

Chapter 6. Confidence intervals and hypothesis testing

- An interval $[u, v]$ constructed using the data \mathbf{y} is said to **cover** a parameter θ if $u \leq \theta \leq v$.
- $[u, v]$ is a 95% **confidence interval** (CI) for θ if the same construction, applied to a large number of draws from the model, would cover θ 95% of the time.
- A **parameter** is a name for any unknown constant in a model. In linear models, each component β_1, \dots, β_p of the **coefficient vector** β is a parameter. The only other parameter is σ , the standard deviation of the measurement error.

- A confidence interval is the usual way to represent the amount of uncertainty in an estimated parameter.
- The parameter is not random. According to the model, it has a fixed but unknown value.
- The observed interval $[u, v]$ is also not random.
- An interval $[U, V]$ constructed using a vector of random variables \mathbf{Y} defined in a probability model is random.
- If the model is appropriate, then it is reasonable to treat the observed confidence interval $[u, v]$ like a realization from the probability model.

Not quite a confidence interval for a linear model

- Consider estimating β_1 in the linear model $\mathbf{Y} = \mathbb{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ with $\boldsymbol{\epsilon} \sim \text{MVN}(\mathbf{0}, \sigma^2 \mathbb{I})$.
- Recall that $\mathbb{E}[\hat{\beta}_1] = \beta_1$ and $\text{SD}(\hat{\beta}_1) = \sigma \sqrt{[(\mathbb{X}^T \mathbb{X})^{-1}]_{11}}$.

Question 6.1. Find $\mathbb{P}(\hat{\beta}_1 - 1.96 \text{SD}(\hat{\beta}_1) \leq \beta_1 \leq \hat{\beta}_1 + 1.96 \text{SD}(\hat{\beta}_1))$

- The interval $[\hat{\beta}_1 - 1.96 \text{SD}(\hat{\beta}_1), \hat{\beta}_1 + 1.96 \text{SD}(\hat{\beta}_1)]$ is almost a confidence interval. Sadly, we don't know σ .

An approximate confidence interval for a linear model

- An approximate 95% CI for β_1 is

$$[b_1 - 1.96 \text{SE}(b_1), b_1 + 1.96 \text{SE}(b_1)]$$

where $\mathbf{y} = \mathbb{X}\mathbf{b} + \mathbf{e}$ with $\text{SE}(b_1) = s \sqrt{[(\mathbb{X}^T \mathbb{X})^{-1}]_{11}}$.

- The **standard error** $\text{SE}(b_1)$ is an estimated standard deviation of $\hat{\beta}_1$ under the linear model $\mathbf{Y} = \mathbb{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ with $\boldsymbol{\epsilon} \sim \text{MVN}(\mathbf{0}, \sigma^2 \mathbb{I})$.

A CI for association between unemployment and mortality

```
c1 <- summary(lm(L_detrended~U_detrended))$coefficients ; c1

##              Estimate Std. Error  t value    Pr(>|t|)
## (Intercept) 0.2899928 0.09343146  3.103802 0.002812739
## U_detrended 0.1313673 0.06321939  2.077959 0.041606370

beta_U <- c1["U_detrended","Estimate"]
SE_U <- c1["U_detrended","Std. Error"]
z <- qnorm(1-0.05/2) # for a 95% CI using a normal approximation
cat("CI = [", beta_U - z * SE_U, ",", beta_U + z * SE_U, "]")

## CI = [ 0.0074596 , 0.2552751 ]
```

Interpreting and criticizing a p-value

Question 6.2. We appear to have found evidence that each percentage point of unemployment above trend is associated with about 0.13 years of additional life expectancy, since the 95% CI doesn't include zero. Do you believe this discovery? How could you criticize it?

Association is not causation

“Whatever phenomenon varies in any manner whenever another phenomenon varies in some particular manner, is either a cause or an effect of that phenomenon, or is connected with it through some fact of causation.” (*John Stuart Mill, A System of Logic, Vol. 1. 1843. p. 470.*)

- Put differently: If A and B are associated statistically, we can infer that either A causes B , or B causes A , or both have some common cause C .
- A useful mantra: **Association is not causation.**
- Writing a linear model where A depends on B can show association but we need extra work to argue B causes A . We need to rule out A causing B and the possibility of any common cause C .

Association is not causation: a case study

Question 6.3. Discuss the extent to which the observed association between detrended unemployment and life expectancy in our data can and cannot be interpreted causally.

Hypothesis tests

- We try to see patterns in our data. We hope to discover phenomena that will advance science, or help the environment, or reduce sickness and poverty, or make us rich, ...
- How can we tell whether our new theory is like seeing animals or faces in the clouds?
- From Wikipedia: “**Pareidolia** is a psychological phenomenon in which the mind responds to a stimulus ... by perceiving a familiar pattern where none exists (e.g. in random data)”.
- The research community has set a standard: The evidence presented to support a new theory should be unlikely under a **null hypothesis** that the new theory is false. To quantify *unlikely* we need a probability model.

Hypothesis tests and the scientific method

- From a different perspective, a standard view of scientific progress holds that scientific theories cannot be proved correct, they can only be falsified (<https://en.wikipedia.org/wiki/Falsifiability>).
- Accordingly, scientists look for evidence to refute the **null hypothesis** that data can be explained by current scientific understanding.
- If the null hypothesis is inadequate to explain data, the scientist may propose an **alternative hypothesis** which better explains these data.
- The alternative hypothesis will subsequently be challenged with new data.

The scientific method in statistical language

- 1 **Ask a question**
- 2 **Obtain relevant data.**
- 3 **Write a null and alternative hypothesis to represent your question in a probability model.** This may involve writing a linear model so that $\beta_1 = 0$ corresponds to the null hypothesis of “no effect” and $\beta_1 \neq 0$ is a discovered “effect.”
- 4 **Choose a test statistic.** The **sample test statistic** is a quantity computed using the data summarizing the evidence against the null hypothesis. For our linear model example, the least squares coefficient b_1 is a natural sample test statistic for the hypothesis $\beta_1 = 0$.
- 5 **Calculate the p-value**, the probability that a **model-generated test statistic** is at least as extreme as that observed. For our linear model example, the p-value is $P(|\hat{\beta}_1| > |b_1|)$. We can find this probability, when $\beta_1 = 0$, using a normal approximation.
- 6 **Conclusions.** A small p-value (often, < 0.05) is evidence for **rejecting** the null hypothesis. The data analysis may suggest new questions: **Return to Step 1.**

Using confidence intervals to construct a hypothesis test

- It is often convenient to use the confidence interval as a sample test statistic.
- If the confidence interval doesn't cover the null hypothesis, then we have evidence to reject that null hypothesis.
- If we do this test using a 95% confidence interval, we have a 5% chance that we reject the null hypothesis if it is true. This follows from the definition of a confidence interval: whatever the true unknown value of a parameter θ , a model-generated confidence interval covers θ with probability 0.95.

Some notation for hypothesis tests

- The null hypothesis is H_0 and the alternative is H_a .
- We write t for the sample test statistic calculated using the data \mathbf{y} . We write T for the model-generated test statistic, which is a random variable constructed by calculating the test statistic using a random vector \mathbf{Y} drawn from the probability model under H_0 .
- The p-value is $\text{pval} = P(|T| \geq |t|)$. Here, we are assuming “extreme” means “large in magnitude.” Occasionally, it may make more sense to use $\text{pval} = P(T \geq t)$.
- We reject H_0 at **significance level** α if $\text{pval} < \alpha$. Common choices of α are $\alpha = 0.05$, $\alpha = 0.01$, $\alpha = 0.001$.

Alternative ways to report a hypothesis test

Question 6.4. When we report the results of a hypothesis test, we can either (i) give the p-value, or (ii) say whether H_0 is rejected at a particular significance level. What are the advantages and disadvantages of each?

Terminology for test statistics

- Recall that a **sample test statistic** is a summary of the data, constructed to test a hypothesis.
- A **model-generated test statistic** is the same summary applied to random variables drawn from a probability model. Usually, this probability model represents the null hypothesis. We can say “model-generated test statistic under H_0 ” to make this explicit.
- Distinguishing between sample test statistics and model-generated ones under a null hypothesis is critical to the logic of hypothesis testing.

Example: testing whether $\beta_1 = 0$ in the linear model $\mathbf{Y} = \mathbb{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$,

- The sample test statistic is $b_1 = [(\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T \mathbf{y}]_1$.
- A model-generated test statistic is $\hat{\beta}_1 = [(\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T \mathbf{Y}]_1$.

A hypothesis test for unemployment and mortality

Question 6.5. Write a formal hypothesis test of the null hypothesis that there is no association between unemployment and mortality. Compute a p-value using a normal approximation. What do you think is an appropriate significance level α for deciding whether to reject the null hypothesis?

Normal approximations versus Student's t distribution

- Notice that `summary(lm(...))` gives `tvalue` and `Pr(>|t|)`.
- The `tvalue` is the estimated coefficient divided by its standard error. This measures how many standard error units the estimated coefficient is from zero.
- `Pr(>|t|)` is similar, but slightly larger, than the p-value coming from the normal approximation.
- R is using Student's t distribution, which makes allowance for chance variation from using s as an approximation to σ when we compute the standard error.
- R uses a t random variable to model the distribution of the statistic t . Giving the full name (Student's t distribution) may add clarity.
- With sophisticated statistical methods, it is often hard to see if they work well just by reading about them. Fortunately, it is often relatively easy to do a **simulation study** to see what is going on.

Simulating from Student's t distribution

- Suppose X and X_1, \dots, X_d are independent identically distributed (iid) normal random variables with mean zero and standard deviation σ .
- Student's t distribution on d degrees of freedom is defined to be the distribution of $T = X/\hat{\sigma}$ where $\hat{\sigma} = \sqrt{\frac{1}{d} \sum_{i=1}^d X_i^2}$.
- A normal approximation would say T is approximately $\text{normal}(0, 1)$ since $\hat{\sigma}$ is an estimate of σ .
- With a computer, we can simulate T many times, plot a histogram, and compare it to the probability density function of the normal distribution and Student's t distribution.

Question 6.6. This is almost the same representation of the t distribution as HW4. What is the difference? Why does it not matter?

- Let's start by simulating a matrix X of iid normal random variables.

```
N <- 50000 ; sigma <- 1 ; d <- 10 ; set.seed(23)
X <- matrix(rnorm(N*(d+1),mean=0,sd=sigma),nrow=N)
```

- Now, we write a function that computes T given X_1, \dots, X_d, X

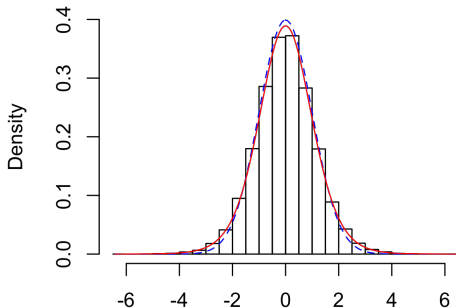
```
T_evaluator <- function(x) x[d+1] / sqrt(sum(x[1:d]^2)/d)
```

- Then, use `apply()` to evaluate T on each row of 'X'.

```
Tsim <- apply(X,1,T_evaluator)
```

- A histogram of these simulations can be compared with the normal density and the t density

```
hist(Tsim,freq=F,main="",
     breaks=30,ylim=c(0,0.4))
x <- seq(length=200,
        min(Tsim),max(Tsim))
lines(x,dnorm(x),
      col="blue",
      lty="dashed")
lines(x,dt(x,df=d),
      col="red")
```



Comparing the normal and t distributions

- Even with as few as $d = 10$ degrees of freedom to estimate σ , the Student's t density looks similar to the normal density.
- Student's t has fatter tails. This is important for the probability of rare extreme outcomes.
- Here, the largest and smallest of the $N = 5 \times 10^4$ simulations are

```
range(Tsim)
## [1] -6.438830  6.480262
```

- Let's check the chance of an outcome more than 5 (or 6) standard deviations from the mean for the normal distribution and the t on 10 degrees of freedom.

```
2*(1-pnorm(5))
## [1] 5.733031e-07
2*(1-pnorm(6))
## [1] 1.973175e-09
```

```
2*(1-pt(5,df=d))
## [1] 0.0005373336
2*(1-pt(6,df=d))
## [1] 0.0001321089
```

Hypotheses about predictions from a linear model

- Consider the sample linear model $\mathbf{y} = \mathbb{X}\mathbf{b} + \mathbf{e}$, where $\mathbb{X} = [x_{ij}]_{n \times p}$.
- We might be interested in predicting outcomes at some new set of explanatory variables $\mathbf{x}^* = (x_1^*, \dots, x_p^*)$, treated as a $1 \times p$ **row vector**.
- Making a prediction involves estimating (i) the expected value of a new outcome; (ii) its variability. In addition, we must make allowance for the statistical uncertainty in these estimates.
- To do inference, we need a probability model. As usual, consider $\mathbf{Y} = \mathbb{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ where $\epsilon_1, \dots, \epsilon_n \sim \text{iid normal}(0, \sigma)$. Also, model a new measurement at \mathbf{x}^* as

$$Y^* = \mathbf{x}^* \boldsymbol{\beta} + \epsilon^*$$

where ϵ^* is another independent draw from the measurement model.

Question 6.7. (a) Why do we want \mathbf{x}^* to be a row vector not a column vector? (b) What is the dimension of $\mathbf{x}^* \boldsymbol{\beta}$?

The expected value of a new outcome and its uncertainty

- According to the model, the expected value of a new outcome at \mathbf{x}^* is

$$E[Y^*] = \mathbf{x}^* \boldsymbol{\beta}.$$

- But, we don't know $\boldsymbol{\beta}$. We estimate $\boldsymbol{\beta}$ by the sample least squares coefficient $\mathbf{b} = (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T \mathbf{y}$, which is modeled as a realization of the model-generated least squares coefficient $\hat{\boldsymbol{\beta}} = (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T \mathbf{Y}$.

- A **sample estimate of the expected value** is the **fitted value** at \mathbf{x}^*

$$\hat{y}^* = \mathbf{x}^* \mathbf{b} = \sum_{j=1}^p x_j^* b_j.$$

- The **model-generated estimate of the expected value** is

$$\hat{Y}^* = \mathbf{x}^* \hat{\boldsymbol{\beta}} = \sum_{j=1}^p x_j^* \hat{\beta}_j.$$

- We can find the mean and variance of \hat{Y}^* . We can use these (together with a normal approximation) to find a confidence interval for $E[Y^*]$. If the model is reasonable, this will tell us the uncertainty in using \hat{y}^* to estimate the sample average of many new outcomes collected at \mathbf{x}^* .

Question 6.8. Use linearity of expectation to show that $E[\hat{Y}^*] = \mathbf{x}^* \boldsymbol{\beta}$

Question 6.9. Use the formula $\text{Var}(\mathbb{A}\mathbf{X}) = \mathbb{A}\text{Var}(\mathbf{X})\mathbb{A}^T$ to show that $\text{Var}[\hat{Y}^*] = \sigma^2 \mathbf{x}^* (\mathbb{X}^T \mathbb{X})^{-1} \mathbf{x}^{*T}$

Question 6.10. Check the dimension of $\text{Var}[\hat{Y}^*]$. Is this correct?

A CI for the expected value of a new outcome

- We can get a confidence interval (CI) for the **linear combination of coefficients** $\mathbf{x}^*\boldsymbol{\beta}$ in a similar way to what we did for a single coefficient.
- A standard error is $\text{SE}(\mathbf{x}^*\mathbf{b}) = s \sqrt{\mathbf{x}^*(\mathbb{X}^T\mathbb{X})^{-1}\mathbf{x}^{*T}}$.
- Then, making a normal approximation, a 95% CI is $[\mathbf{x}^*\mathbf{b} - 1.96 \text{SE}(\mathbf{x}^*\mathbf{b}), \mathbf{x}^*\mathbf{b} + 1.96 \text{SE}(\mathbf{x}^*\mathbf{b})]$.

Example. We consider again the data on freshman GPA, ACT exam scores and percentile ranking of each student within their high school for 705 students at a large state university. We seek to predict using the probability model considered in the midterm exam, where freshman GPA is modeled to depend linearly on ACT score and high school ranking.

```
gpa <- read.table("gpa.txt",header=T); gpa[1,]
```

```
##    ID  GPA High_School ACT Year
## 1   1 0.98          61  20 1996
```


Worked example 6.1. Find a 95% confidence interval for the expected freshman GPA among students with an ACT score of 20 ranking at the 40th percentile in his/her high school.

```
lm1 <- lm(GPA~ACT+High_School,data=gpa)
x <- c(1,20,40)
pred <- x%*%coef(lm1)
V <- summary(lm1)$cov.unscaled
s <- summary(lm1)$sigma
SE_pred <- sqrt(x%*%V%*%x)*s
c <- qnorm(0.975)
cat("CI = [", round(pred-c*SE_pred,3),
    ", ", round(pred+c*SE_pred,3), "]" )

## CI = [ 2.344 , 2.532 ]
```

Question 6.11. How would you check whether your answer is plausible? How would you check the R calculation has done what you want it to do?

A prediction interval for a new outcome

- A 95% **prediction interval** for a new outcome of a linear model with explanatory variables \mathbf{x}^* covers the outcome with probability 95%.
- The prediction interval allows for the uncertainty around the mean, modeled as **measurement error** in the outcome.
- The prediction interval aims to cover $Y^* = \mathbf{x}^* \boldsymbol{\beta} + \epsilon^*$ whereas the confidence interval for the mean only aims to cover $E[Y^*] = \mathbf{x}^* \boldsymbol{\beta}$.
- Since ϵ^* is independent of $\mathbf{x}^* \hat{\boldsymbol{\beta}}$ we have

$$\begin{aligned}\text{Var}[Y^* - \mathbf{x}^* \hat{\boldsymbol{\beta}}] &= \text{Var}[Y^* - \mathbf{x}^* \boldsymbol{\beta}] + \text{Var}[\mathbf{x}^* \boldsymbol{\beta} - \mathbf{x}^* \hat{\boldsymbol{\beta}}] \\ &= \sigma^2 + \sigma^2 \mathbf{x}^* (\mathbb{X}^T \mathbb{X})^{-1} \mathbf{x}^{*T}\end{aligned}$$

- This suggests using a standard error for prediction of

$$\text{SE}_{\text{pred}} = s \sqrt{1 + \mathbf{x}^* (\mathbb{X}^T \mathbb{X})^{-1} \mathbf{x}^{*T}}$$

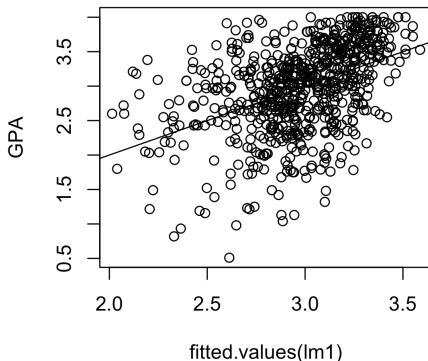
- A 95% prediction interval, using a normal approximation, is

$$[\mathbf{x}^* \mathbf{b} - 1.96 \text{SE}_{\text{pred}}, \mathbf{x}^* \mathbf{b} + 1.96 \text{SE}_{\text{pred}}]$$

Using the t distribution for predictions

- We could use a t quantile instead of a normal approximation.
- Just as for parameter confidence intervals, since we use the sample standard deviation s in place of the true standard deviation σ , a t distribution is more accurate.
- With 705 observations, the normal quantile $1.96 = \text{qnorm}(0.975)$ is identical to $1.96 = \text{qt}(0.975, \text{df}=702)$ up to 3 significant figures.

```
plot(x=fitted.values(lm1),y=gpa$GPA,ylab="GPA")  
abline(a=0,b=1)
```



Question 6.12. Is the linear model a good fit for the data? What cautions do you recommend when using this model for prediction?

Worked example 6.2. Find a 95% prediction interval for the freshman GPA of an incoming student with an ACT score of 20 ranking at the 40th percentile in his/her high school.

```
lm1 <- lm(GPA~ACT+High_School,data=gpa)
x <- c(1,20,40)
pred <- x%*%coef(lm1)
V <- summary(lm1)$cov.unscaled
s <- summary(lm1)$sigma
SE_pred <- sqrt(x%*%V%*%x + 1)*s
c <- qnorm(0.975)
cat("prediction interval = [", round(pred-c*SE_pred,3),
    ",", round(pred+c*SE_pred,3), "]" )

## prediction interval = [ 1.322 , 3.553 ]
```

Question 6.13. Where does this calculation differ from the confidence interval?