

Econometrics 1 *Applied Econometrics with R*

Supplement: Review of Statistics

黄嘉平

中国经济特区研究中心 讲师

办公室：文科楼2613

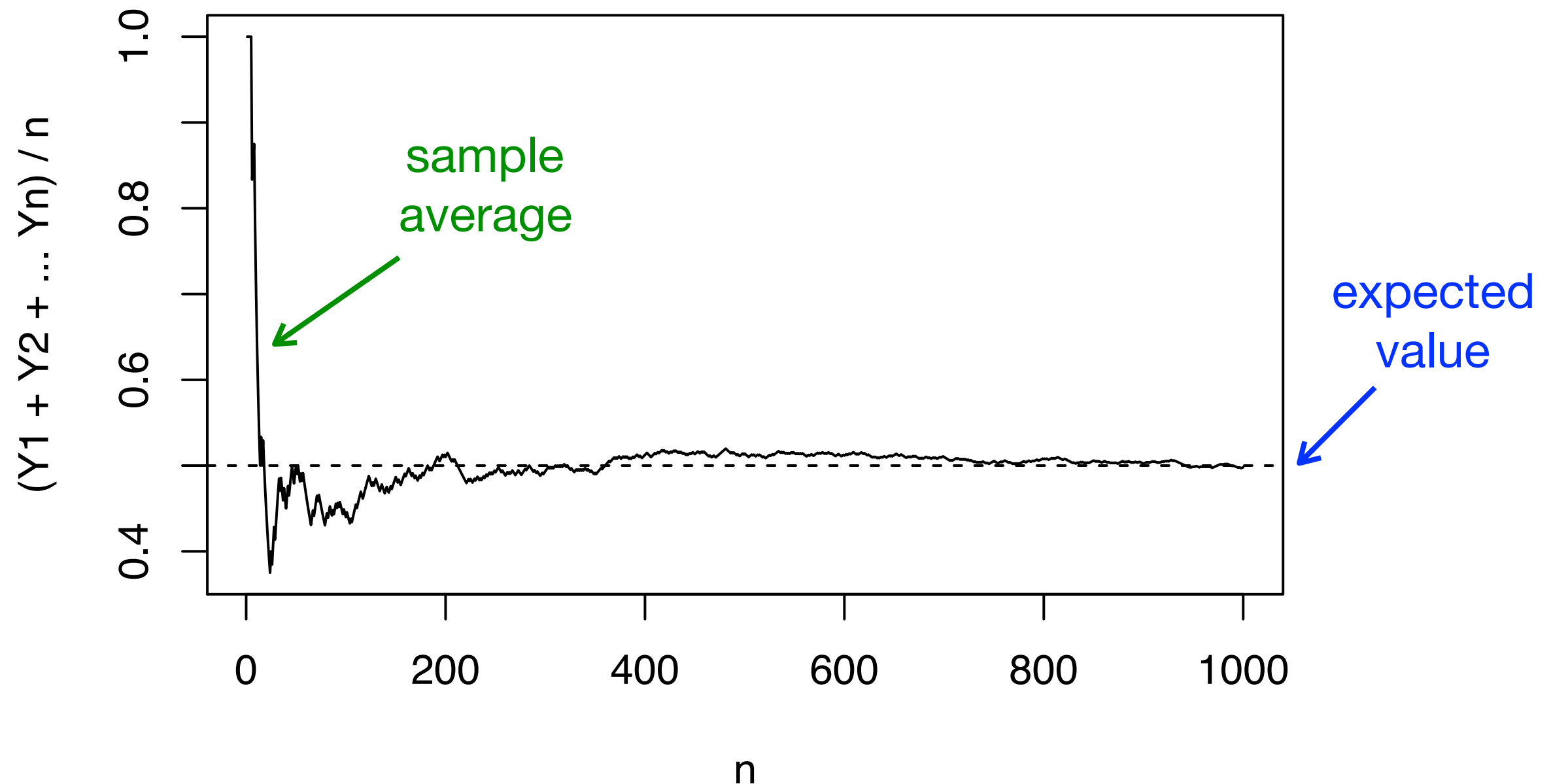
E-mail: huangjp@szu.edu.cn

Tel: (0755) 2695 0548

Website: <https://huangjp.com>

Demonstrating the LLN

- The sample mean of n Bernoulli random variables (flipping a fair coin n times).



Basics of Statistics

What does statistics do?

- A typical statistical question:

What is the mean of the distribution of earnings of recent college graduates?

- It is too expensive, sometimes even impossible, to do an exhaustive survey to know the answer of such questions.
- Though we cannot know the exact answer, we can use statistical methods to reach tentative conclusions — to draw statistical inferences — about characteristics of the full population based on simple *random sampling* data.

What does statistics do?

- Statistics (statistical tools) helps us answer questions about unknown characteristics of distributions in *population* of interest.
- Three most used statistical methods:

Estimation

Hypothesis testing

Confidence intervals

Estimation

Estimation: a “best guess” numerical value

- Suppose you want to know the mean value of Y (that is, μ_Y) in a population.
- You have a sample of n i.i.d. observations Y_1, \dots, Y_n .
- A natural way of estimate the value of μ_Y is to calculate the sample mean \bar{Y} , which is an *estimator* of μ_Y .
- There are many other possible estimators, for example, the first observation Y_1 .

Estimator and estimates

- An **estimator** is a *function* of a sample of data to be drawn randomly from a population. Thus it is a random variable.
- An **estimate** is the numerical *value* of the estimator when it is actually computed using data from a specific sample. Thus it is a nonrandom number.

Desirable characteristics of an estimator $\hat{\mu}_Y$

- *Unbiasedness* — The bias of $\hat{\mu}_Y$ is $E(\hat{\mu}_Y) - \mu_Y$. We say $\hat{\mu}_Y$ is an unbiased estimator of μ_Y if $E(\hat{\mu}_Y) = \mu_Y$.
- *Consistency* — $\hat{\mu}_Y$ is a consistent estimator of μ_Y if

$$\hat{\mu}_Y \xrightarrow{p} \mu_Y$$

as the sample size increases.

- *Efficiency* — Let $\tilde{\mu}_Y$ be another estimator of μ_Y . Suppose both $\tilde{\mu}_Y$ and $\hat{\mu}_Y$ are unbiased. $\hat{\mu}_Y$ is more efficient than $\tilde{\mu}_Y$ if $\text{var}(\hat{\mu}_Y) < \text{var}(\tilde{\mu}_Y)$.

Properties of \bar{Y}

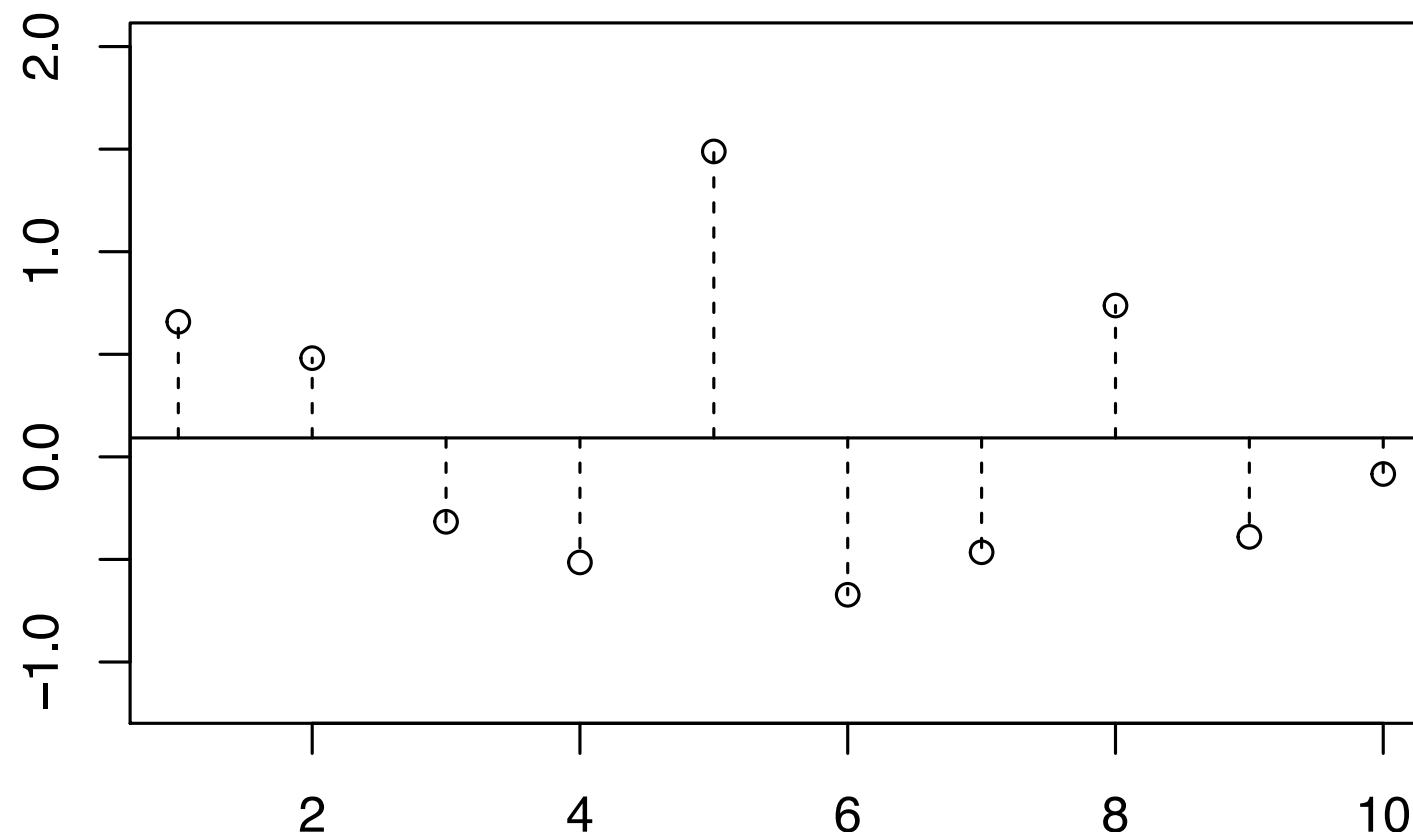
- \bar{Y} is unbiased.
- \bar{Y} is consistent (the law of large numbers).
- \bar{Y} is the **Best Linear Unbiased Estimator (BLUE)** of μ_Y . That is, it is the most efficient (best) estimator among all estimators that are unbiased and are linear functions of Y_1, \dots, Y_n .
- \bar{Y} is the least squares estimator of μ_Y .

Least squares estimator

- The estimator m that minimizes function

$$\sum_{i=1}^n (Y_i - m)^2$$

is called the least squares estimator.



Practice

- Write a program to find the least squares estimator of the mean of a standard normal distributed population using 100 random samples, and compare it with \bar{Y} .

Hint:

1. Generate 100 standard normal random samples.
2. Make a guess of an interval that the population mean may drop in.
3. Construct a grid on that interval.
4. Each value of the grid can be seen as a candidate of the least squares estimator.

```
y <- rnorm(100)
ybar <- mean(y)

mstep <- 0.0001
mgrid <- seq(-5, 5, mstep)
          # candidates of LSE
msls <- rep(0, length(mgrid))
          # sum of linear squares

for (i in 1:length(mgrid)) {
  msls[i] <- sum((y - mgrid[i])^2)
}

mindex <- which.min(msls)
          # the index of the minimum in "msls"
lse <- mgrid[mindex]
          # the least squares estimator
```

Hypothesis Testing

Hypothesis

- A hypothesis here is a question about the characteristics of the population that can be answered by “yes” or “no”.
- **Null hypothesis** — the specific hypothesis to be tested.

H_0 : Is the population mean equal to 0?

- **Alternative hypothesis** — the hypothesis which is true if the null hypothesis is not.

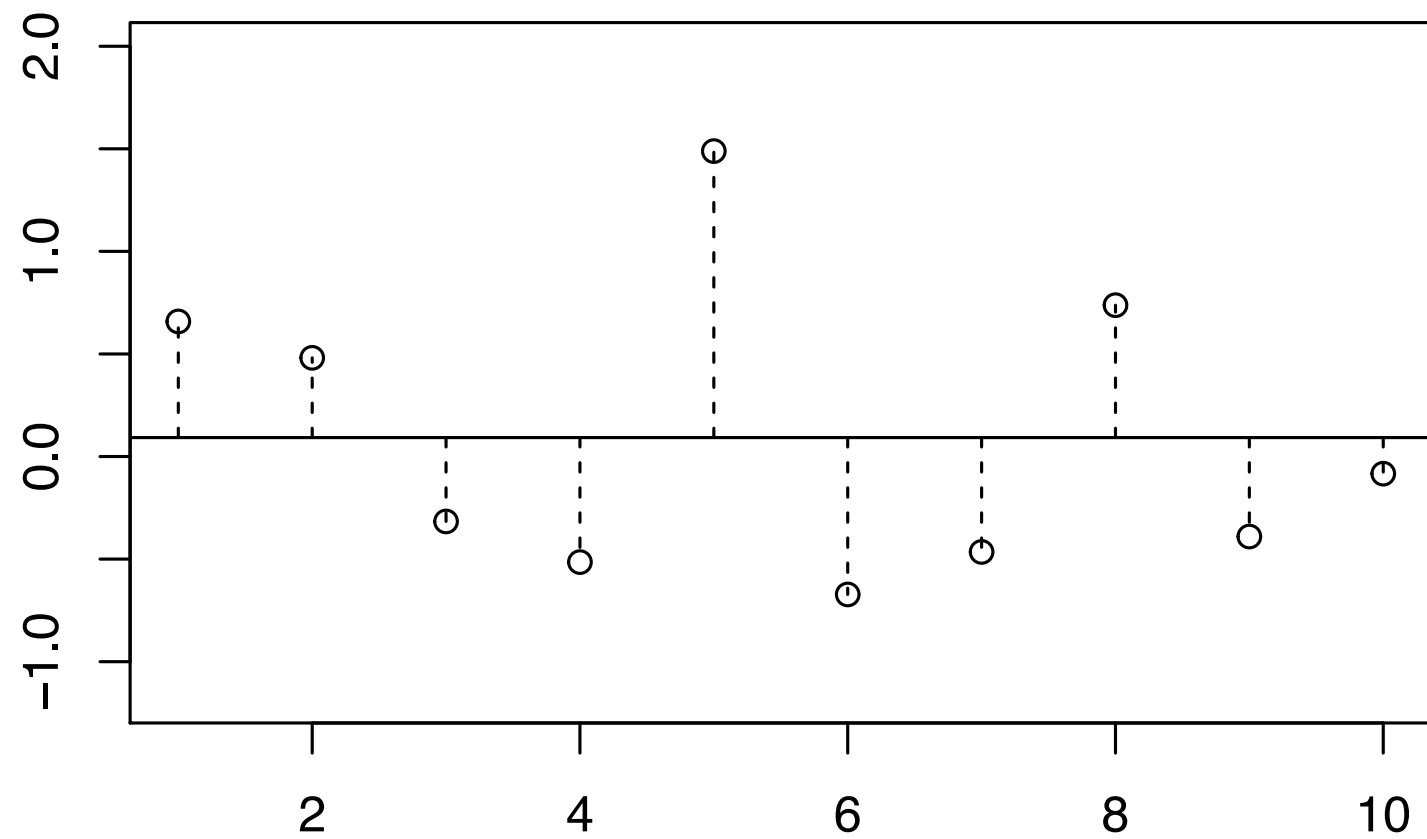
H_1 : Is the population mean not equal to 0?

Hypothesis testing

- To either *reject* or *not reject* the null hypothesis based on a random sample of data.
- The null hypothesis is *rejected* if the alternative hypothesis is likely to be true.
- If the alternative hypothesis is not true, it means that the null hypothesis is failed to be rejected by the current sample.
 - Either the null hypothesis is true,
or the evidence is not strong enough to reject it.

Example

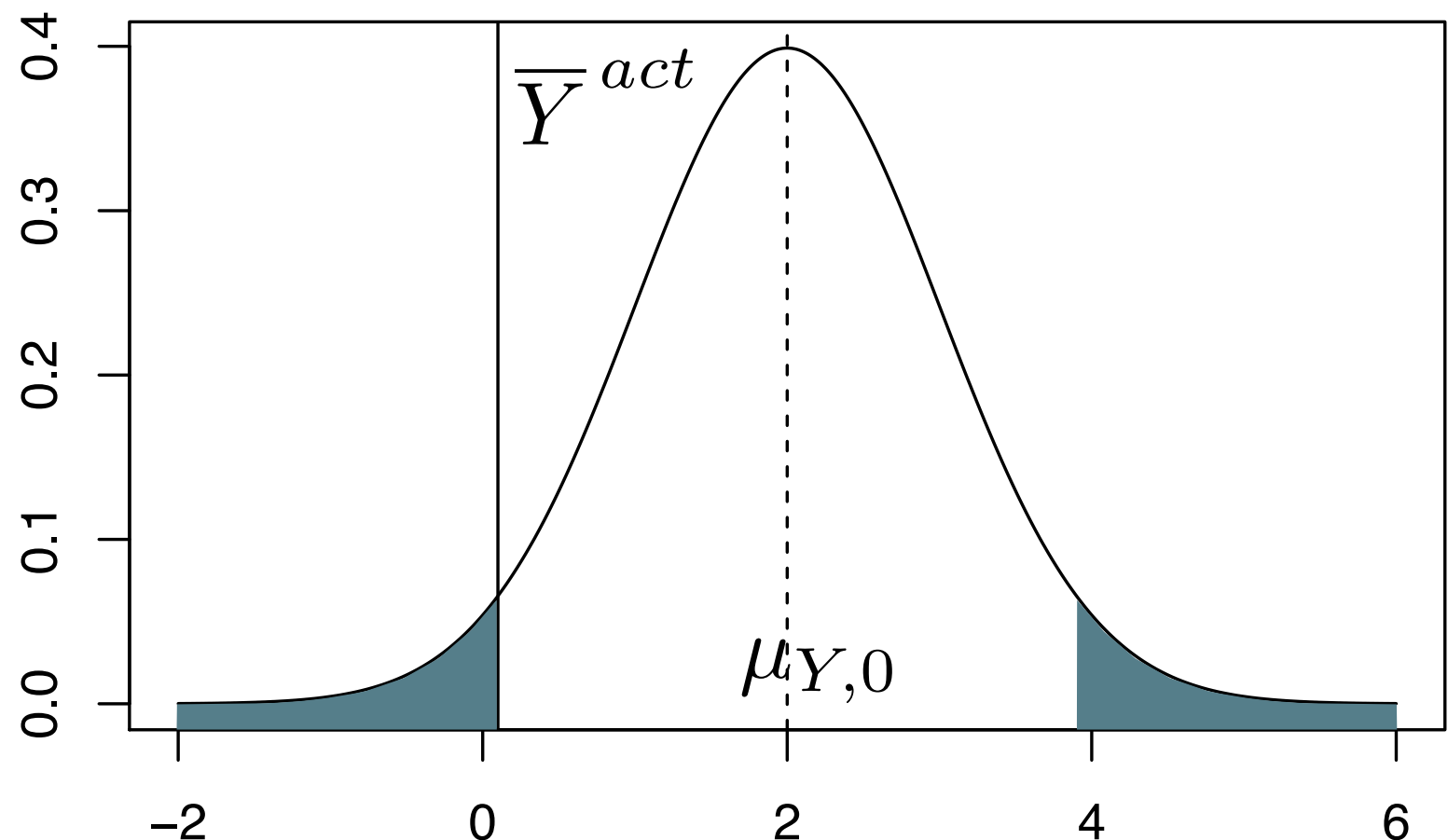
- $H_0: E(Y) = 2$
- $H_1: E(Y) \neq 2$



The p -value

- The p -value, also called the *significance probability*, is the probability of drawing a statistic at least as adverse to the null hypothesis as the one you actually computed in your sample, assuming the null hypothesis is correct.

$$H_0 : E(Y) = \mu_{Y,0}$$



The p -value

- If the p -value is extremely small, it means that our observation is very unlikely to be obtained from the population, assuming the null hypothesis is correct. Thus, the null hypothesis is very likely to be wrong.
- Let \bar{Y}^{act} denote the actual computed sample average of the data set at hand. $H_0 : E(Y) = \mu_{Y,0}$

$$p\text{-value} = \Pr_{H_0} \left[|\bar{Y} - \mu_{Y,0}| > |\bar{Y}^{act} - \mu_{Y,0}| \right]$$

- How to evaluate the p -value depends on whether the population variance is known.

Calculation the p -value when σ_Y is known

- Under the central limit theorem and the null hypothesis,

$$\bar{Y} \sim N(\mu_{Y,0}, \sigma_Y^2/n)$$

- So the standardized r.v. $\frac{\bar{Y} - \mu_{Y,0}}{\sigma_Y/\sqrt{n}}$ has a standard normal distribution.

$$\begin{aligned} p\text{-value} &= \Pr_{H_0} \left[\left| \frac{\bar{Y} - \mu_{Y,0}}{\sigma_Y/\sqrt{n}} \right| > \left| \frac{\bar{Y}^{act} - \mu_{Y,0}}{\sigma_Y/\sqrt{n}} \right| \right] \\ &= 2\Phi \left(- \left| \frac{\bar{Y}^{act} - \mu_{Y,0}}{\sigma_Y/\sqrt{n}} \right| \right) \end{aligned}$$

When σ_Y is unknown

- If σ_Y is unknown, we need to estimate it before calculation the p -value.
- Sample variance — an estimator of the population variance.
- Sample standard deviation — an estimator of the population standard deviation.
- Standard error of the sample mean — an estimator of of the standard deviation of the sampling distribution of the sample mean.

- The sample variance s_Y^2

$$s_Y^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2$$

- The sample standard deviation is s_Y
- The standard error of \bar{Y}

$$SE(\bar{Y}) = s_Y / \sqrt{n}$$

Calculation the p -value when σ_Y is unknown

- We can use s_Y instead of σ_Y

$$p\text{-value} = 2\Phi\left(-\left|\frac{\bar{Y}^{act} - \mu_{Y,0}}{s_Y / \sqrt{n}}\right|\right) = 2\Phi\left(-\left|\frac{\bar{Y}^{act} - \mu_{Y,0}}{SE(\bar{Y})}\right|\right)$$

- The t -statistic

$$t = \frac{\bar{Y} - \mu_{Y,0}}{SE(\bar{Y})}$$

Large-sample distribution of the t -statistic

- By the central limit theorem, when the sample size n is large, the distribution of t is well approximated by the standard normal distribution $N(0, 1)$.
- The p -value can be rewritten in terms of the t -statistic actually computed:

$$t^{act} = \frac{\bar{Y}^{act} - \mu_{Y,0}}{SE(\bar{Y})}$$

$$p\text{-value} = 2\Phi(-|t^{act}|)$$

General procedure of hypothesis testing for the population mean

1. Specify a null hypothesis and an alternative hypothesis

$$H_0 : E(Y) = \mu_{Y,0} \quad H_1 : E(Y) \neq \mu_{Y,0}$$

2. Calculate the t -statistic and the p -value
3. Choose a level of significance — e.g. 5%
4. If the p -value < 0.05 , the null hypothesis is rejected.

The corresponding critical value to 0.05 is 1.96 under the standard normal distribution, we can also say the null hypothesis is rejected if $|t^{act}| > 1.96$.

Practice

- The R command for this kind of hypothesis testing is `t.test(...)`
- Read the help of `t.test` and learn how to use it.
- Generate a sample of 100 observations that follows a chi-squared distribution with 3 degrees of freedom.
- Test the hypotheses

$$H_0 : E(Y) = \mu_{Y,0} \qquad H_1 : E(Y) \neq \mu_{Y,0}$$

where $\mu_{Y,0} = 1.5$. Find the estimate, the t -statistic, and the p -value. What is your conclusion?

```
> y <- rchisq(100, 3)
> mu0 <- 1.5
> ty <- t.test(y, mu = mu0)
      # with null hypothesis "mu = 0"

> ty
```

One Sample t-test

```
data: y
t = 4.5992, df = 99, p-value = 1.256e-05
alternative hypothesis: true mean is not equal to 1.5
95 percent confidence interval:
 1.975775 2.697797
sample estimates:
mean of x
 2.336786
```

The `t.test` command use an exact Student t distribution instead of normal approximation. We address this point later.

Confidence Intervals

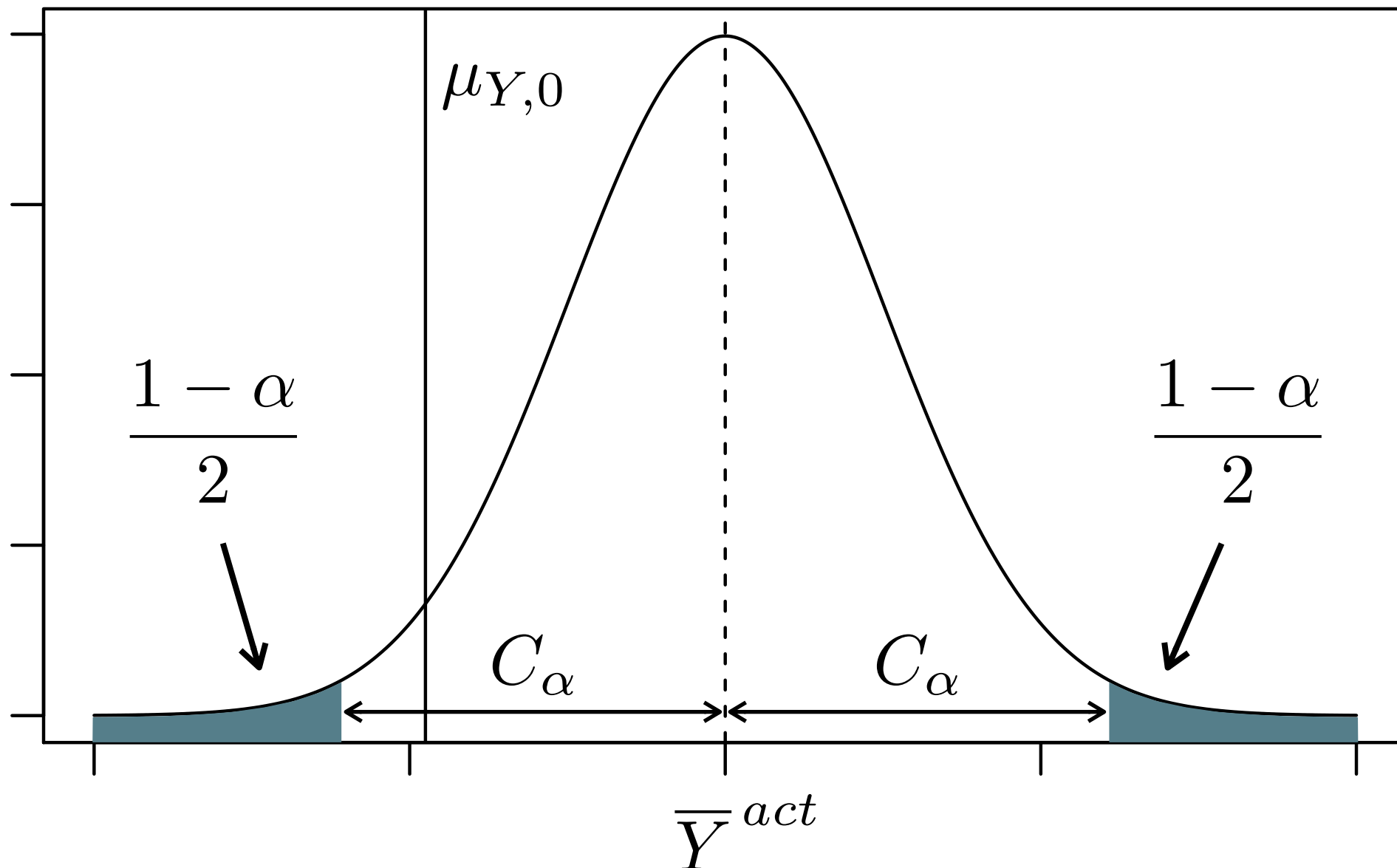
Confidence intervals

- What can we say about the population mean without specifying a hypothesis?
- It is possible to *use data from a random sample* to construct a set of values that contains the true population mean μ_Y with a certain pre-specified probability. This set is called a *confidence set*, and the pre-specified probability is called the *confidence level*.
- The confidence set for μ_Y is an interval, thus it is also called a *confidence interval*.

Confidence interval v.s. p -value

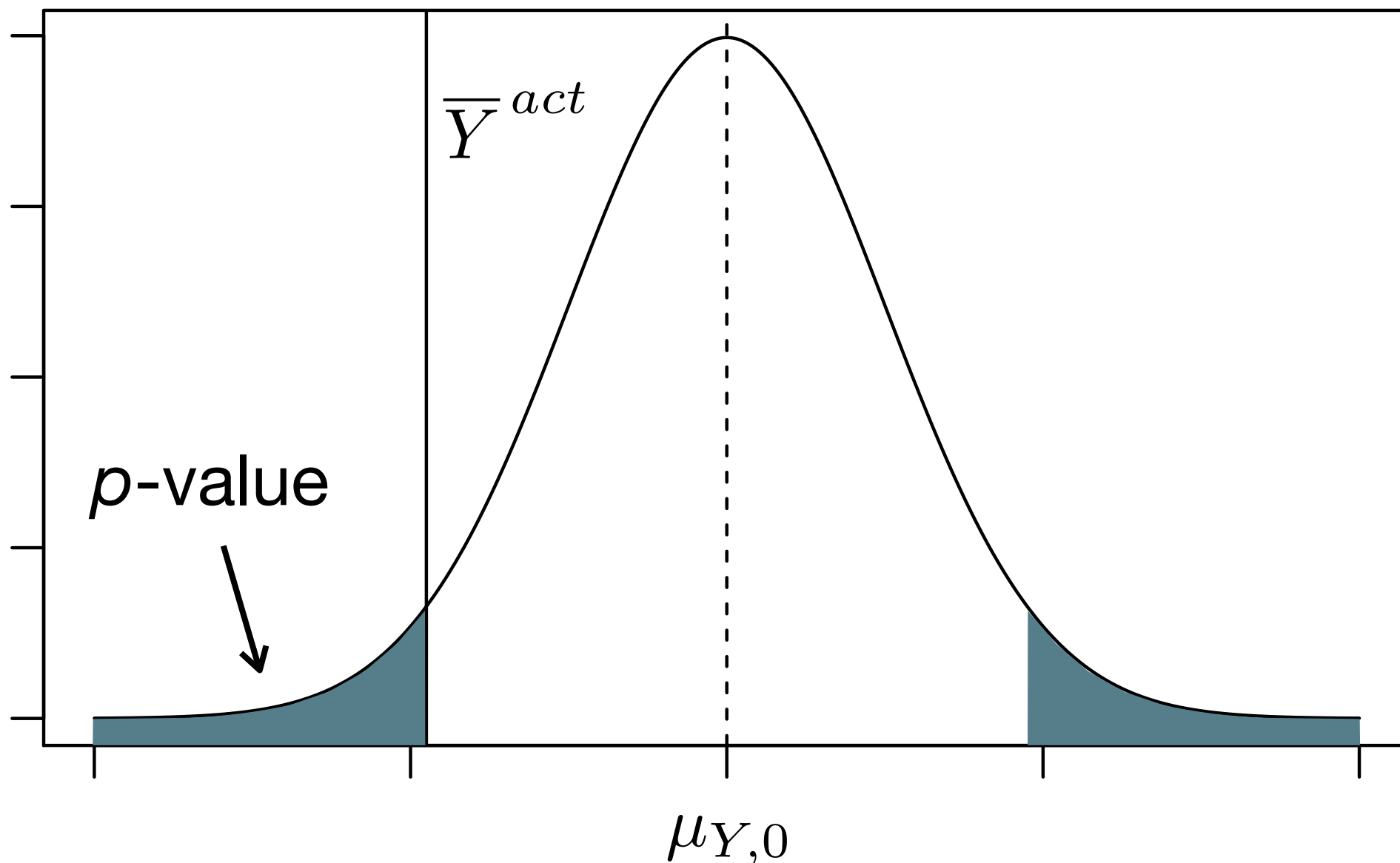
α confidence interval of μ_Y is

$$[\bar{Y} - C_\alpha \times SE(\bar{Y}), \bar{Y} + C_\alpha \times SE(\bar{Y})]$$



Confidence interval v.s. p -value

$$p\text{-value} = \Pr_{H_0} \left[|\bar{Y} - \mu_{Y,0}| > |\bar{Y}^{act} - \mu_{Y,0}| \right]$$



Confidence intervals when sample size is large

- α confidence interval of μ_Y

$$[\bar{Y} - C_\alpha \times SE(\bar{Y}), \bar{Y} + C_\alpha \times SE(\bar{Y})]$$

- If the sample size is large, by the central limit theorem, C_α is calculated from the standard normal distribution.

$$C_\alpha = \Phi^{-1}\left(\alpha + \frac{1-\alpha}{2}\right)$$

Confidence intervals when sample size is large

- The 95% confidence interval of μ_Y

$$[\bar{Y} - 1.96SE(\bar{Y}), \bar{Y} + 1.96SE(\bar{Y})]$$

- The 90% confidence interval of μ_Y

$$[\bar{Y} - 1.64SE(\bar{Y}), \bar{Y} + 1.64SE(\bar{Y})]$$

- The 99% confidence interval of μ_Y

$$[\bar{Y} - 2.58SE(\bar{Y}), \bar{Y} + 2.58SE(\bar{Y})]$$

The p -value and Confidence Interval
when Sample is Small*

The t -statistic when sample size is small

- The t -statistic

$$t = \frac{\bar{Y} - \mu_{Y,0}}{SE(\bar{Y})}$$

is approximated by normal distribution using the central limit theorem for large samples.

- If the sample size is small, say $n < 30$, the use of central limit theorem may leads to a poor approximation, and the exact distribution of the t -statistic depends on the population distribution of Y , which can be complicated.

The t -statistic when sample size is small

- If the population distribution is known to be normal, the t -statistic has an exact Student t distribution with $n - 1$ degrees of freedom.
- For small samples with normal population distribution assumed, the hypothesis testing (p -value) and confidence intervals (C_α) should be evaluated with the Student t distribution.

Practice

- Generate a sample of 25 observations that follows a normal distribution with mean being 1.2 and variance being 4.
- Calculate the 80% confidence interval of the population mean from the generated sample using the exact formula for small samples (t distribution).
- Repeat the above quest but use the large sample approximation (normal distribution), compare the result with the previous.

```
y <- rnorm(25, 1.2, 2)
estimate <- mean(y)
se <- sd(y) / sqrt(length(y))

# using exact t distribution
df_t <- length(y) - 1
confint_t <- c(estimate - qt(0.9, df_t) * se,
               estimate + qt(0.9, df_t) * se)

# using normal approximation
confint_n <- c(estimate - qnorm(0.9) * se,
               estimate + qnorm(0.9) * se)
```

Practice: hypothesis testing of population means

- Practice the command `t.test` in hypothesis testing.

- A sample of data:

1.95,	0.31,	0.47,	1.54,	1.64,
2.99,	0.53,	1.21,	0.83,	1.45,
3.46,	2.23,	1.17,	1.16,	0.36,
1.76,	0.19,	0.43,	1.78,	1.56

- Test the hypothesis

$$H_0 : E(Y) = 1 \quad H_1 : E(Y) \neq 1$$

(*t*-statistic, standard error, *p*-value)

Hypothesis testing of population means

- The sample is generated from an F distribution with d.f. = (3, 6)
- Search “F distribution” in wikipedia, and find the theoretical mean and variance of $F_{m,n}$.
- Redo your hypothesis testing with these new information. Compare the p -values obtained from t .test, large-sample formulas with unknown/known population variance. What have you learned?

p -value for large samples

- The p -value when the population mean is unknown

$$p\text{-value} = 2\Phi\left(-\left|\frac{\bar{Y}^{act} - \mu_{Y,0}}{s_Y / \sqrt{n}}\right|\right) = 2\Phi\left(-\left|\frac{\bar{Y}^{act} - \mu_{Y,0}}{SE(\bar{Y})}\right|\right)$$

- The p -value when the population mean is known

$$\begin{aligned} p\text{-value} &= \Pr_{H_0} \left[\left| \frac{\bar{Y} - \mu_{Y,0}}{\sigma_Y / \sqrt{n}} \right| > \left| \frac{\bar{Y}^{act} - \mu_{Y,0}}{\sigma_Y / \sqrt{n}} \right| \right] \\ &= 2\Phi\left(-\left|\frac{\bar{Y}^{act} - \mu_{Y,0}}{\sigma_Y / \sqrt{n}}\right|\right) \end{aligned}$$

```

y <- c(1.95, 0.31, 0.47, 1.54, 1.64,
       2.99, 0.53, 1.21, 0.83, 1.45,
       3.46, 2.23, 1.17, 1.16, 0.36,
       1.76, 0.19, 0.43, 1.78, 1.56)
mu0 <- 1
ty <- t.test(y, mu = mu0)

# theoretical moments for F distribution with d.f. = (3,6)
d1 <- 3
d2 <- 6
pmean <- d2 / (d2 - 2)
pvar <- 2 * d2^2 * (d1 + d2 - 2) / (d1 * (d2 - 2)^2 * (d2 - 4))

# p-values under large sample assumption
estimate <- mean(y)
se <- sd(y) / sqrt(length(y))
tstat <- (estimate - mu0) / se

pvalue_un <- 2*pnorm(- abs(tstat))
pvalue_kn <- 2*pnorm(- abs((estimate - mu0) /
                           sqrt(pvar / length(y))))

```

A summary

- Sample size = 20. Sample estimate is 1.351.
- Population distribution is not normal, and population variance is known.
- `t.test` (which use Student t distribution for t -statistic and assume the population distribution is normal) gives p -value = 0.0920
- Large-sample formulas with unknown (known) population variance gives p -value = 0.0759 (0.4933)

References

1. Stock, J. H. and Watson, M. M., *Introduction to Econometrics*, 3rd Edition, Pearson, 2012.