# Review

Nicky Wakim

2023-01-08

# Lecture Overview

- Quick basics
- Important Distributions
- Statistical inference: Estimation
- Statistical inference: Hypothesis testing
- Error Rates and Power

# What did we learn in 511?

- In 511, we talked about *categorical* and *continuous* outcomes (dependent variables)

- We also talked about their relationship with 1-2 *continuous* or *categorical* exposure (independent variables or predictor)

- We had many good ways to assess the relationship between an outcome and exposure:

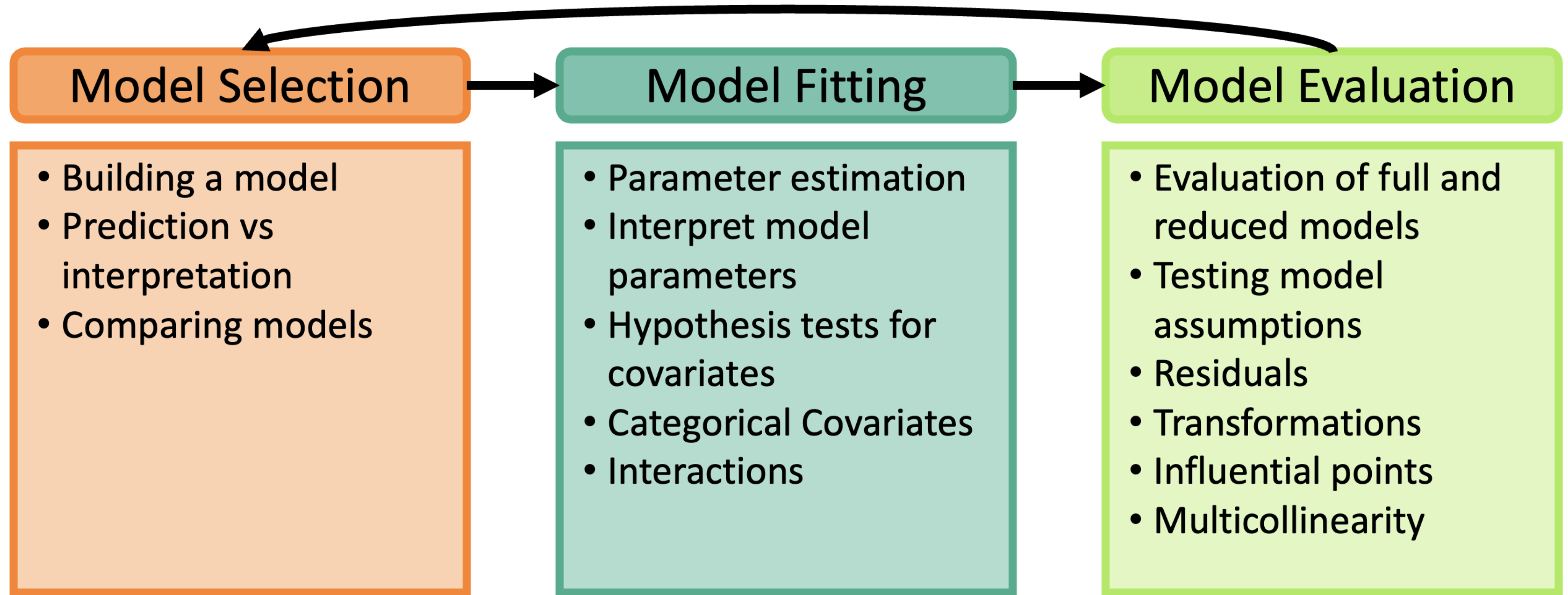|  | Continuous Outcome | Categorical Outcome |
| --- | --- | --- |
| Continuous Exposure | Correlation, simple linear regression | ?? |
| Categorical Exposure | t-tests, paired t-tests, 2 sample t-tests, ANOVA | proportion t-test, Chi-squared goodness of fit test, Fisher's Exact test, Chi-squared test of independence, etc. |

# What did we learn in 511?

- You set up a really **important foundation**

  - Including distributions, mathematical definitions, hypothesis testing, and more!

- Tests and statistical approaches learned are incredibly helpful!

- While you had to learn a lot of different tests and approaches for each combination of categorical/continuous exposure with categorical/continuous outcome

  - **Those tests cannot handle more complicated data**

- **What happens when other variables influence the relationship between your exposure and outcome?**

  - Do we just ignore them?

# What will we learn in this class?

- We will be building towards models that can handle many variables!

  - **Regression** is the building block for modeling multivariable relationships

- In Linear Models we will *build, interpret, and evaluate* linear regression models

# Process of regression data analysis

**Model Selection**

- Building a model
- Prediction vs interpretation
- Comparing models

**Model Fitting**

- Parameter estimation
- Interpret model parameters
- Hypothesis tests for covariates
- Categorical Covariates
- Interactions

**Model Evaluation**

- Evaluation of full and reduced models
- Testing model assumptions
- Residuals
- Transformations
- Influential points
- Multicollinearity

# Main sections of the course

1. Review

2. Intro to SLR: estimation and testing

   - Model fitting

3. Intro to MLR: estimation and testing

   - Model fitting

4. Diving into our predictors: categorical variables, interactions between variable

   - Model fitting

5. Key ingredients: model evaluation, diagnostics, selection, and building

   - Model evaluation and Model selection

# Main sections of the course

## 1. Review

2. Intro to SLR: estimation and testing

- Model fitting

3. Intro to MLR: estimation and testing

- Model fitting

4. Diving into our predictors: categorical variables, interactions between variable

- Model fitting

5. Key ingredients: model evaluation, diagnostics, selection, and building

- Model evaluation and Model selection

# Before we begin

- Meike has some really good online notes, code, and work on <span style="color:red">her BSTA 511 page</span>

# Learning Objectives

1. Identify important descriptive statistics and visualize data from a continuous variable

2. Identify important distributions that will be used in 512/612

3. Use our previous tools in 511 to estimate a parameter and construct a confidence interval

4. Use our previous tools in 511 to conduct a hypothesis test

5. Define error rates and power

# Learning Objectives

**1. Identify important descriptive statistics and visualize data from a continuous variable**

2. Identify important distributions that will be used in 512/612

3. Use our previous tools in 511 to estimate a parameter and construct a confidence interval

4. Use our previous tools in 511 to conduct a hypothesis test

5. Define error rates and power

# Quick basics

# Some Basic Statistics "Talk"

- Random variable $Y$
  - Sample $Y_i, i = 1, \ldots, n$
- Summation:
  $$\sum_{i=1}^{n} Y_i = Y_1 + Y_2 + \ldots + Y_n$$
- Product:
  $$\prod_{i=1}^{n} Y_i = Y_1 \times Y_2 \times \ldots \times Y_n$$

# Descriptive Statistics: continuous variables

**Measures of central tendency**

- Sample mean

$$\bar{x} = \frac{x_1 + x_2 + \ldots + x_n}{n} = \frac{\sum_{i=1}^{n} x_i}{n}$$

- Median

**Measures of variability (or dispersion)**

- Sample variance
    - Average of the squared deviations from the sample mean
- Sample standard deviation

$$s = \sqrt{\frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \ldots + (x_n - \bar{x})^2}{n - 1}} = \sqrt{\frac{\sum_{i=1}^{n} (x_i - \bar{x})^2}{n - 1}}$$

- IQR
    - Range from 1st to 3rd quartile

# Descriptive Statistics: continuous variables (R code)

## Measures of central tendency

- Sample mean

```
1  mean( sample )
```

- Median

```
1  median( sample )
```

## Measures of variability (or dispersion)

- Sample variance

```
1  var( sample )
```

- Sample standard deviation

```
1  sd( sample )
```

- IQR

```
1  IQR( sample )
```
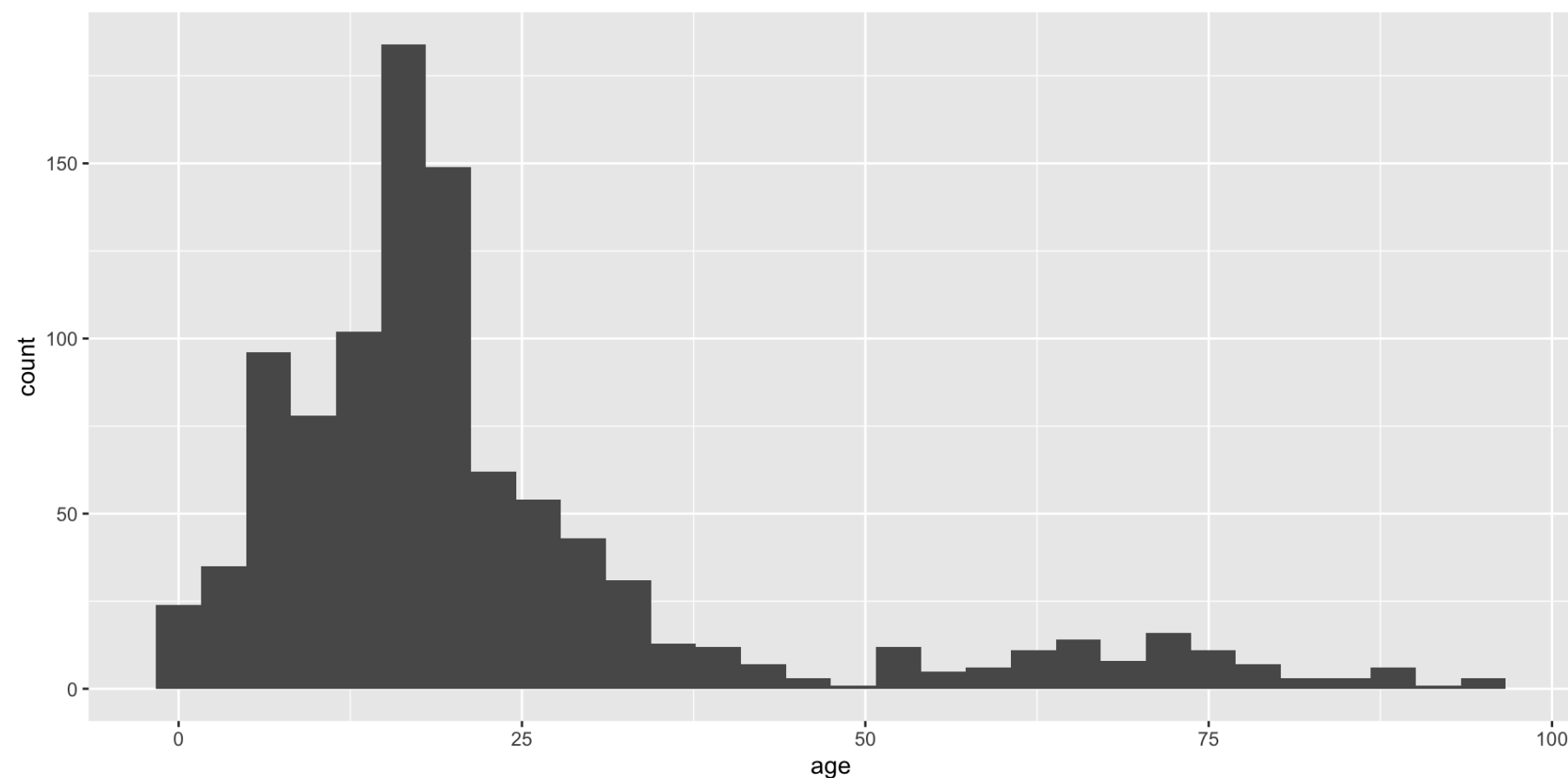
# Data visualization

- Using the library `ggplot2` to visualize data

- We will load the package:
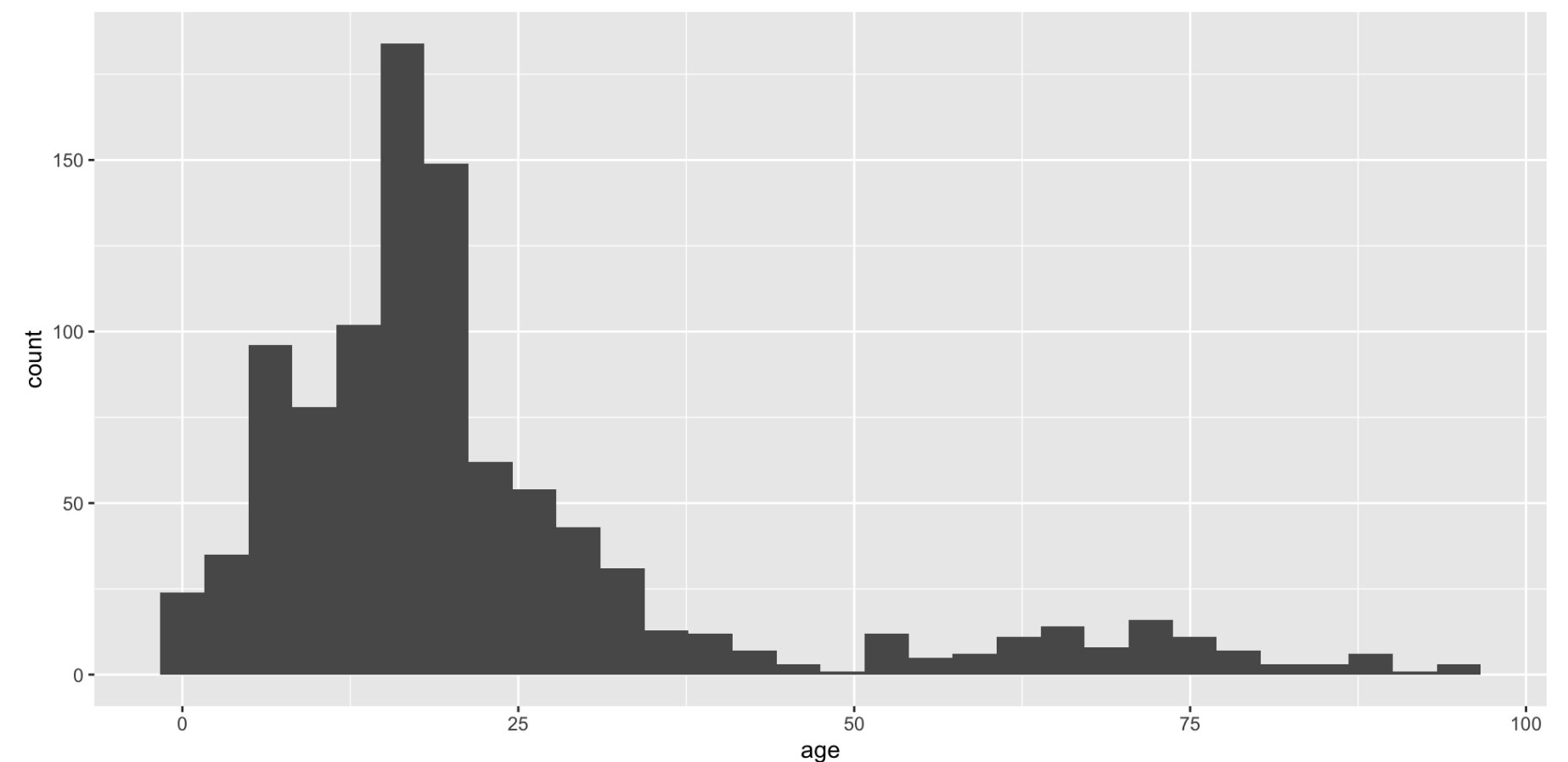
```
1  library(ggplot2)
```

# Histogram using **ggplot2**

We can make a basic graph for a continuous variable:

```
1  ggplot(data = dds.discr,
2         aes(x = age)) +
3    geom_histogram()
```

```
1  ggplot() +
2    geom_histogram(data = dds.discr,
3                   aes(x = age))
```
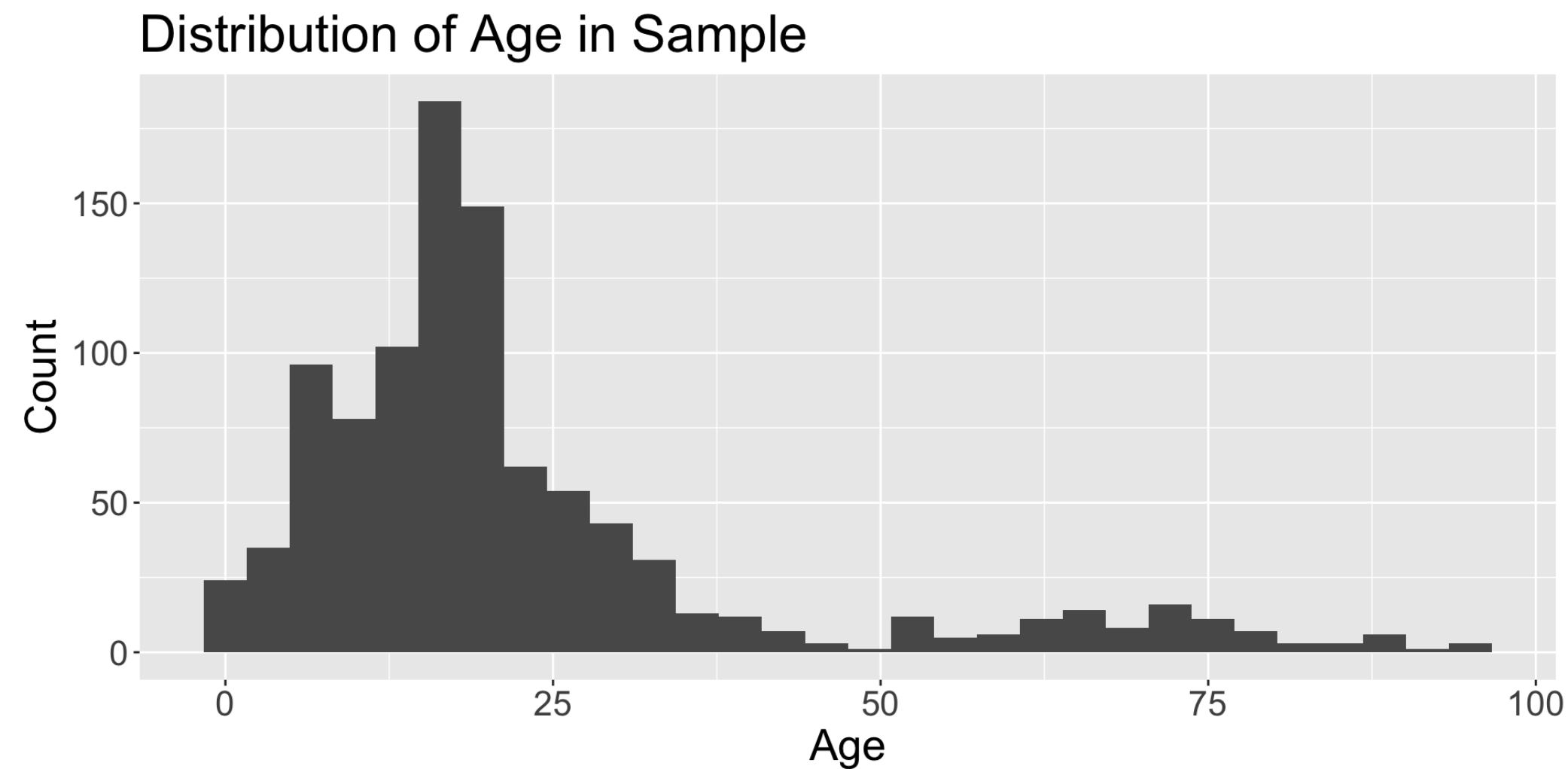




Some more information on histograms using `ggplot2`

# Spruced up histogram using `ggplot2`

We can make a more formal, presentable graph:

```
1  ggplot(data = dds.discr,
2        aes(x = age)) +
3    geom_histogram() +
4    theme(text = element_text(size=20)) +
5    labs(x = "Age",
6        y = "Count",
7        title = "Distribution of Age in Sample")
```
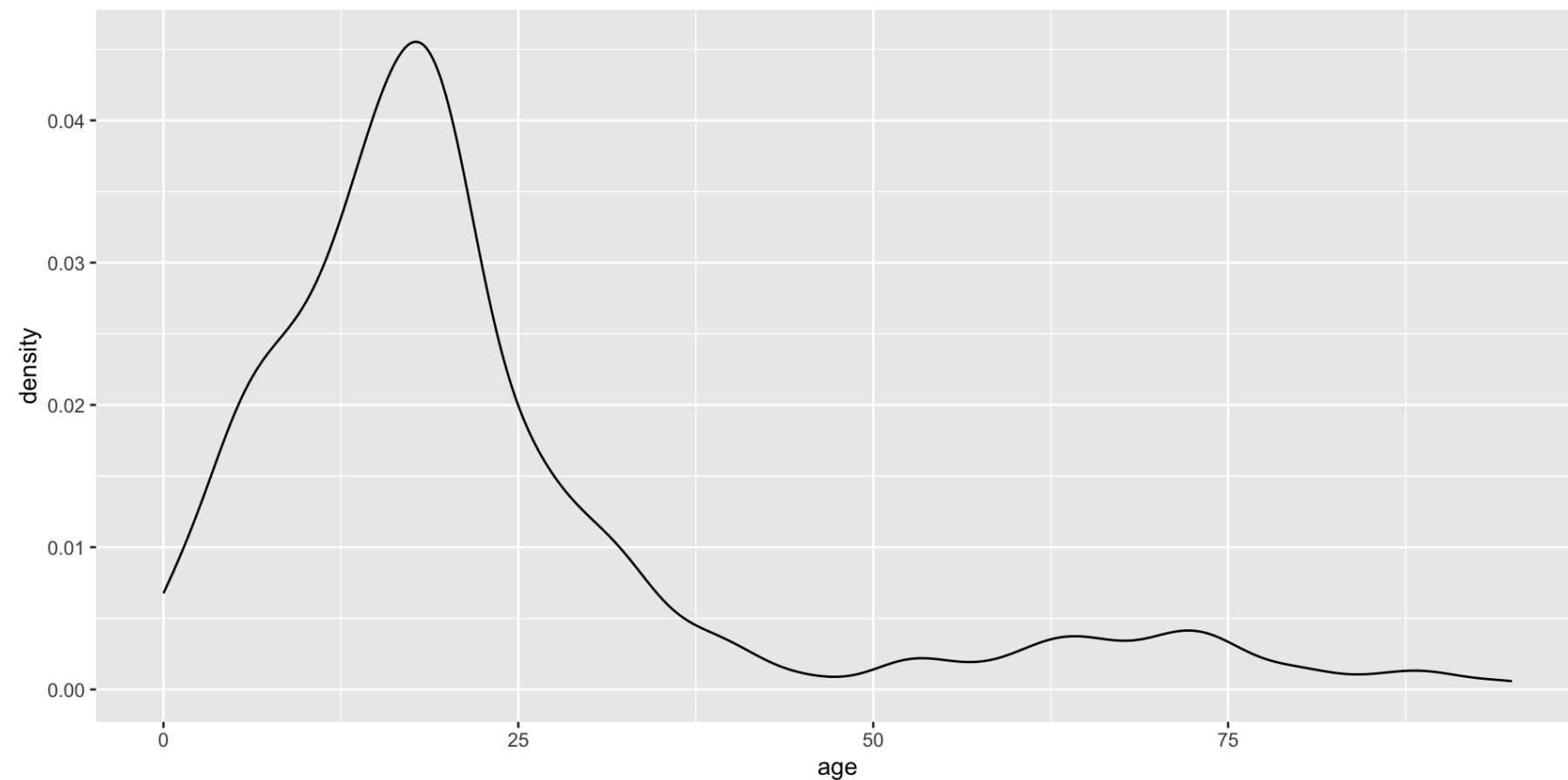


I would like you to turn in homework, labs, and project reports with graphs like these.
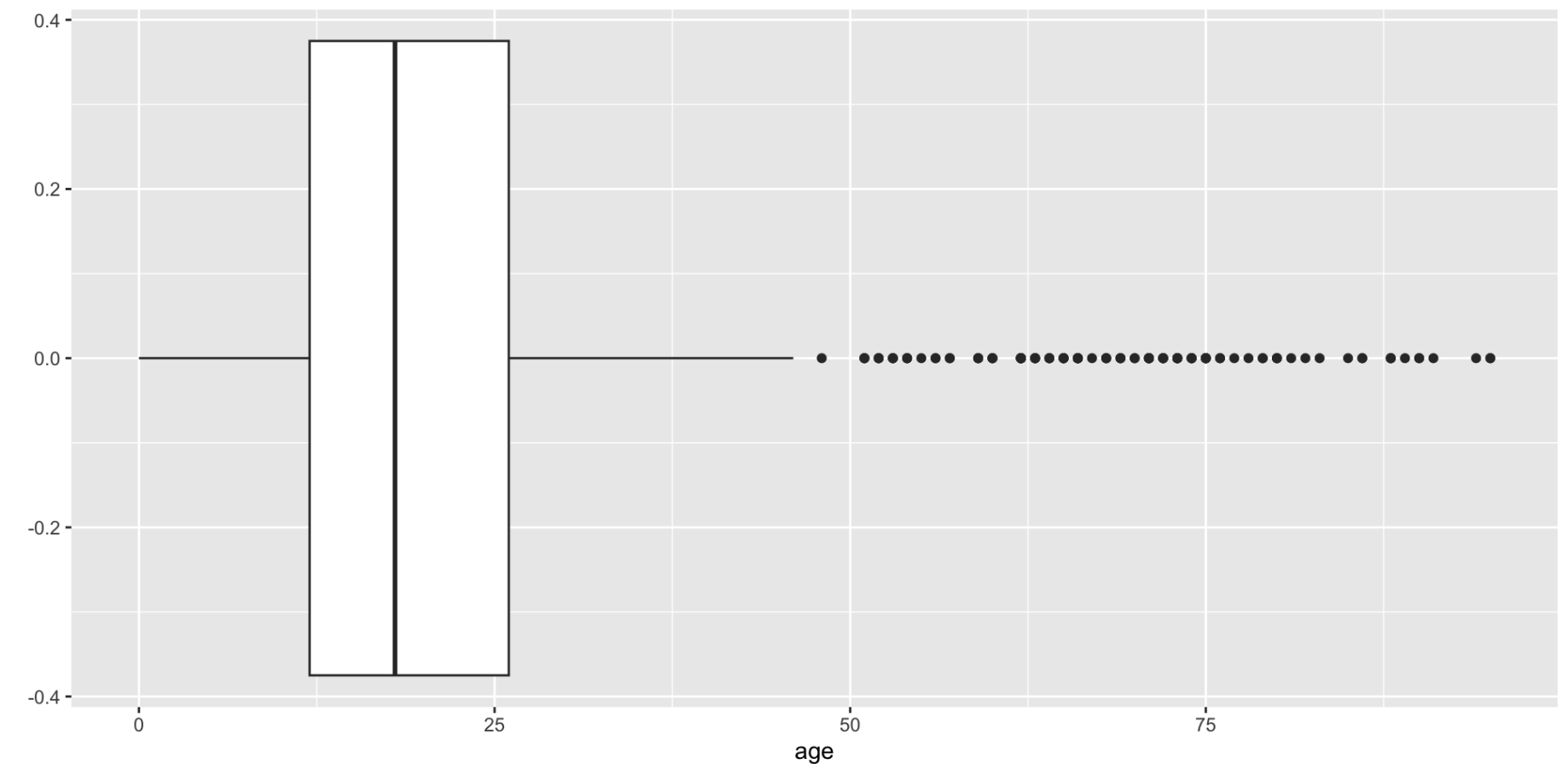
# Other basic plots from **ggplot2**

We can also make a density and boxplot for the continuous variable with ggplot2

```
1  ggplot(data = dds.discr,
2          aes(x = age)) +
3    geom_density()
```

```
1  ggplot(data = dds.discr,
2          aes(x = age)) +
3    geom_boxplot()
```

# Learning Objectives

1. Identify important descriptive statistics and visualize data from a continuous variable

2. Identify important distributions that will be used in 512/612

3. Use our previous tools in 511 to estimate a parameter and construct a confidence interval

4. Use our previous tools in 511 to conduct a hypothesis test

5. Define error rates and power

# Important Distributions

# Distributions that will be used in this class

- Normal distribution

- Chi-square distribution

- t distribution

- F distribution

# Normal Distribution

- Notation: $Y \sim N(\mu, \sigma^2)$

- Arguably, the most important distribution in statistics

- If we know $E(Y) = \mu, \, Var(Y) = \sigma^2$ then

  - 2/3 of $Y$'s distribution lies within 1 $\sigma$ of $\mu$

  - 95% … … is within $\mu \pm 2\sigma$

  - $> 99\%$ … … lies within $\mu \pm 3\sigma$

- Linear combinations of Normal's are Normal
  e.g., $(aY + b) \sim N(a\mu + b, \, a^2\sigma^2)$
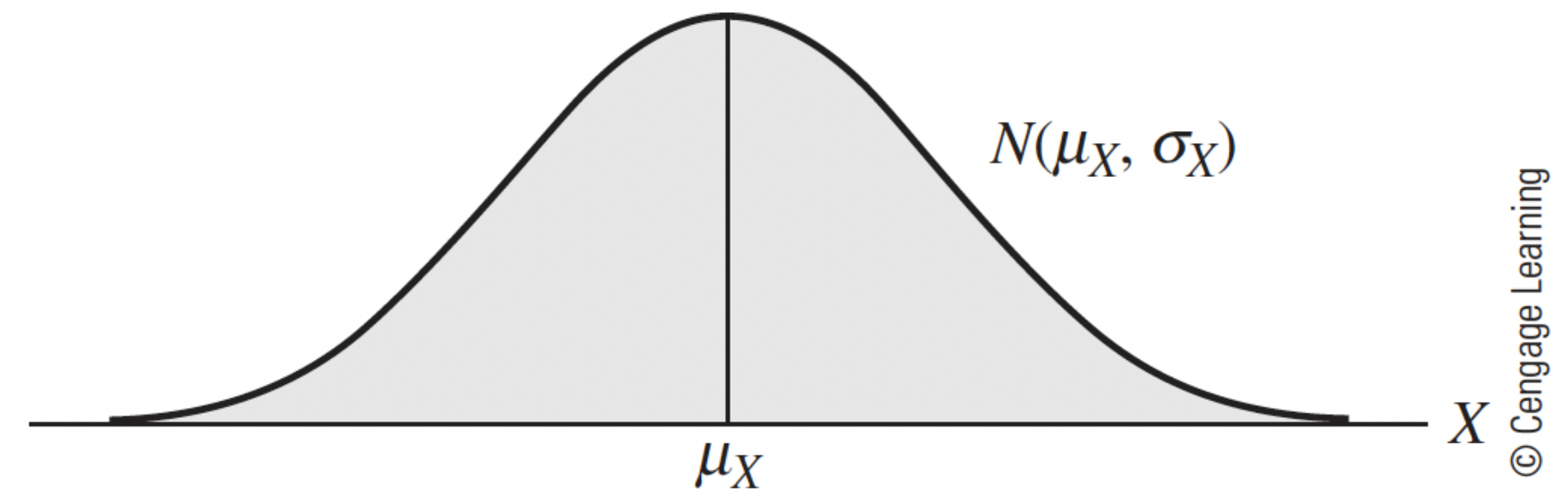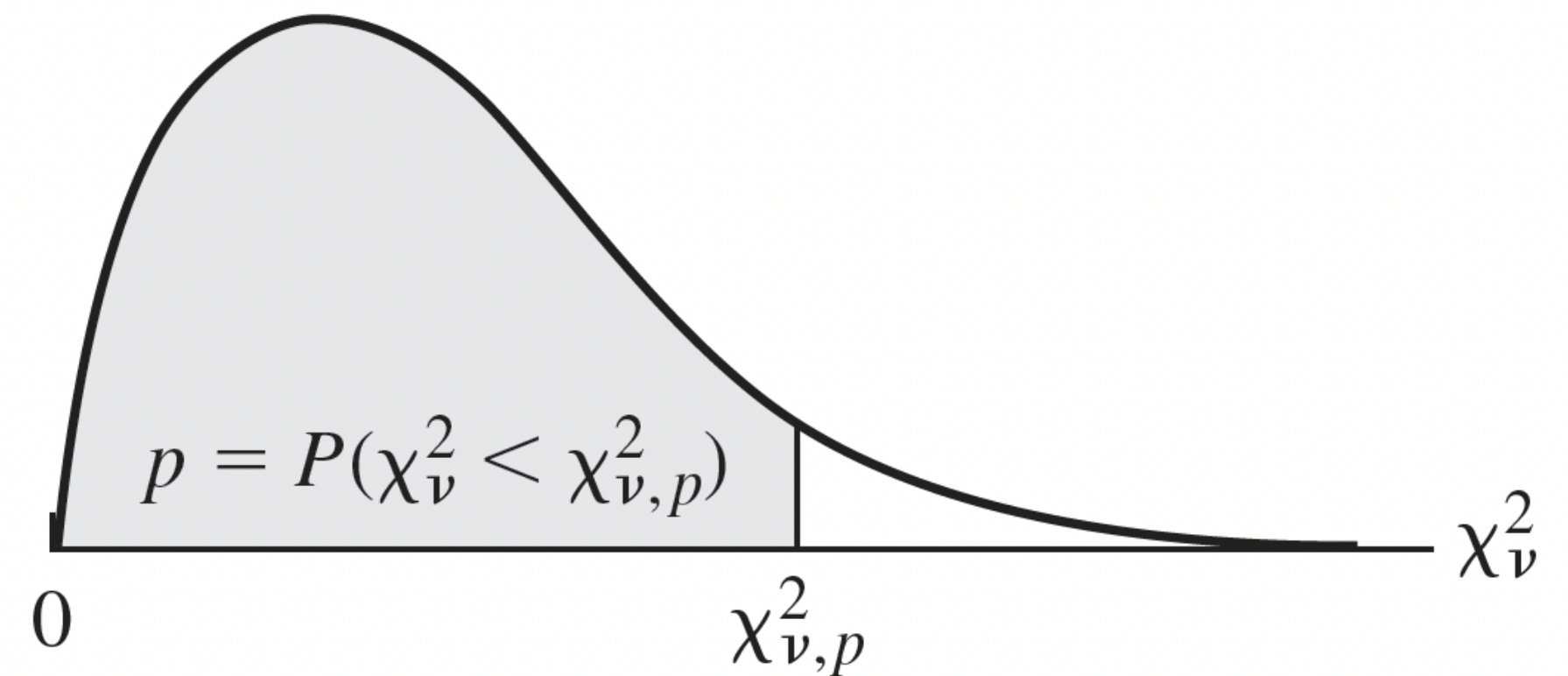
- Standard normal: $Z = \frac{Y - \mu}{\sigma} \sim N(0, 1)$



$N(\mu_X, \sigma_X)$

$\mu_X$

$X$

© Cengage Learning

**FIGURE 3.4** A normal distribution

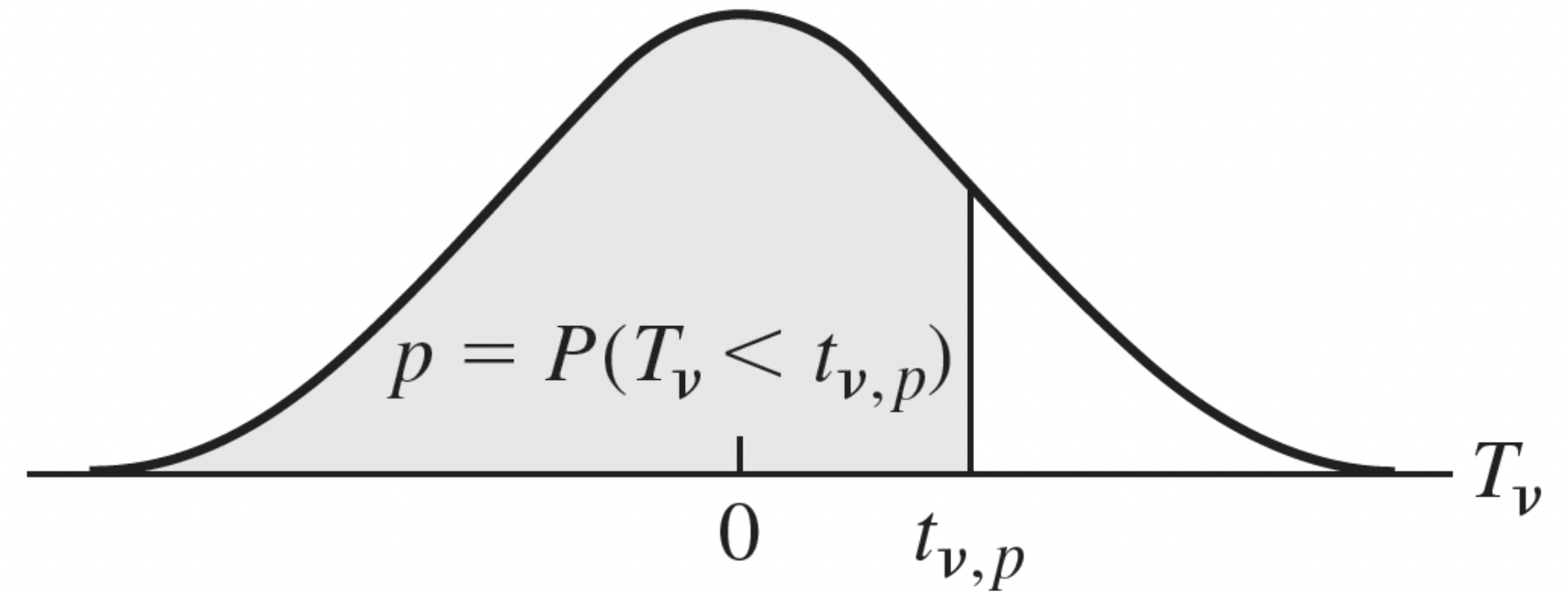# Chi-squared distribution: *models sampling variance*

- Notation: $X \sim \chi^2_{df}$ OR $X \sim \chi^2_{\nu}$

  - Degrees of freedom (df): $df = n - 1$

  - $X$ takes on only positive values

- If $Z_i \sim N(0, 1)$, then $Z_i^2 \sim \chi^2_1$

  - A standard normal distribution squared is the Chi squared distribution with df of 1.

- Used in hypothesis testing and CI's **for variance or standard deviation**

  - Sample variance (and SD) is random and thus can be modeled by a probability distribution: Chi-sqaured

- Chi-squared distribution used to model the **ratio of** the *sample variance* $s^2$ to *population variance* $\sigma^2$:

  - $$\frac{(n-1)s^2}{\sigma^2} \sim \chi^2_{n-1}$$

$$p = P(\chi^2_{\nu} < \chi^2_{\nu, p})$$

$0 \qquad \chi^2_{\nu, p} \qquad \chi^2_{\nu}$

(b) $\chi^2$ distribution

# Student's t Distribution

- Notation: $T \sim t_{df}$ OR $T \sim t_{n-1}$

  - Degrees of freedom (df): $df = n - 1$

  - $T = \dfrac{\bar{x} - \mu_x}{\dfrac{s}{\sqrt{n}}} \sim t_{n-1}$

- In linear modeling, used for inference on individual regression parameters

  - Think: our estimated coefficients ($\hat{\beta}$)



$$p = P(T_v < t_{v,p})$$

(a) Student's *t* distribution

# F-Distribution

- Model ratio of sample variances

  - Ratio of variances is important for hypothesis testing of regression models

- If $X_1^2 \sim \chi_{df1}^2$ and $X_2^2 \sim \chi_{df2}^2$, where $X_1^2 \perp X_2^2$, then:

$$\frac{X_1^2/df\,1}{X_2^2/df\,2} \sim F_{df\,1,df\,2}$$



$$p = P(F_{v_1,v_2} < F_{v_1,v_2,p})$$

(c) $F$ distribution

- only takes on positive values

- Important relationship with t distribution: $T^2 \sim F_{1,v}$

  - The square of a t-distribution with $df = v$

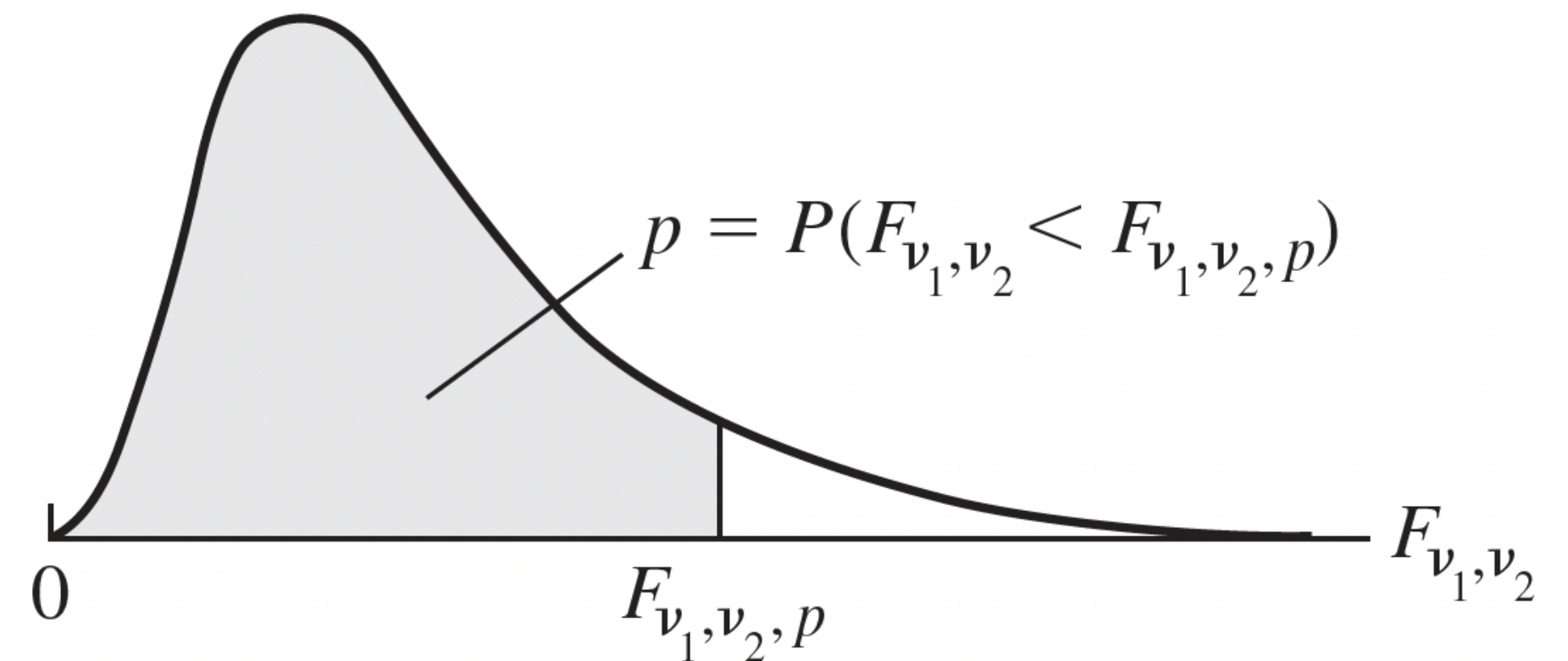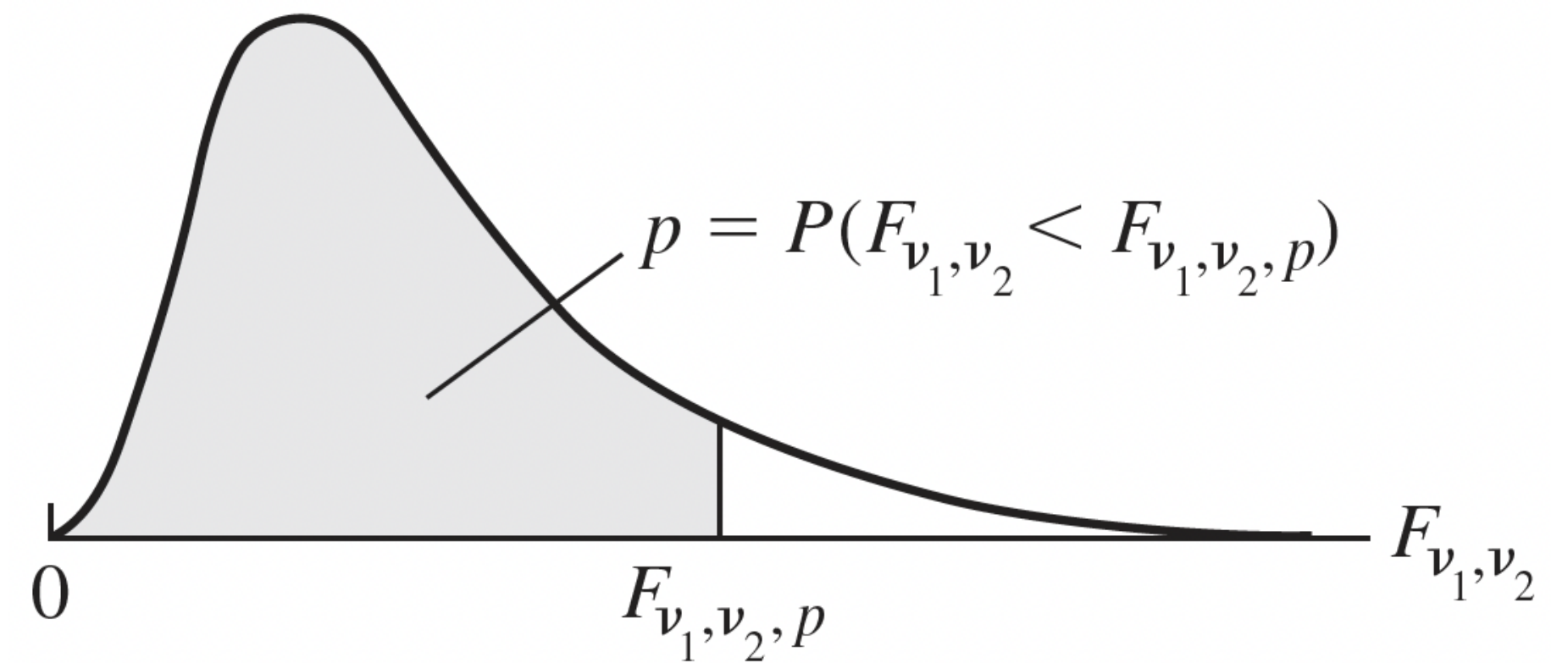  - is an F-distribution with numerator df ($df_1 = 1$) and denominator df ($df_2 = v$)

# F-Distribution

- Model ratio of sample variances
  - Ratio of variances is important for hypothesis testing of regression models
- If $X_1^2 \sim \chi_{df\,1}^2$ and $X_2^2 \sim \chi_{df\,2}^2$, where $X_1^2 \perp X_2^2$, then:

$$\frac{X_1^2/df\,1}{X_2^2/df\,2} \sim F_{df\,1,df\,2}$$

- only takes on positive values

- Important relationship with t distribution: $T^2 \sim F_{1,\nu}$
  - The square of a t-distribution with $df = \nu$
  - is an F-distribution with numerator df $(df_1 = 1)$ and denominator df $(df_2 = \nu)$



$$p = P(F_{v_1,v_2} < F_{v_1,v_2,p})$$

$0$      $F_{v_1,v_2,p}$      $F_{v_1,v_2}$

(c) $F$ distribution

Is there a relationship between our chi-squared and F-distribution?

Recall, $\dfrac{(n-1)s^2}{\sigma^2} \sim \chi_{n-1}^2 \, .$

The F-distribution for a ratio of variances between two models is: $F = \dfrac{s_1^2 \sigma_2^2}{s_2^2 \sigma_1^2} \sim F_{n_1-1,n_2-1}$

# R code for probability distributions

Here is a site with the various probability distributions and their R code.

- It also includes practice with R code to see what each function outputs

| Distribution | Functions | | | |
|---|---|---|---|---|
| Beta | pbeta | qbeta | dbeta | rbeta |
| Binomial (including Bernoulli) | pbinom | qbinom | dbinom | rbinom |
| Birthday | pbirthday | qbirthday | | |
| Cauchy | pcauchy | qcauchy | dcauchy | rcauchy |
| Chi-Square | pchisq | qchisq | dchisq | rchisq |
| Discrete Uniform | sample | | | |
| Exponential | pexp | qexp | dexp | rexp |
| F | pf | qf | df | rf |
| Gamma | pgamma | qgamma | dgamma | rgamma |
| Geometric | pgeom | qgeom | dgeom | rgeom |
| Hypergeometric | phyper | qhyper | dhyper | rhyper |
| Logistic | plogis | qlogis | dlogis | rlogis |
| Log Normal | plnorm | qlnorm | dlnorm | rlnorm |
| Multinomial | | | dmultinom | rmultinom |
| Negative Binomial | pnbinom | qnbinom | dnbinom | rnbinom |
| Normal | pnorm | qnorm | dnorm | rnorm |
| Poisson | ppois | qpois | dpois | rpois |
| Kolmogorov-Smirnov Test Statistic | psmirnov | qsmirnov | | rsmirnov |
| Student t | pt | qt | dt | rt |
| Studentized Range | ptukey | qtukey | dtukey | rtukey |
| Continuous Uniform | punif | qunif | dunif | runif |
| Weibull | pweibull | qweibull | dweibull | rweibull |
| Wilcoxon Rank Sum Statistic | pwilcox | qwilcox | dwilcox | rwilcox |
| Wilcoxon Signed Rank Statistic | psignrank | qsignrank | dsignrank | rsignrank |
| Wishart | | | | rWishart |

# Learning Objectives

1. Identify important descriptive statistics and visualize data from a continuous variable

2. Identify important distributions that will be used in 512/612

3. Use our previous tools in 511 to estimate a parameter and construct a confidence interval

4. Use our previous tools in 511 to conduct a hypothesis test

5. Define error rates and power

# Statistical inference: Estimation

# Confidence interval for one mean

The confidence interval for population mean μ:

$$\bar{x} \pm t^* \frac{s}{\sqrt{n}}$$

- where $t^*$ is the critical value for the 95% (or other percent) corresponding to the t-distribution and dependent on $df = n - 1$

We can use R to find the critical t-value, $t^*$

For example the critical value for the 95% CI with $n = 10$ subjects is...

```
1  qt(0.975, df=9)
```
```
[1] 2.262157
```

- Recall, that as the $df$ increases, the t-distribution converges towards the Normal distribution

# Confidence interval for one mean

The confidence interval for population mean μ:

$$\bar{x} \pm t^* \frac{s}{\sqrt{n}}$$

- where $t^*$ is the critical value for the 95% (or other percent) corresponding to the t-distribution and dependent on $df = n - 1$

> ## We can use R to find the critical t-value, $t^*$
>
> For example the critical value for the 95% CI with $n = 10$ subjects is...
>
> ```
> 1  qt(0.975, df=9)
> ```
>
> `[1] 2.262157`
>
> - Recall, that as the $df$ increases, the t-distribution converges towards the Normal distribution

We can also use `t.test` in R to calculate the confidence interval if we have a dataset.

```
1  t.test(dds.discr$age)
```

```
        One Sample t-test

data:  dds.discr$age
t = 39.053, df = 999, p-value < 2.2e-16
alternative hypothesis: true mean is not equal
to 0
95 percent confidence interval:
 21.65434 23.94566
sample estimates:
```

# Confidence interval for two independent means

The confidence interval for difference in independent population means, $\mu_1$ and $\mu_2$ :

$$\bar{x}_1 - \bar{x}_2 \pm t^* \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

- where $t^*$ is the critical value for the 95% (or other percent) corresponding to the t-distribution and dependent on $df = n_1 + n_2 - 2$

# Here's a decent source for other R code for tests in 511

Website from UCLA

# Learning Objectives

1. Identify important descriptive statistics and visualize data from a continuous variable

2. Identify important distributions that will be used in 512/612

3. Use our previous tools in 511 to estimate a parameter and construct a confidence interval

4. Use our previous tools in 511 to conduct a hypothesis test

5. Define error rates and power

# Statistical inference: Hypothesis testing

# Steps in hypothesis testing

1. Check the assumptions regarding the properties of the underlying variable(s) being measured that are needed to justify use of the testing procedure under consideration.
2. State the null hypothesis $H_0$ and the alternative hypothesis $H_A$.
3. Specify the significance level $\alpha$.
4. Specify the test statistic to be used and its distribution under $H_0$.

## Critical region method

5. Form the decision rule for rejecting or not rejecting $H_0$ (i.e., specify the rejection and nonrejection regions for the test, based on both $H_A$ and $\alpha$).
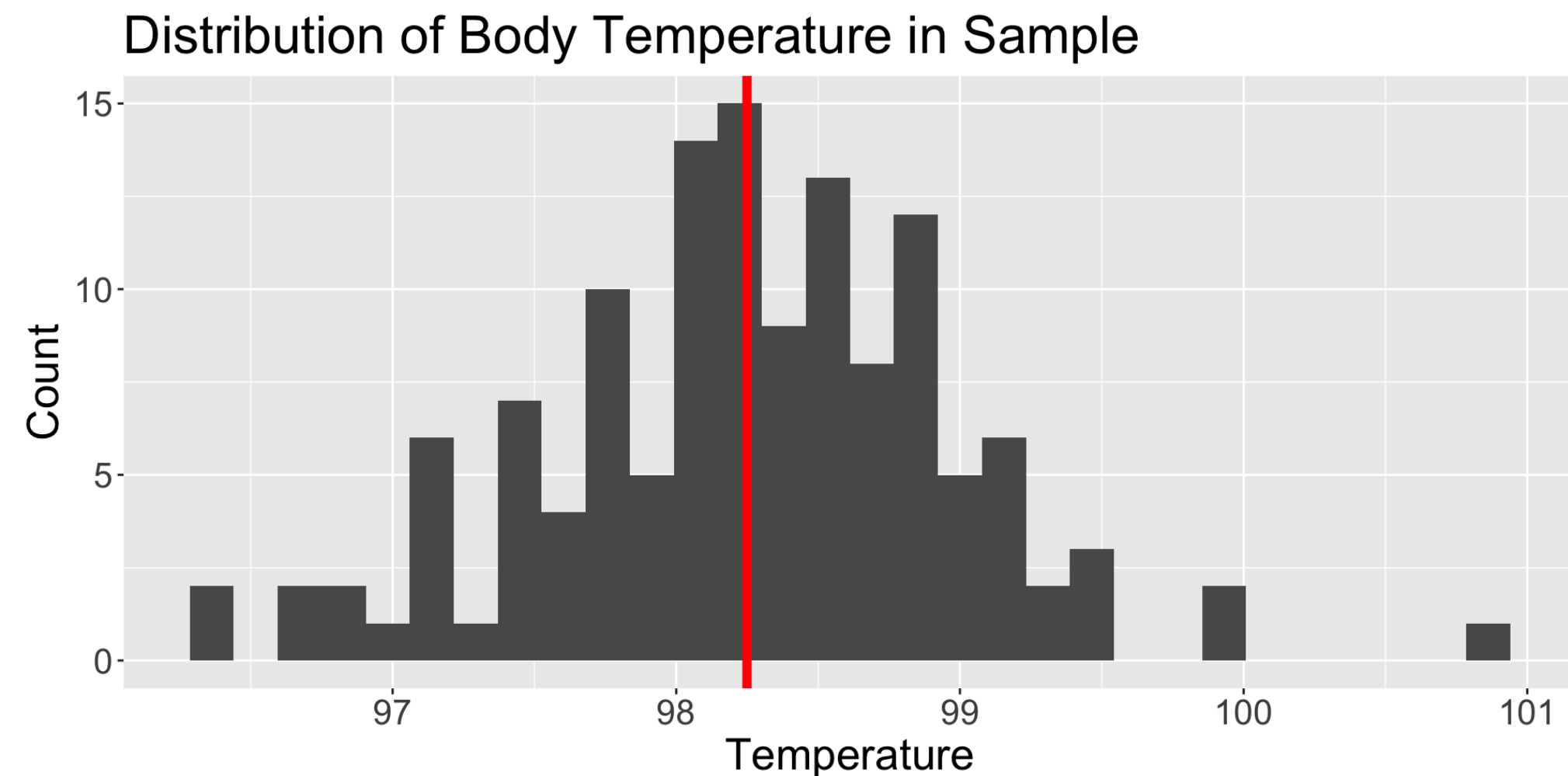6. Compute the value of the test statistic from the observed data.

## p-value method

5. Compute the value of the test statistic from the observed data.
6. Calculate the p-value

7. Draw conclusions regarding rejection or nonrejection of H0.

# Example: one sample t-test

```r
1   BodyTemps = read.csv("data/BodyTemperatures.csv")
2
3   ggplot(data = BodyTemps,
4          aes(x = Temperature)) +
5     geom_histogram() +
6     theme(text = element_text(size=20)) +
7     labs(x = "Temperature", y = "Count",
8          title = "Distribution of Body Temperature in Sample") +
9     geom_vline(aes(xintercept = mean(BodyTemps$Temperature, na.rm = T)),
10               color = "red", linewidth = 2)
```



Distribution of Body Temperature in Sample

# Example: one sample t-test using *p-value approach*

We want to see what the mean population body temperature is.

2. State the null and alternative hypotheses:

| $H_0 : \mu = 98.6$ | $H_0$: The population mean body temperature is 98.6 degrees F |
|---|---|
| $H_A : \mu \neq 98.6$ | $H_A$: The population mean body temperature is **not** 98.6 degrees F |

3. The significance level is $\alpha = 0.05$

4. The test statistic, $t_{\bar{x}}$ follows a student's t-distribution with $df = n - 1 = 129$

5. The test statistic is: $t_{\bar{x}} = \dfrac{\bar{x} - \mu_0}{\dfrac{s}{\sqrt{n}}}$ and with the data: $t_{\bar{x}} = \dfrac{98.25 - 98.6}{\dfrac{0.73}{\sqrt{130}}} = -5.45$

6. Calculate the p-value: $p - value = P(t \leq -5.45) + P(t \geq 5.45)$

```
1  2*pt(-5.4548, df = 130-1, lower.tail=T)
```
```
[1] 2.410889e-07
```

7. Conclusion: We reject the null hypothesis. There is sufficient evidence that the (population) mean body temperature after is different from 98.6 degree ($p - value < 0.001$).

# Example: one sample t-test using *critical values approach*

We want to see what the mean population body temperature is.

2. State the null and alternative hypotheses:

$$H_0 : \mu = 98.6 \qquad\qquad H_0: \text{The population mean body temperature is 98.6 degrees F}$$

$$H_A : \mu \neq 98.6 \qquad\qquad H_A: \text{The population mean body temperature is \textbf{not} 98.6 degrees F}$$

3. The significance level is $\alpha = 0.05$

4. The test statistic, $t_{\bar{x}}$ follows a student's t-distribution with $df = n - 1 = 129$

5. Decision rule (critical value): For $\alpha = 0.05$, $2 * P(t \geq t^*) = 0.05$

```
1 qt(0.05/2, df = 130-1, lower.tail=F)
```
[1] 1.978524

6. The test statistic is: $t_{\bar{x}} = \dfrac{\bar{x} - \mu_0}{\dfrac{s}{\sqrt{n}}}$ and with the data: $t_{\bar{x}} = \dfrac{98.25 - 98.6}{\dfrac{0.73}{\sqrt{130}}} = -5.45$

7. Conclusion: We reject the null hypothesis. There is sufficient evidence that the (population) mean body temperature after is different from 98.6 degree ( 95% CI: 98.12, 98.38).

# How did we get the 95% CI?

- The `t.test` function can help us answer this, and give us the needed information for both approaches.

```r
1  BodyTemps = read.csv("data/BodyTemperatures.csv")
2
3  t.test(x = BodyTemps$Temperature,
4         # alternative = "two-sided",
5         mu = 98.6)
```

```
    One Sample t-test

data:  BodyTemps$Temperature
t = -5.4548, df = 129, p-value = 2.411e-07
alternative hypothesis: true mean is not equal to 98.6
95 percent confidence interval:
 98.12200 98.37646
sample estimates:
mean of x
```

# Learning Objectives

1. Identify important descriptive statistics and visualize data from a continuous variable

2. Identify important distributions that will be used in 512/612

3. Use our previous tools in 511 to estimate a parameter and construct a confidence interval

4. Use our previous tools in 511 to conduct a hypothesis test

5. Define error rates and power

# Error Rates and Power

# Outcomes of our hypothesis test

| TABLE 3.1 | Outcomes of hypothesis testing | |
|---|---|---|
| **Hypothesis Chosen** | **True State of Nature** | |
| | $H_0$ | $H_A$ |
| $H_0$ | Correct decision | False negative decision (Type II error) |
| $H_A$ | False positive decision (Type I error) | Correct decision |

© Cengage Learning

# Prabilities of outcomes

- Type 1 error is $\alpha$
  - The probability that we falsly reject the null hypothesis (but the null is true!!)
- Power is $1 - \beta$
  - The probability of correctly rejecting the null hypothesis

| TABLE 3.2 Probabilities of outcomes of hypothesis testing | | |
|---|---|---|
| **Hypothesis Chosen** | **True State of Nature** | |
| | $H_0$ | $H_A$ |
| $H_0$ | $1 - \alpha$ | $\beta$ |
| $H_A$ | $\alpha$ | $1 - \beta$ |

© Cengage Learning

# What I think is the most intuitive way to look at it



Null hypothesis (H₀) distribution

Alternative hypothesis (H₁) distribution

$1 - \alpha$

$1 - \beta$

$\beta$   $\alpha$

Type II error rate   Type I error rate