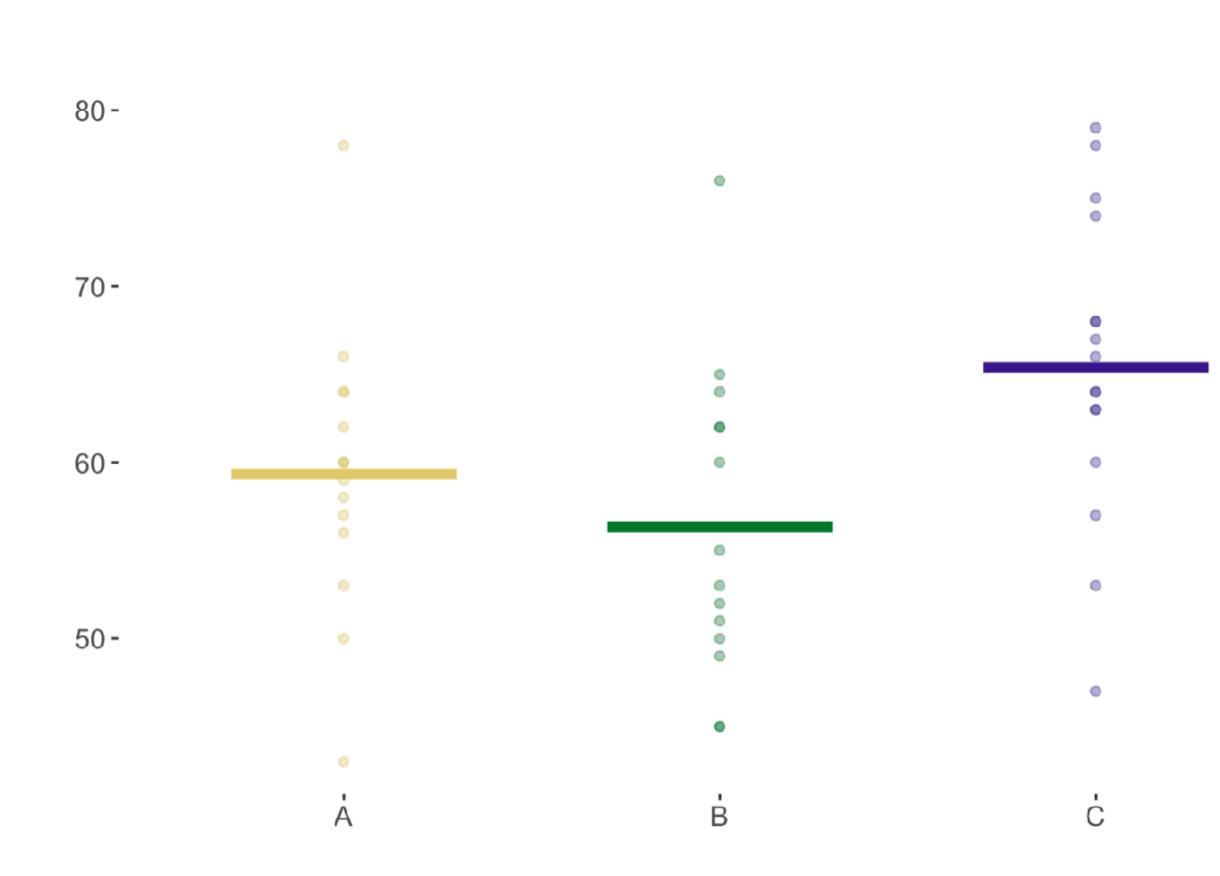


۲

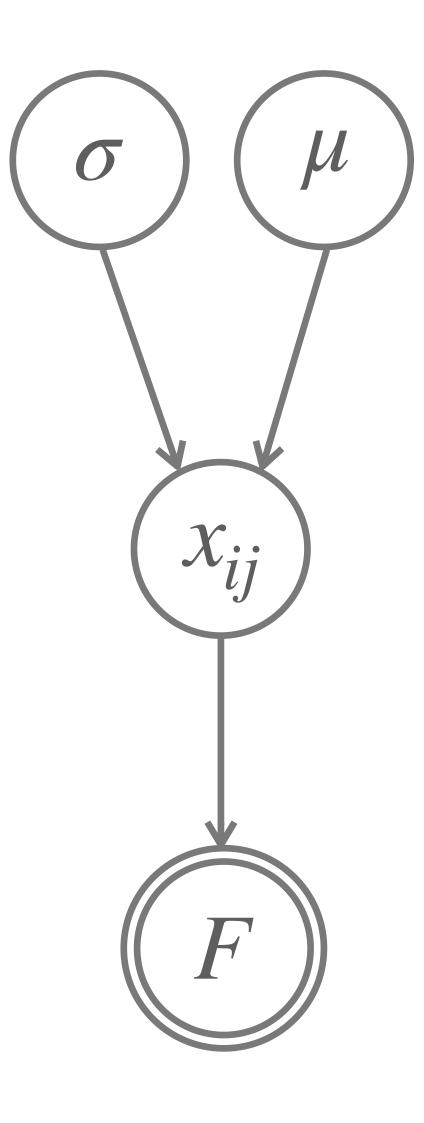
Is it plausible to assume that these measures stem from the same normal distribution?



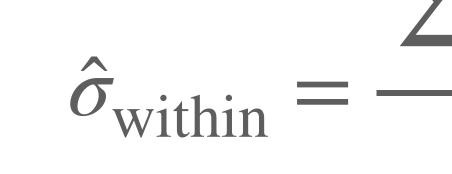
- we could run *t*-tests between different groups
- chance of α error rises with
 each comparison
 - common corrections apply
- gets tedious with large k



FREQUENTIST MODEL FOR ANOVA [ONE-WAY]



 $x_{ij} \sim \text{Norma}$





Sampling d

 $F \sim F$ -distribu

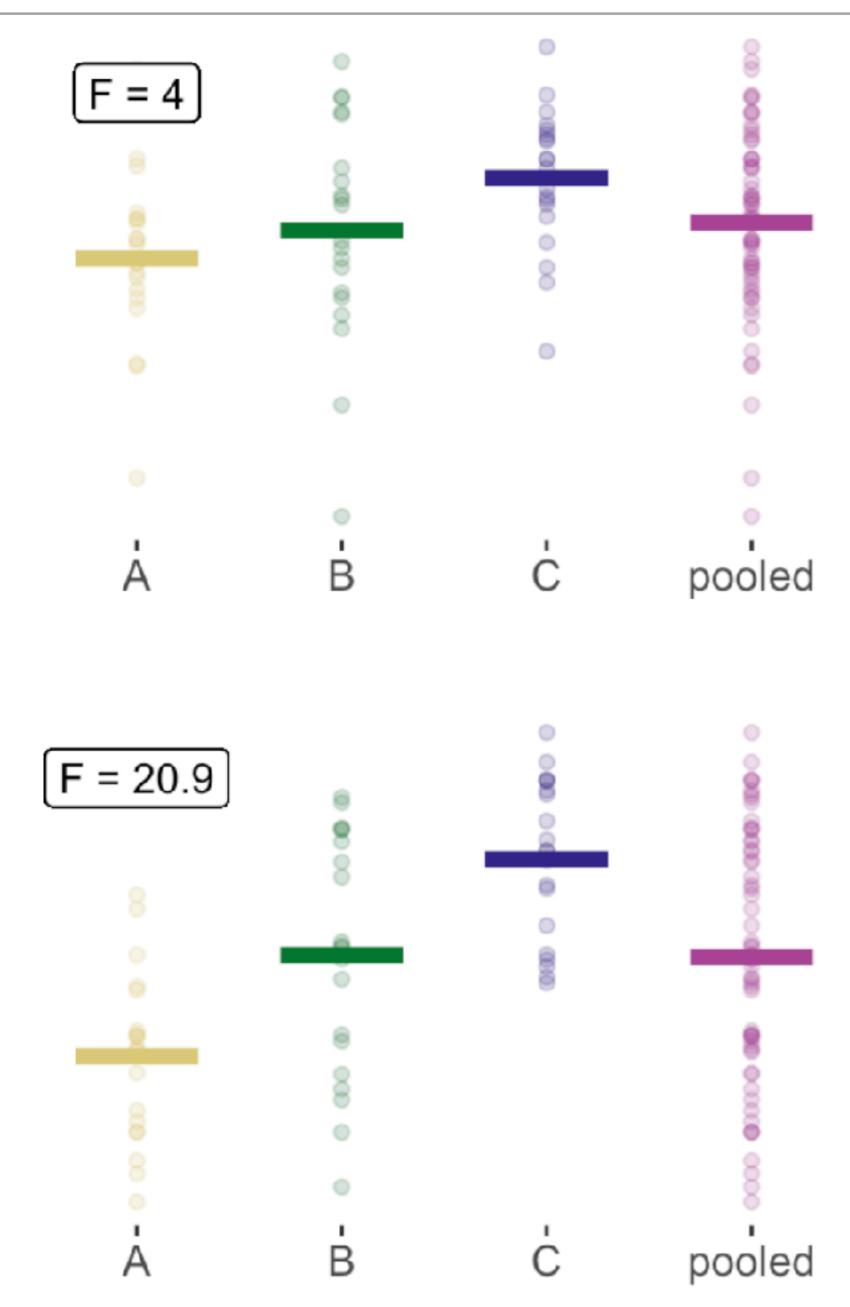
al
$$(\mu, \sigma)$$

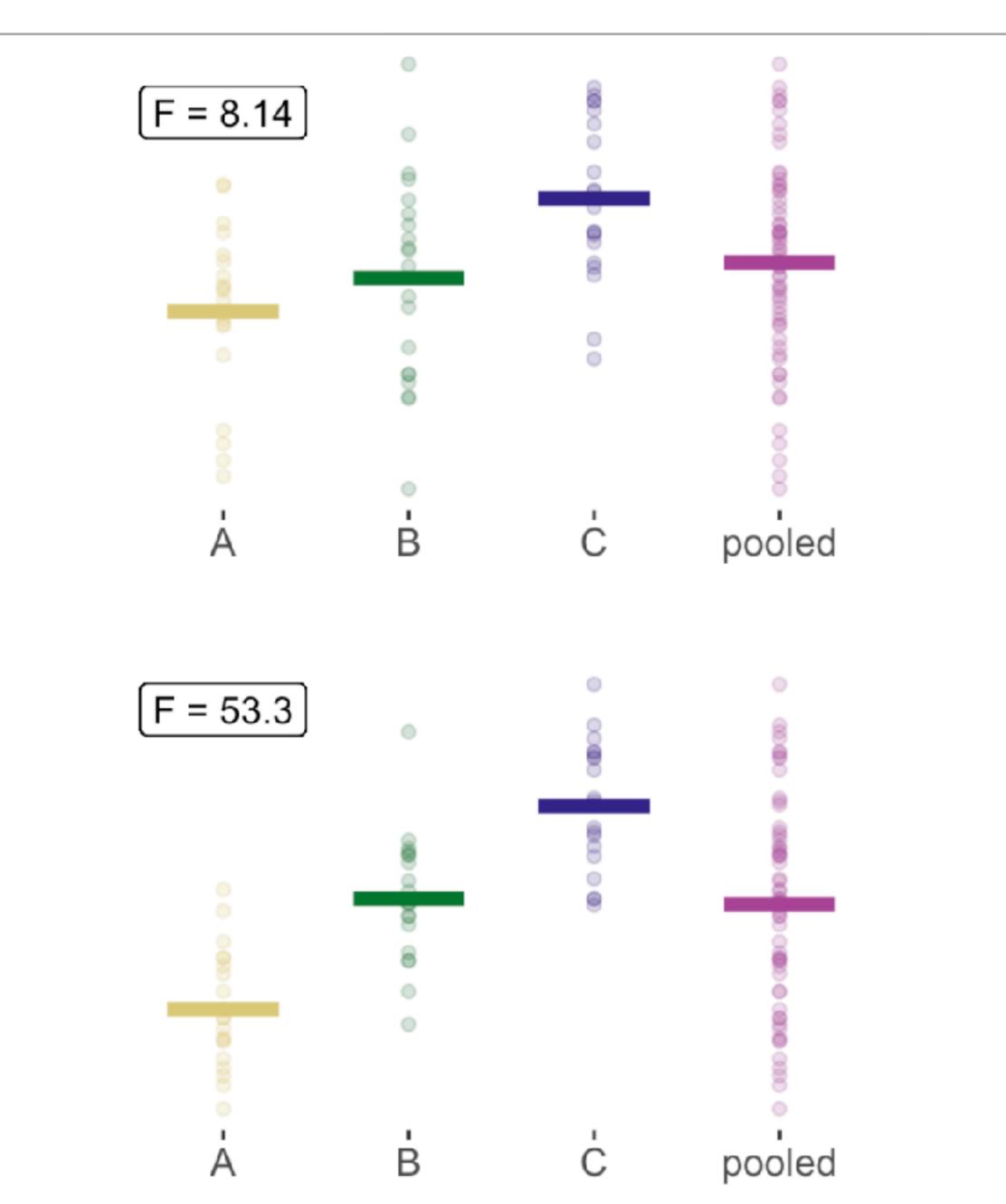
$$F = \frac{\hat{\sigma}_{between}}{\hat{\sigma}_{within}}$$

$$\frac{\sum_{j=1}^{k} \sum_{i=1}^{n_j} (x_{ij} - \bar{x}_j)^2}{\sum_{i=1}^{k} (n_i - 1)}$$

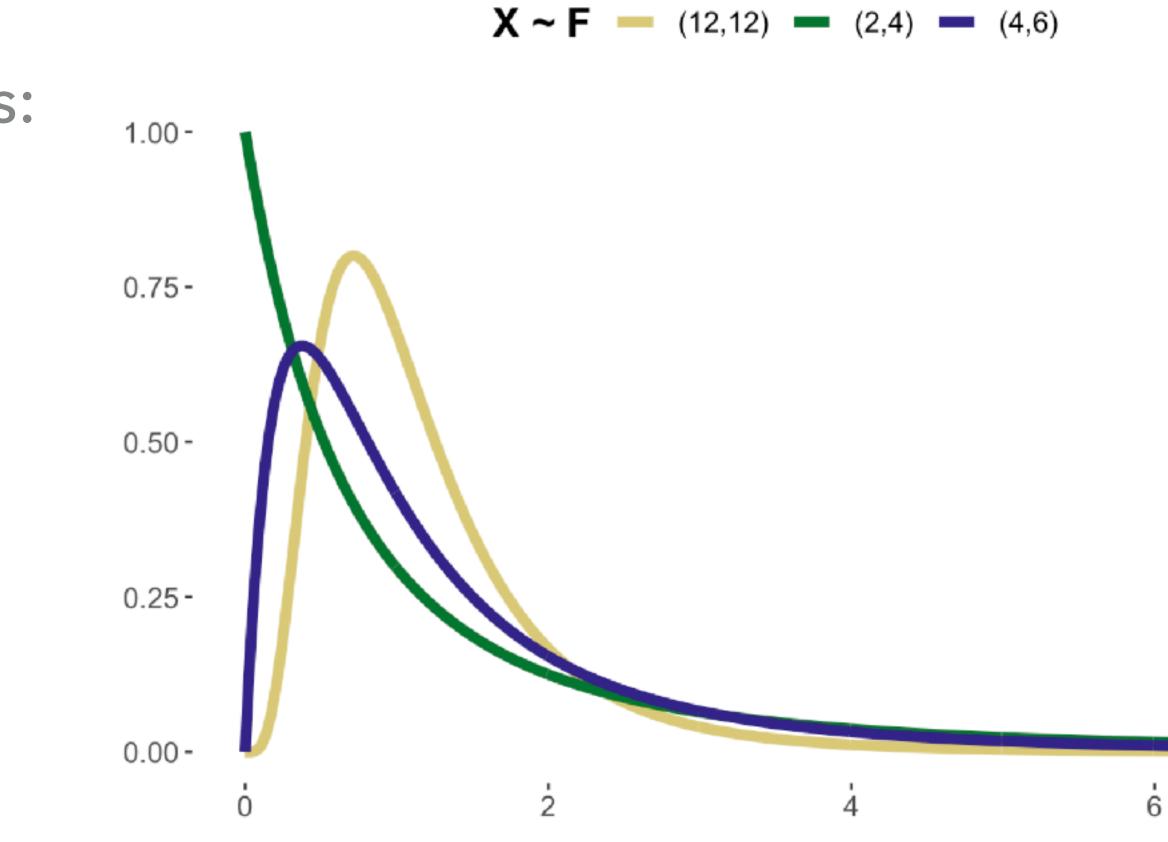
$$\frac{\sum_{j=1}^{k} n_j (\bar{x}_j - \bar{x})^2}{k - 1}$$
istribution:
ation $\left(k - 1, \sum_{i=1}^{k} (n_i - 1)\right)$

F-STATISTIC EXAMPLES

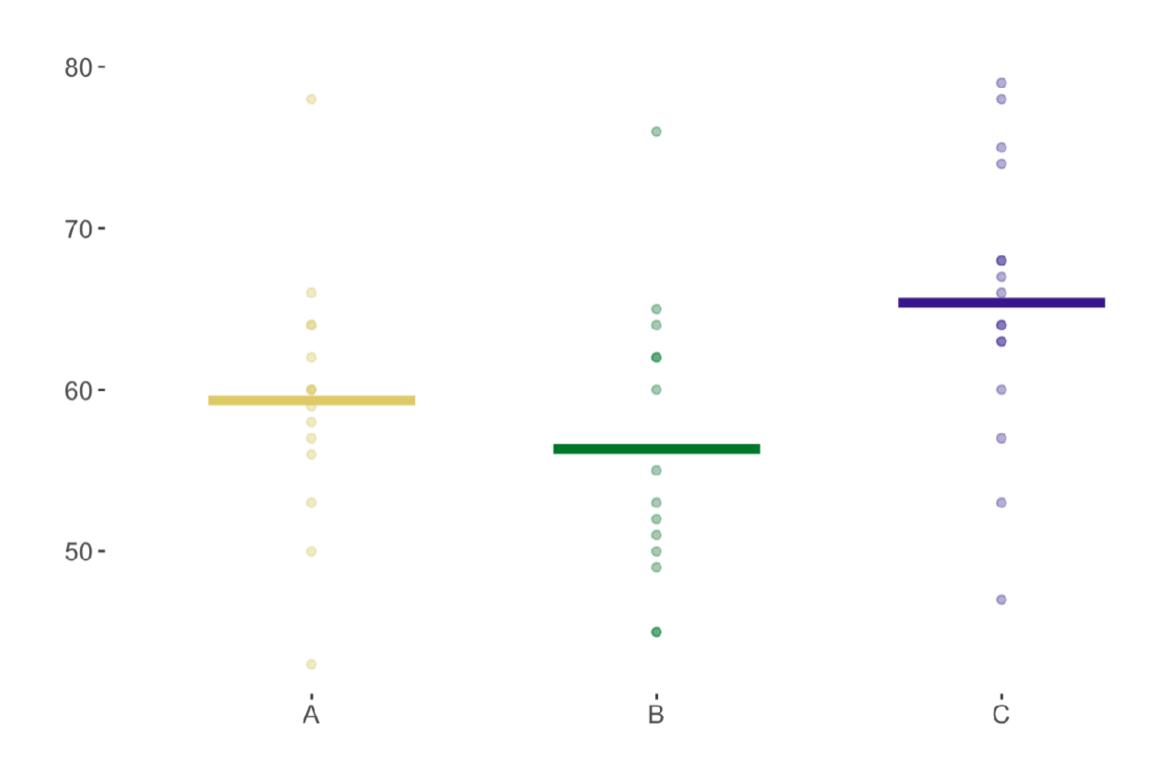




two χ²-distributed random variables: x ~ χ²-distribution(m) y ~ χ²-distribution(n) derived RV: Z = X/m/Y/n it follows (by construction) that: z ~ F-distribution(m, n)

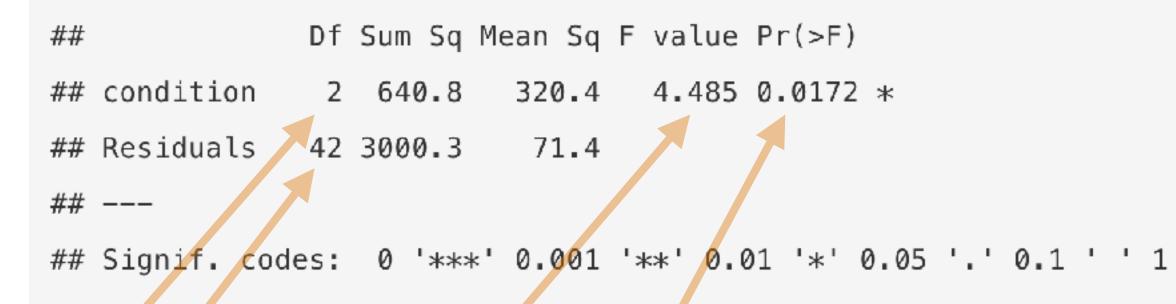


EXAMPLE



Based on a one-way ANOVA, we find evidence against the assumption of equal means across all groups ($F(2,42) \approx 4.485, \, p \approx 0.0172$.)

aov(formula = value ~ condition, anova_data) %>% summary()







varieties of frequentist testing

THREE VARIETIES OF FREQUENTIST TESTING

	FISHER	NEYMAN/PEARSON	HYBRID NHST*
explicit & serious alternative H _a	X		
when to set-up statistical model	after data collection	before data collection	after data collection
goal of statistical analysis	quantify evidence against H ₀	decide action: adopt H_0 or H_a	decide action: adopt H_0 or $\neg H_0$
power calculation			

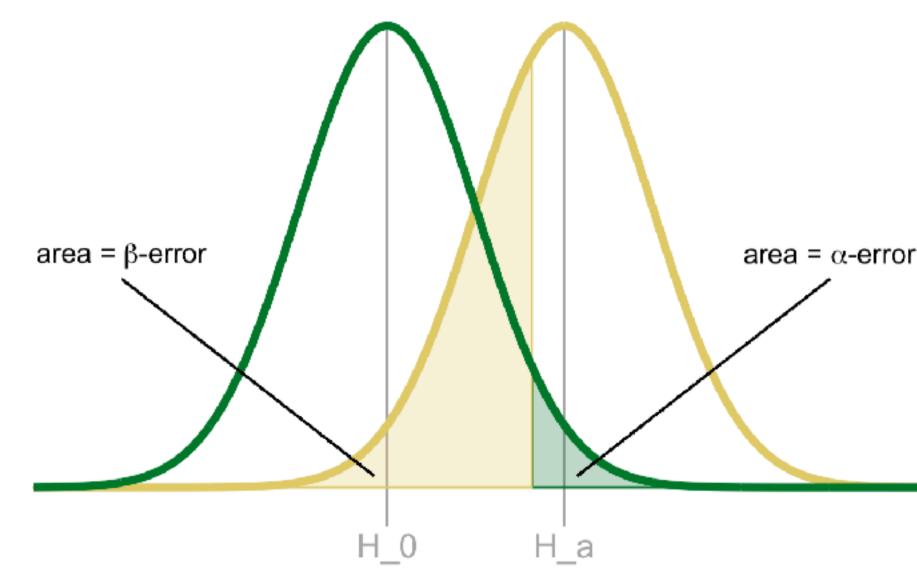
* this is a worst-case portrait of modern NHST ; this is not how it should be done





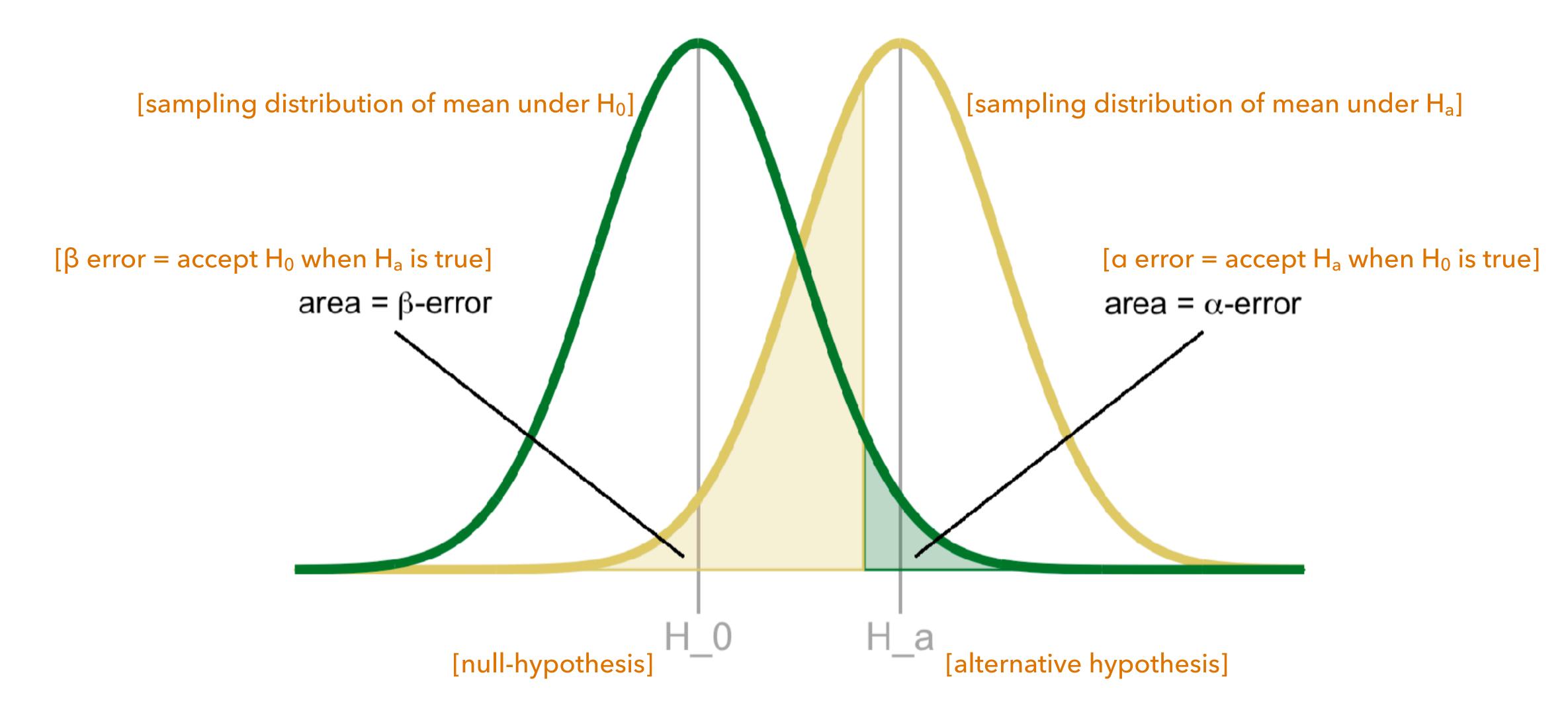
NEYMAN/PEARSON APPROACH [INFORMAL GIST]

- procedure in N/P approach:
 - fix H₀ and H_a (based on prior research)
 - determine desired α- and β-error level
 - calculate sample size N necessary for β given α
 - run the experiment
 - determine significance based on a-level
 - make a dichotomous decision:
 - accept H_a if test is significant
 - accept H₀ otherwise

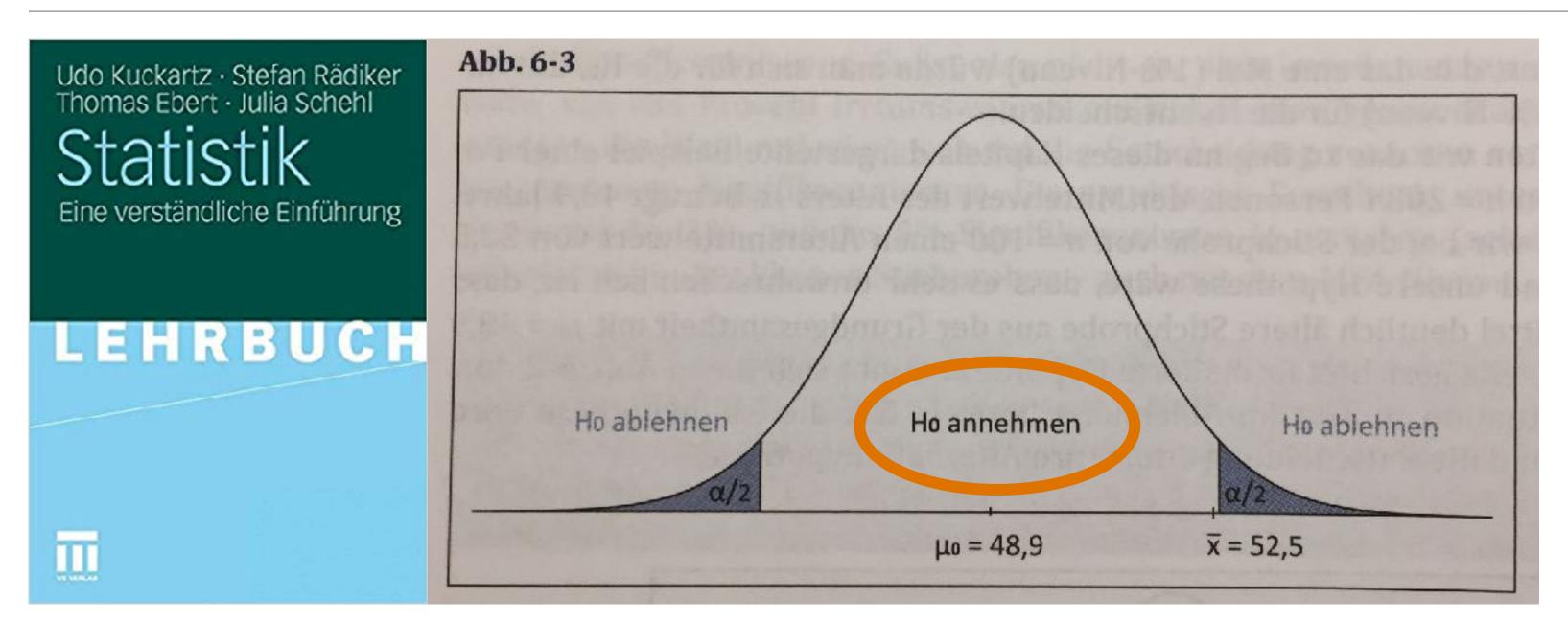


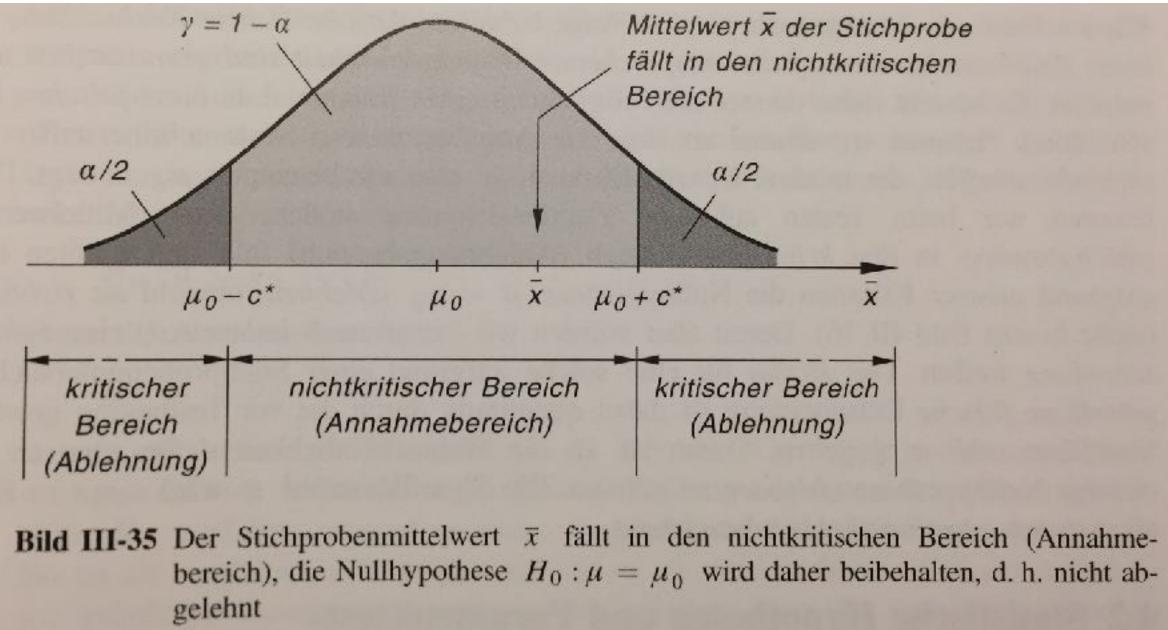
LONG-TERM ERROR CONTROL IN NEYMAN/PEARSON APPROACH

[more data = tighter curves!! = lower β]



EXAMPLES FROM TEXTBOOKS





neither textbook talks about fixing Ha and/or calculating power of a test

bereich), die Nullhypothese $H_0: \mu = \mu_0$ wird daher beibehalten, d. h. nicht ab-

Lothar Papula

Mathematik für Ingenieure und Naturwissenschaftler Band 3

Vektoranalysis, Wahrscheinlichkeitsrechnung, Mathematische Statistik, Fehler- und Ausgleichsrechnung

7. Auflage





THREE VARIETIES OF FREQUENTIST TESTING

	FISHER	NEYMAN/PEARSON	HYBRID NHST*
explicit & serious alternative H _a	X		
when to set-up statistical model	after data collection	before data collection	after data collection
goal of statistical analysis	quantify evidence against H ₀	decide action: adopt H_0 or H_a	decide action: adopt H_0 or $\neg H_0$
power calculation			

* this is a worst-case portrait of modern NHST ; this is not how it should be done





