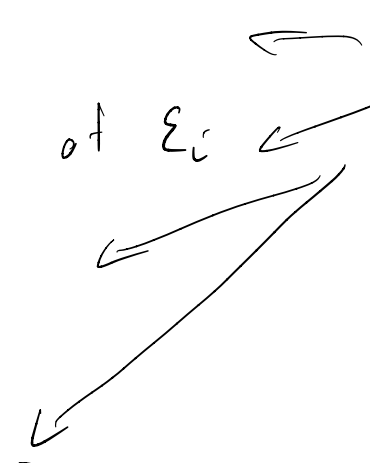# Lecture 3: Beginning statistical simulations

# A new question

In STA 112, you learned about the simple linear regression model:

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

**Question:** What assumptions does this model make?

- Normality of $\varepsilon_i$
- Constant variance of $\varepsilon_i$
- independence of $\varepsilon_i$
- linearity (shape)
- $\varepsilon_i$ have mean 0
- randomness

$$\varepsilon_i \overset{iid}{\sim} N(0, \sigma^2)$$

# A new question

In STA 112, you learned about the simple linear regression model:

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

**Question:** How important is it that $\varepsilon_i \sim N(0, \sigma^2)$? Does it matter if the errors are *not* normal?

# Activity

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

**Activity:** With a neighbor, brainstorm how you could use simulation to assess the importance of the normality assumption (you do not need to write code!).

- How would you simulate data?

- What result would you measure for each run of the simulation?

# Activity

Normal

Exp(1)

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$
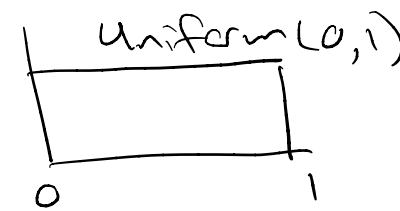
How would you study the importance of the normality assumption?

95% CI: $\hat{\beta_1} \pm t^*_{n-2} SE(\hat{\beta_1})$

depends on $\varepsilon_i \sim$ Normal

- do the simulations many times

- simulate data with different distributions for $\varepsilon_i$

  e.g. Poisson distribution, binomial distribution, exponential, $\chi^2$, Normal

  compare performance for different distributions of $\varepsilon_i$

  e.g. · SSE

  · confidence intervals for $\beta_1$

- Construct 95% CIs for $\beta_1$ and check whether they capture the $\beta_1$ 95% of time in practice

# Simulating data

Uniform(0,1)

Normal

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

```
1  n <- 100 # sample size
2  beta0 <- 0.5 # intercept        ← β_0
3  beta1 <- 1 # slope              ← β_1
4
5  x <- runif(n, min=0, max=1)
6  noise <- rnorm(n, mean=0, sd=1)
7  y <- beta0 + beta1*x + noise
```

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

↑

Can have whatever distribution we want

- `runif(n, min=0, max=1)` samples $X_i$ uniformly between 0 and 1

number of samples

- `rnorm(n, mean=0, sd=1)` samples $\varepsilon_i \sim N(0,1)$

simulate
n observations
from Uniform(0,1)

r unif

"random" ↗    abbrev. for
distribution name

other dists:  rexp  , rchisq , rpois , etc.

# Fit a model

```r
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 $x_s$ & $y_s$

$\leftarrow$ fit linear regression model

response      explanatory

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

Coefficients:
(Intercept)           x
     0.2971      1.3073
```

$\hat{\beta}_0$          $\hat{\beta}_1$

True value for $\beta_1 = 1$

compare $\hat{\beta}_1$ to $\beta_1$

compare a CI for $\beta_1$
to the value of $\beta_1$

# 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 0.6777885 1.936822
```

*calculate* a *CI*

← 95% CI

*fitted model*  *coefficient of interest* $(\beta_1)$

*vector*

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

$$CI = (0.68, 1.94) \quad \text{contains } \beta_1 = 1$$

$$0.68 < 1 \quad \& \quad 1.94 > 1 \quad (TRUE)$$

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

# Calculate confidence interval

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

```
        2.5 %    97.5 %
x 0.6777885 1.936822
```

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

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

```
[1] TRUE
```

# Repeat!

```r
nsim <- 1000
n <- 100 # sample size
beta0 <- 0.5 # intercept
beta1 <- 1 # slope
results <- rep(NA, nsim)

for(i in 1:nsim){
  x <- runif(n, min=0, max=1)
  noise <- rnorm(n, mean=0, sd=1)
  y <- beta0 + beta1*x + noise

  lm_mod <- lm(y ~ x)
  ci <- confint(lm_mod, "x", level = 0.95)

  results[i] <- ci[1] < 1 & ci[2] > 1
}
mean(results)
```

*Handwritten annotations:*
- Set up parameters (lines 1–4)
- ⟵ results vector to store results (line 5)
- Sample data (lines 8–10)
- fit model & calculate 95% CI (lines 12–13)
- ⟵ check if CI contains $\beta_1$, store result (line 15)
- ⟵ observed coverage (line 17)

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

*Handwritten:* expect $\approx 0.95$

# Repeat!

```r
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: change distribution of $\xi_i$

```
[1] 0.948
```

- 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-s24.github.io/class_activities/ca_lecture_3.html