

# Lecture 3: Designing simulations

# Last time

Potential  
results

Coverage:

$\frac{\varepsilon_i \sim \text{Normal}}{95\%}$

$\frac{\varepsilon_i \sim \text{Exp}}{95\%}$

$\frac{\varepsilon_i \sim \chi^2}{95\%}$

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

How would you study the importance of the normality assumption?

One approach:

- Simulate data with different distributions for  $\varepsilon_i$   
e.g. Normal,  $\chi^2$ , exponential, etc.
- Fit the linear regression model and calculate a 95% confidence interval for  $\beta_1$ . (Ideally, 95% of these intervals contain  $\beta_1$ .)

If coverage is  
 $< 95\%$ , or  
 $> 95\%$ , the  
normality  
is (important)

- Repeat many times; does the confidence interval have the coverage (e.g. do 95% of the intervals actually contain  $\beta_1$ )?
- Compare coverage for different distribution for  $\varepsilon_i$

# Simulating data

To start, simulate data for which the normality assumption holds:

```
1 n <- 100 # sample size
2 beta0 <- 0.5 # intercept  $\leftarrow \beta_0 = 0.5$ 
3 beta1 <- 1 # slope  $\leftarrow \beta_1 = 1$ 
4
5 x <- runif(n, min=0, max=1)
6 noise <- rnorm(n, mean=0, sd=1)
7 y <- beta0 + beta1*x + noise
```

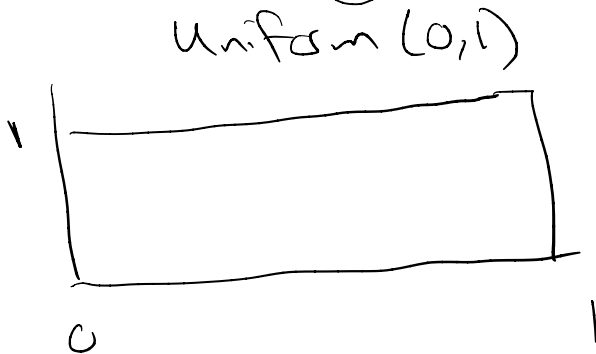
$(x_1, y_1), \dots, (x_n, y_n)$  ( $n$  observations)

$x_i \sim \text{Uniform}(0,1)$

$\varepsilon_i \sim \text{Normal}(0,1)$

$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$

- `runif(n, min=0, max=1)` samples  $X_i$  uniformly between 0 and 1
- `rnorm(n, mean=0, sd=1)` samples  $\varepsilon_i \sim N(0,1)$



# Fit a model

```
1 n <- 100 # sample size
2 beta0 <- 0.5 # intercept
3 beta1 <- 1 # slope
4
5 x <- runif(n, min=0, max=1)
6 noise <- rnorm(n, mean=0, sd=1)
7 y <- beta0 + beta1*x + noise
8
9 lm_mod <- lm(y ~ x)
10 lm_mod
```

generating data  
( $x_s$  and  $y_s$ )

response  
explanatory

Call:

```
lm(formula = y ~ x)
```

Coefficients:

(Intercept)

0.2836

$\hat{\beta}_0$

x  
1.4302

$\hat{\beta}_1$

95% CI for  $\beta_1$ :

$$\hat{\beta}_1 \pm t_{n-2}^* SE \hat{\beta}_1$$

$n=2$

# Calculate confidence interval

```
1 lm_mod <- lm(y ~ x)
2
3 ci <- confint(lm_mod, "x", level = 0.95)
4 ci
```

2.5 %    97.5 %  
x 0.6883911 2.172003

fitted model  
95% CI  
coefficient of interest (eg.  $\beta_1$ )

- **Question:** How can we check whether the confidence interval contains the true  $\beta_1$ ?

$$\beta_1 = 1$$

$$0.688 < 1 \quad \& \quad 2.172 > 1 \quad (\text{TRUE})$$

$$ci[1] < 1 \quad \& \quad ci[2] > 1$$

# Calculate confidence interval

```
1 lm_mod <- lm(y ~ x)
2
3 ci <- confint(lm_mod, "x", level = 0.95)
4 ci
```

```
      2.5 %    97.5 %
x 0.6883911 2.172003
```

- **Question:** How can we check whether the confidence interval contains the true  $\beta_1$  ?

```
1 ci[1] < 1 & ci[2] > 1
```

```
[1] TRUE
```

$ci[1] < \beta_1$       &       $ci[2] > \beta_1$

# Repeat!

```
1 nsim <- 1000
2 n <- 100 # sample size
3 beta0 <- 0.5 # intercept
4 beta1 <- 1 # slope
5 results <- rep(NA, nsim)
6
7 for(i in 1:nsim){
8   x <- runif(n, min=0, max=1)
9   noise <- rnorm(n, mean=0, sd=1)
10  y <- beta0 + beta1*x + noise
11
12  lm_mod <- lm(y ~ x)
13  ci <- confint(lm_mod, "x", level = 0.95)
14
15  results[i] <- ci[1] < 1 & ci[2] > 1
16 }
17 mean(results)
```

} Sample data at each iteration

} fit model, calculate a 95% CI

← check if CI contains  $\beta_1$ , store result

- What fraction of the time should the confidence interval contain  $\beta_1$ ?

expect  $\approx 0.95$

# Repeat!

```
1 nsim <- 1000
2 n <- 100 # sample size
3 beta0 <- 0.5 # intercept
4 beta1 <- 1 # slope
5 results <- rep(NA, nsim)
6
7 for(i in 1:nsim){
8   x <- runif(n, min=0, max=1)
9   noise <- rnorm(n, mean=0, sd=1)
10  y <- beta0 + beta1*x + noise
11
12  lm_mod <- lm(y ~ x)
13  ci <- confint(lm_mod, "x", level = 0.95)
14
15  results[i] <- ci[1] < 1 & ci[2] > 1
16 }
17 mean(results)
```

next step: try a  
different distribution  
for  $\epsilon_i$

```
[1] 0.952
```

- What should we do next?



# Class activity

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

That is, how important is the assumption that  $\varepsilon_i \sim N(0, \sigma^2)$ ?

Continue simulation from last time, but experiment with different values of  $n$  and different distributions for the noise term.

[https://sta279-f23.github.io/class\\_activities/ca\\_lecture\\_3.html](https://sta279-f23.github.io/class_activities/ca_lecture_3.html)

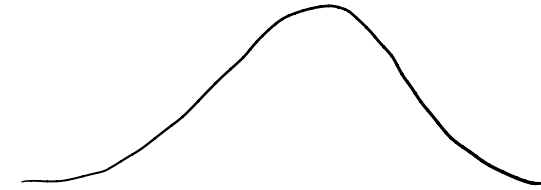
# Class activity

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

How does confidence interval coverage change when you change the distribution of  $\varepsilon_i$ ?

# Class activity

```
1 nsim <- 1000
2 n <- 100 # sample size
3 beta0 <- 0.5 # intercept
4 beta1 <- 1 # slope
5 results <- rep(NA, nsim)
6
7 for(i in 1:nsim){
8   x <- runif(n, min=0, max=1)
9   noise <- rchisq(n, 1) ←
10  y <- beta0 + beta1*x + noise
11
12  lm_mod <- lm(y ~ x)
13  ci <- confint(lm_mod, "x", level = 0.95)
14
15  results[i] <- ci[1] < 1 & ci[2] > 1
16 }
17 mean(results)
```



$\chi^2_1$

```
[1] 0.963
```

Exp(L)  $\approx 95\%$

$\chi^2_1$   $\approx 95\%$

Normal = 95%



# ADEMP: A useful framework for simulation studies

- **Aims:** Why are we doing the study?
- **Data generation:** How are the data simulated?
- **Estimand/target:** What are we estimating for each simulated dataset?
- **Methods:** What methods are we using for model fitting, estimation, etc?
- **Performance measures:** How do we measure performance of our chosen methods?

# ADEMP

For the normal errors simulation study:

- **Aims:** Assess importance of the normality assumption
- **Data generation:**  $X_i \sim \text{Uniform}(0,1)$        $Y_i = 0.5 + X_i + \varepsilon_i$
- **Estimand/target:**  $\beta_1$        $\varepsilon_i \sim \text{Uniform}(0,1)$  or  $\varepsilon_i \sim \text{Exp}(1)$   
or  $\varepsilon_i \sim N(0,1)$  or  $\varepsilon_i \sim \chi^2_{1,\dots}$
- **Methods:** Fit linear model in R, calculate 95% CIs for  $\beta_1$
- **Performance measures:** observed coverage of confidence intervals

