# SOC6302 Statistics for Sociologists

Monica Alexander

Week 9: Simple Linear Regression II

# Announcements

- A3 and RQ and EDA due next week
- A3 optional
- RQ: primary and second questions, ideally
  - e.g. Question: are people born outside of Canada more likely to start having children later compared to those born in Canada?
  - Does the answer to this question persist after we take into account education?
- Note that your dependent variable must be a continuous variable
- But explanatory variables can be anything

# Overview

- Hypothesis testing of coefficients
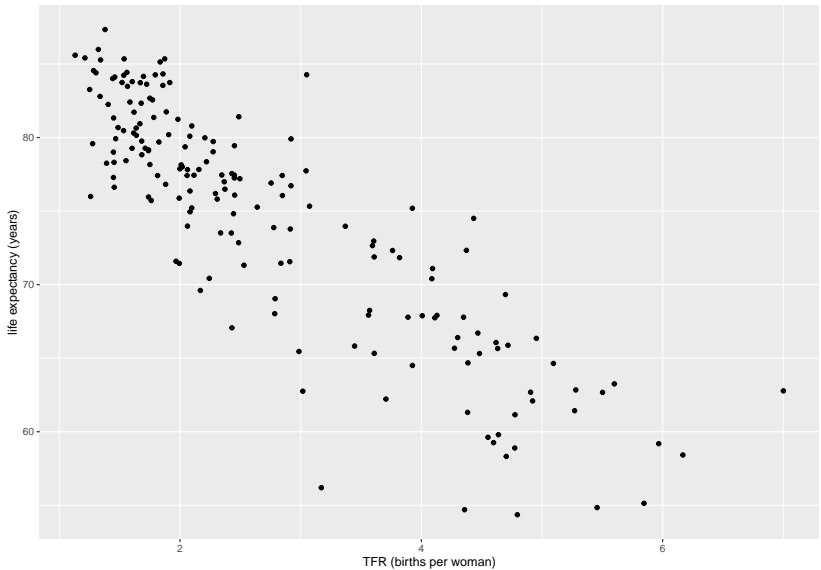- Log transforms
- Start multiple linear regression

# Review of SLR set-up

- $Y_i$ is the response variable, and $X_i$ is the explanatory variable

Example:

- Research question: In 2017, how does the expected value of life expectancy differ or change across countries with different levels of fertility?
- In other words, is life expectancy associated with fertility, and if so, how?

# Scatter plot

# Fit SLR in R

```
country_ind_2017 <- country_ind %>% filter(year==2017)
slr_mod <- lm(life_expectancy~tfr, data = country_ind_2017)
summary(slr_mod)
```

```
##
## Call:
## lm(formula = life_expectancy ~ tfr, data = country_ind_2017)
##
## Residuals:
##      Min      1Q   Median       3Q      Max
## -16.0718  -2.3864   0.3132   2.6537  11.3498
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)  89.2394     0.7085  125.95   <2e-16 ***
## tfr          -5.3526     0.2326  -23.02   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.994 on 174 degrees of freedom
## Multiple R-squared:  0.7527, Adjusted R-squared:  0.7513
## F-statistic: 529.7 on 1 and 174 DF,  p-value: < 2.2e-16
```

# Sums of squares

- Variation of data $Y_i$ around the observed mean $\bar{Y}_i$
  - Total sum of squares SST: $(Y_i - \bar{Y}_i)^2$
  - Total variation in $Y_i$
- Variation of fitted values $\hat{Y}_i$ around observed mean $\bar{Y}_i$
  - Model sum of squares SSM: $(\hat{Y}_i - \bar{Y}_i)^2$
  - Variation explained by our $X$'s
- Variation of data $Y_i$ around fitted values $\hat{Y}_i$
  - Residual sum of squares SSR: $(Y_i - \hat{Y}_i)^2$
  - Variation not explained by $X$'s

$$SST = SSM + SSR$$

# $R^2$

$$SST = SSM + SSR$$

$$R^2 = \frac{SSM}{SST} = 1 - \frac{SSR}{SST}$$

The proportion of total variation in $Y_i$ explained by covariates $X_i$.

# Hypothesis testing

# Fit SLR in R

```
country_ind_2017 <- country_ind %>% filter(year==2017)
slr_mod <- lm(life_expectancy~tfr, data = country_ind_2017)
summary(slr_mod)
```

```
##
## Call:
## lm(formula = life_expectancy ~ tfr, data = country_ind_2017)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -16.0718  -2.3864   0.3132   2.6537  11.3498
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)  89.2394     0.7085  125.95   <2e-16 ***
## tfr          -5.3526     0.2326  -23.02   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.994 on 174 degrees of freedom
## Multiple R-squared:  0.7527,	Adjusted R-squared:  0.7513
## F-statistic: 529.7 on 1 and 174 DF,  p-value: < 2.2e-16
```

# SLR fit

- ▶ The estimate of $\hat{\beta}_1$ tells us there's a negative association between TFR and life expectancy as estimated from the data
- ▶ But how sure of this are we? It's not a perfect relationship, and there is some noise
- ▶ It was reasonably clear from our scatterplot, but what if our scatter plot had looked different?

# Intuition of hypothesis testing

- ▶ We are assuming there's some underlying $\beta_1$ that we're trying to find (this assumes the truth is a linear relationship)
- ▶ We get an estimate of $\beta_1$ (called $\hat{\beta}_1$) based on data we collect
- ▶ But this estimate could be right, almost right, or completely wrong compared to the truth
- ▶ In regression we are usually interested in deciding whether we believe $\beta_1$ is non-zero (i.e. there is a linear association between our two variables)
- ▶ The degree to which we believe this depends on what the data look like

# Intuition of hypothesis testing

- If the data look a lot like a linear relationship, then we conclude that there's enough evidence to suggest a non-zero relationship and that our estimate is probably right
- The more randomness there is in the data, the less likely we are to believe our estimate is the truth
- Hypothesis testing (based on t-tests) is a way of accounting for this uncertainty and making inferences about the relationships between variables

# Intuition of hypothesis testing

How do we account for the uncertainty in the data before making decisions about whether $\beta_1$ is zero or not?

- ▶ The regression model has a bunch (five) of assumptions underlying it
- ▶ If we assume these are true, then it turns out we know what the probability distribution of possible values of $\hat{\beta}_1$ look like
- ▶ If a lot of probability density in this distribution is near zero (read: if zero is likely), then we would conclude there's not enough evidence to suggest a linear relationship
- ▶ And vice versa

# To dos

To do:

- ▶ Learn assumptions
- ▶ Write down distribution for $\hat{\beta}_1$
- ▶ Do hypothesis testing
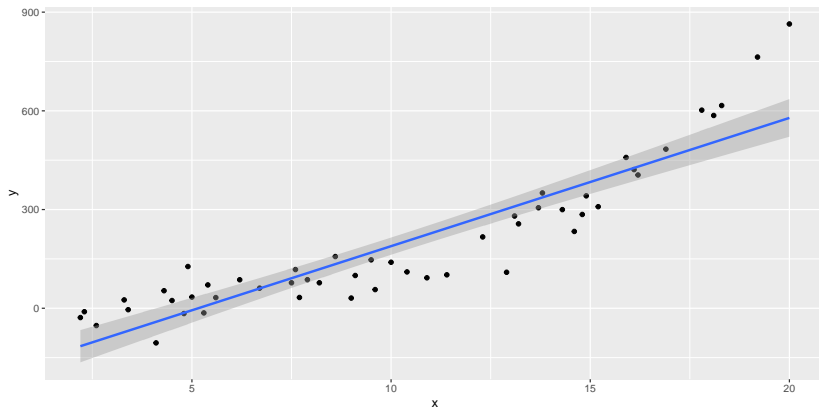- ▶ Celebrate, eat cake, graduate

# Assumption 1

**1. no model mis-specification**

- ▶ This assumption means that the dependent variable must be a simple linear function of the explanatory variable.

# Assumption 1

## 1. no model mis-specification

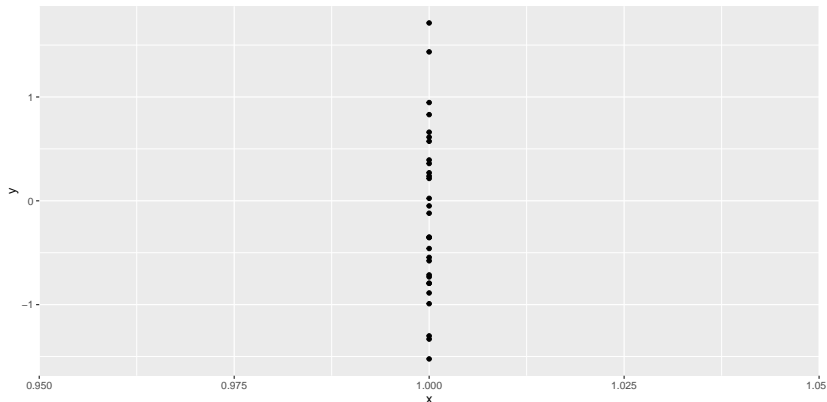Example violation ($y$ is a function of $x^3$)

# Assumption 2

## 2. $X_i$ is not a constant

If there is no variation in $X_i$, then there is not a unique solution to $\hat{\beta}_1$
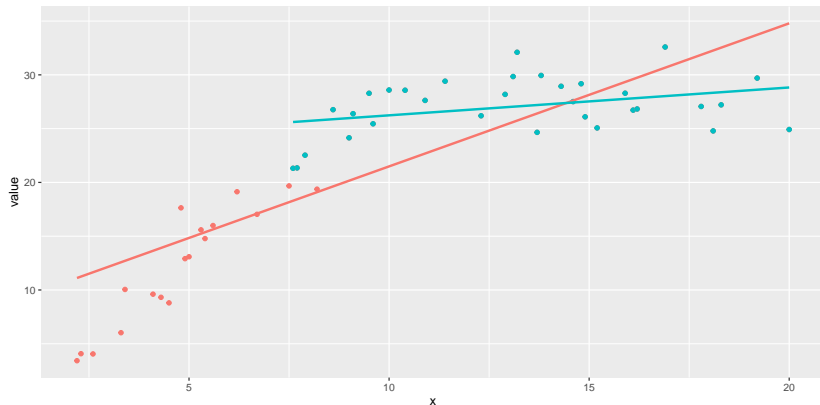
Example violation:

**3.** $\{X_i, Y_i\}$ **are from a simple random sample**

This assumption implies that all members of a population have an equal probability of selection, that all possible samples of size *n* have an equal probability of selection, and that each observation is independent of all the others

# Assumption 3

**3. $\{X_i, Y_i\}$ are from a simple random sample**

Example of violation: sample in which $\{X_i, Y_i\}$ are only observed if $Y_i > 20$

**4. The variance of $\varepsilon_i$ is the same across all values of $X_i$**
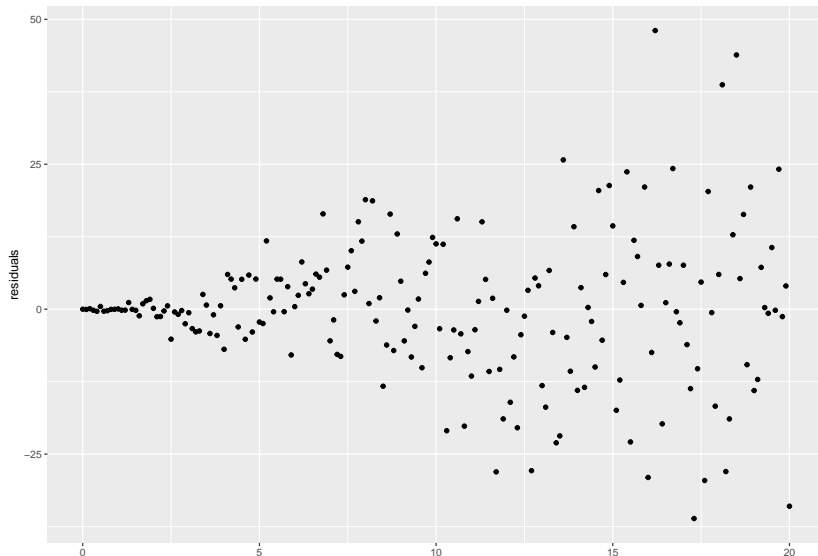
That is,

$$\text{Var}\left(\varepsilon_i \mid X_i\right) = \sigma^2$$

▶ If, for example, the variance of $\varepsilon_i$ is larger for higher values of $X_i$, then this assumption is violated

▶ When the error variance is constant across $X_i$, it is called "homoscedastic"

▶ When the error variance is non-constant across $X_i$, it is called "heteroscedastic"

# Assumption 4

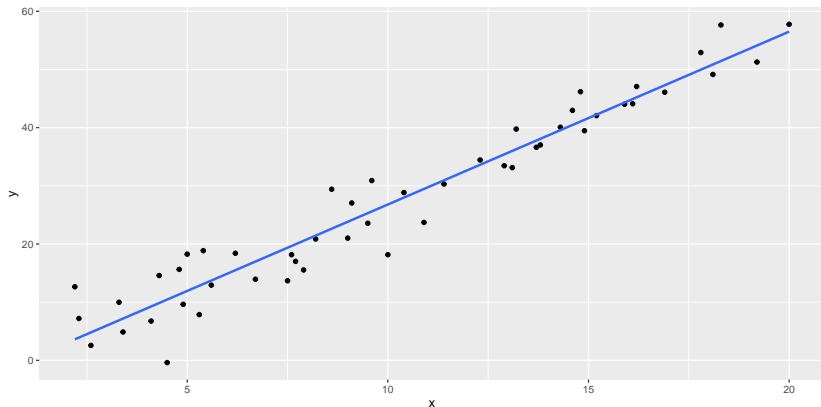## 4. The variance of $\varepsilon_i$ is the same across all values of $X_i$

Example violation (heteroskedasticity):

# Assumption 5

**5. The normality assumption:** $\varepsilon_i$ **is normally distributed**

- ▶ If, for example, the distribution of $\varepsilon_i$ is skewed or has "heavy" tails, then this assumption is violated

# Assumptions

The first four assumptions are the **Gauss-Markov assumptions**

1. GM1: No model mis-specification (linearity in covariates)
2. GM2: $X_i$ is not constant
3. GM3: Simple random sampling
4. GM4: Constant variance of $\varepsilon_i$
5. Normality assumption of $\varepsilon_i$

Taken together, these assumptions imply that

$$y|x \sim N\left(\beta_0 + \beta_1 x, \sigma^2\right)$$

# Assumptions taken as correct, standardizing

▶ If we take the five assumptions above as given, it turns out that the distribution of possible values of our estimate $\hat{\beta}_1$ around the true value $\beta_1$ is knowm

▶ In particular, we are going to look at a transformed version of $\hat{\beta}_1$:

$$\frac{\widehat{\beta}_1 - \beta_1}{se\left(\widehat{\beta}_1\right)}$$

where $se\left(\widehat{\beta}_1\right)$ is the standard error of $\hat{\beta}_1$.

▶ This should look familiar!

# The t-statistic

Let's give this quantity a name:

$$T_{\widehat{\beta}_1} = \frac{\widehat{\beta}_1 - \beta_1}{se\left(\widehat{\beta}_1\right)}$$

Given the five assumptions discussed, this follows a t-distribution with $n - (k + 1)$ degrees of freedom, where $k$ is the number of covariates/explanatory variables in the model (for simple linear regression, $k = 1$).

# What is the standard error?

The standard error of $\hat{\beta}_1$, is

$$\text{se}\left(\hat{\beta}_1\right) = \sqrt{\frac{\hat{\sigma}^2}{\sum_i \left(X_{i1} - \bar{X}_{i1}\right)^2 \left(1 - R_1^2\right)}}$$

where

$$\hat{\sigma}^2 = \frac{\sum_i \hat{\varepsilon}_i^2}{n - (k+1)} = \frac{SSR}{df}$$

and $R_1^2$ is the $R^2$ from a regression of $X_1$ against all other variables in the model

# What is the standard error?

Don't try and remember the formulas from the previous slide. Just remember that the standard error of $\hat{\beta}_1$ is proportional to the sum of squares of residuals.

▶ a larger error variance (i.e., greater unexplained variation in the outcome) is associated with a larger $\text{se}\left(\hat{\beta}_1\right)$ and vice versa

What does the standard error do to the distribution of $T_{\widehat{\beta}_1}$?

# Hypothesis testing: more intuition

- In regression, we are interested to see if there's evidence to suggest that $\beta_1$ is different enough from zero.
- Pretend for a moment that the true value of $\beta_1$ is zero. In this world (the null hypothesis world), our $T_{\widehat{\beta}_1}$ is just

$$\frac{\widehat{\beta}_1}{se\left(\widehat{\beta}_1\right)}$$

- In the null hypothesis world, this thing should be t-distributed (i.e. centered at zero with some variation around that)
- So if we calculate this thing and it's really different from zero (i.e. where the distribution is centered), then it's unlikely it came from this distribution, and we can probably reject the world in which $\beta_1$ is zero
- If this thing is not very different from zero, then we may not reject this world

# Hypothesis testing: more intuition

- ▶ We are dealing with randomness, and so there's always a chance that the value we see is from the null hypothesis world in which $\beta_1$ is zero
- ▶ But the farther away it is from zero, the less likely that's true
- ▶ The size of $T_{\widehat{\beta}_1}$ depends not only on the magnitude of $\widehat{\beta}_1$ but also the magnitude of the standard error of $\widehat{\beta}_1$
- ▶ So the stronger the relationship (the bigger the $\widehat{\beta}_1$) the less likely we are going to believe the null hypothesis
- ▶ But also for less noisy data (the smaller the standard error) the less likely we are going to believe the null hypothesis

# Hypothesis testing: more formal language

Say we run an SLR.

▶ The slope coefficient $\beta_1$ is an unknown population quantity, which we have estimated with data from a random sample of that population

▶ We can test hypotheses about this unknown population quantity based on the fact that the $T_{\widehat{\beta_1}}$ follows a t-distribution with $n - 2$ degrees of freedom

▶ With knowledge of the probability distribution of $T_{\widehat{\beta_1}}$ we can make probabilistic statements about the chances of observing any particular value of $T_{\widehat{\beta_1}}$ given a hypothesized value for the unknown parameter

▶ In particular, we are often interested in testing to see whether there is evidence to suggest that $\beta_1 \neq 0$ i.e. the slope coefficient is not zero i.e. there is evidence of a relationship between our dependent and independent variable

# The t-test steps

To test hypotheses about the value of $\beta_1$, we use a t-test (as the SE-standardized estimate follows a t-distribution). The steps of a t-test are:

1. State your null and alternative hypotheses about $\beta_1$

▶ The null hypothesis is denoted $H_0$
▶ The alternative hypothesis is denoted $H_1$
▶ e.g. $H_0 : \beta_1 = b$ and $H_1 : \beta_1 \neq b$

2. Choose the level of type-I error, $\alpha$, which gives the probability of rejecting the null hypothesis when it is actually true

▶ For example, $\alpha$ is most commonly chosen to be 0.05 i.e. the type-I error rate is 5%

# The t-test steps (ctd)

3. Compute the t-test statistic

$$
t_{\widehat{\beta}_1} = \frac{\left(\widehat{\beta}_1 - b\right)}{\text{se}\left(\widehat{\beta}_1\right)}
$$

4. Compute the p-value, which gives the probability of observing a test statistic as or even more extreme than $t_{\widehat{\beta}_1}$ under the assumption that the null hypothesis is true

5. Make a decision (reject the null if the p-value is less than $\alpha$, and fail to reject otherwise)
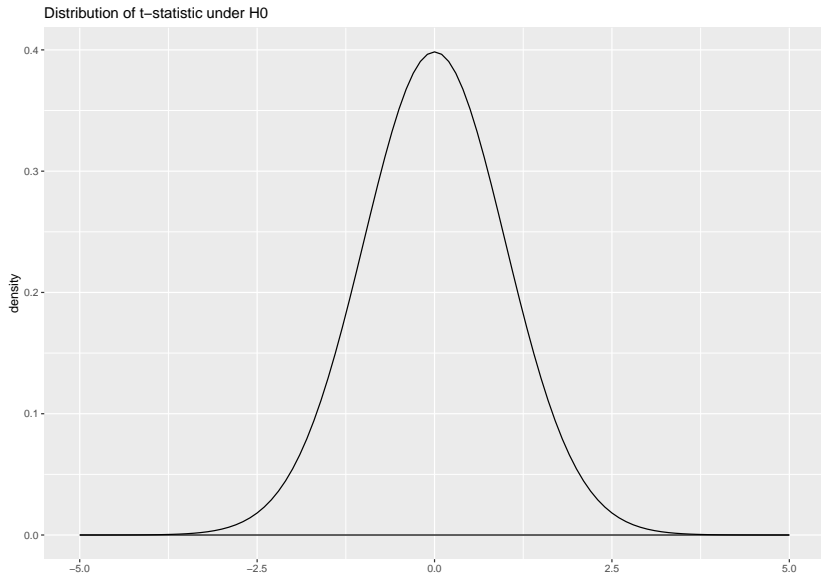
# The t-test in R

The `lm` summary put put shows the calculations for $t_{\widehat{\beta_1}}$ and corresponding p-value. Specifically these calculations test whether $H_0 : \beta_1 = 0$ and $H_1 : \beta_1 \neq 0$.

```
slr_mod <- lm(life_expectancy~tfr, data = country_ind_2017)
summary(slr_mod)
```

```
##
## Call:
## lm(formula = life_expectancy ~ tfr, data = country_ind_2017)
##
## Residuals:
##      Min      1Q   Median      3Q     Max
## -16.0718  -2.3864   0.3132   2.6537  11.3498
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  89.2394     0.7085  125.95   <2e-16 ***
## tfr          -5.3526     0.2326  -23.02   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.994 on 174 degrees of freedom
## Multiple R-squared:  0.7527, Adjusted R-squared:  0.7513
## F-statistic: 529.7 on 1 and 174 DF,  p-value: < 2.2e-16
```
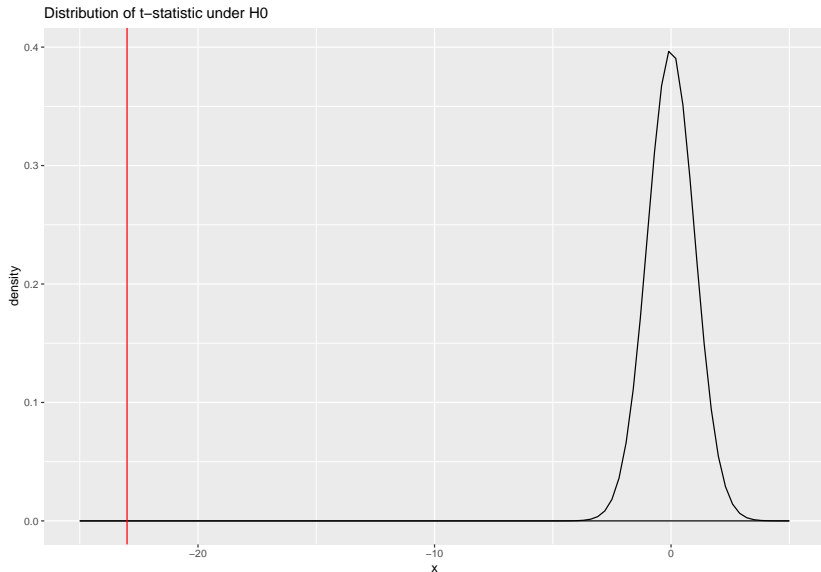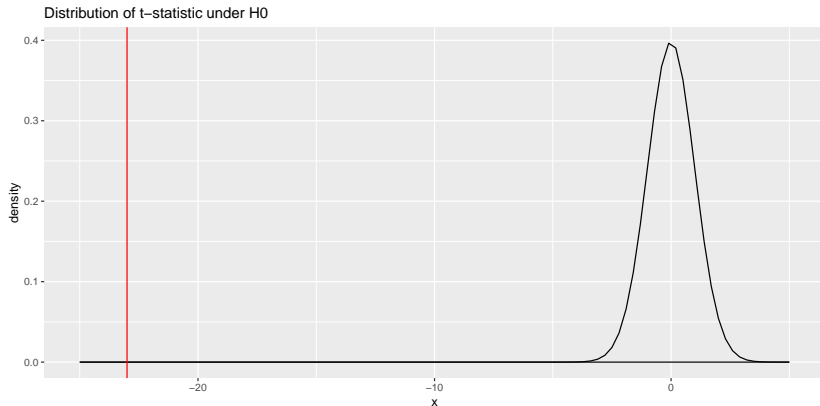
What should we conclude?

# Logic of the t-test



Distribution of t–statistic under H0

# Logic of the t-test

We calculated $t_{\widehat{\beta}_1} = -23$



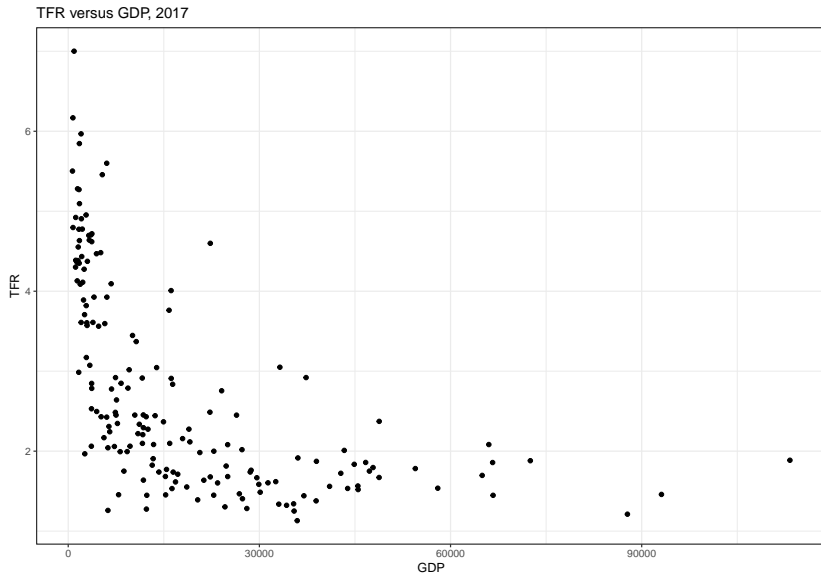Distribution of t−statistic under H0

# Logic of the t-test

- ▶ We calculated $t_{\widehat{\beta}_1} = -23$
- ▶ Under the null hypothesis, the probability of observing this value is very small—thus, we conclude the null hypothesis is likely false

Distribution of t–statistic under H0

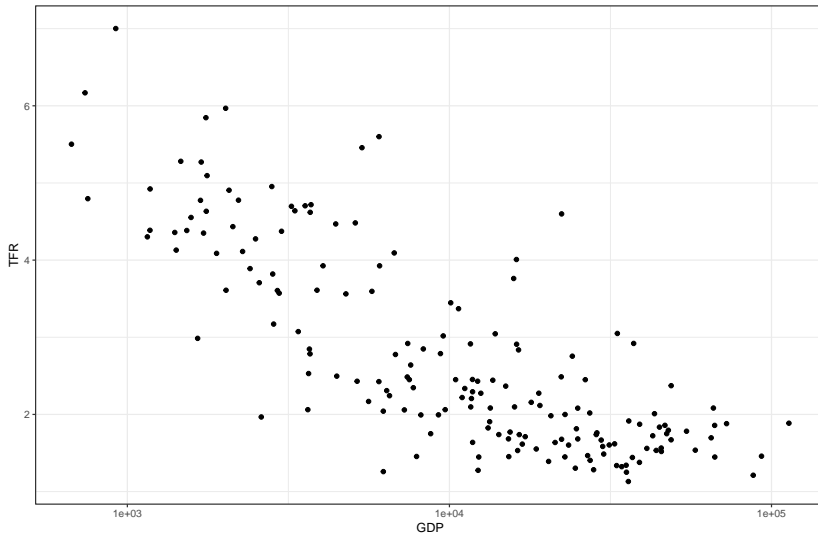# Regression with transformed variables

# Motivation



TFR versus GDP, 2017

# Motivation



TFR versus GDP, 2017
GDP plotted on log scale

# Variable transformations

- ▶ Sometimes we may want to allow for nonlinearities in our models
- ▶ A common way to deal with this is to perform a nonlinear transformation on one or more of the explanatory variables **AND/OR** on the response variable
- ▶ The interpretation of parameter estimates is less intuitive after transforming the explanatory variables and/or the response variable, although some transformations lend themselves to simple interpretations (i.e., the log transform)

# Log transforms

- By far the most common transformation is the natural log transform
- Either $\log Y$ or $\log X$ (or both)
- Luckily, the log transform has a meaningful coefficient interpretation

We will look at

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

# Log transforms: response variable

For response variables, when the model is

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

A one unit change in $X_i$ leads to a $(\exp(\beta_1) - 1)\,100$ percent change in $Y_i$, on average, holding other factors constant.

## Response variable: approximation

It turns out that $\exp(z) \approx 1 + z$ for small values of $z$.

So an approximate interpretation is

$$100\beta_1 \left(\Delta X_i\right) = \%\Delta Y_i$$

where $\Delta$ stands for "change".

▶ Thus, a one unit increase in $X_i$ is associated with a $100 \cdot \beta_1\%$ change in $Y_i$, on average, holding other factors constant

# Log transforms: expanatory variables

For explanatory variables, when the model is

$$Y_i = E(Y_i \mid \log X_i) + \varepsilon_i$$
$$= \beta_0 + \beta_1 \log X_i + \varepsilon_i$$

The interpretation is

$$\frac{\beta_1}{100} (\%\Delta X_i) = \Delta Y_i$$

where $\Delta$ stands for "change".

- Thus, a one percent (1%) increase in $X_k$ is associated with a $\frac{\beta_1}{100}$ unit change in $Y_i$, on average, holding other factors constant

# Log transforms: both variables

When both the response and explanatory variable is transformed, so the model is

$$\log Y_i = E\left(Y_i \mid \log X_i\right) + \varepsilon_i$$
$$= \beta_0 + \beta_1 \log X_i + \varepsilon_i$$

We are going to utilize the first approximation here, and say

The interpretation is

$$\beta_1\left(\%\Delta X_i\right) = \%\Delta Y_i$$

▶ Thus, a one percent (1%) increase in $X_1$ is associated with a $\beta_1$ % change in $Y_i$, on average, holding other factors constant

# Example

```
country_ind <- country_ind %>%
  mutate(log_tfr = log(tfr)) # log of GDP

summary(lm(log_tfr ~ gdp, data = country_ind))
```

```
##
## Call:
## lm(formula = log_tfr ~ gdp, data = country_ind)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -0.8901 -0.3104 -0.0297  0.3282  1.2304
##
## Coefficients:
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept)  1.189e+00  1.307e-02   90.93   <2e-16 ***
## gdp         -1.381e-05  5.081e-07  -27.17   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.3837 on 1582 degrees of freedom
## Multiple R-squared:  0.3182, Adjusted R-squared:  0.3178
## F-statistic: 738.4 on 1 and 1582 DF,  p-value: < 2.2e-16
```

- A 10^5 unit increase in GDP is associated with a 14% decrease in TFR

# Example

```
country_ind <- country_ind %>%
  mutate(log_gdp = log(gdp)) # log of GDP

summary(lm(tfr ~ log_gdp, data = country_ind))
```

```
##
## Call:
## lm(formula = tfr ~ log_gdp, data = country_ind)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -2.3279 -0.6231 -0.0883  0.4864  3.7178
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 11.39650    0.17194   66.28   <2e-16 ***
## log_gdp     -0.92940    0.01862  -49.91   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.8813 on 1582 degrees of freedom
## Multiple R-squared:  0.6116, Adjusted R-squared:  0.6114
## F-statistic:  2491 on 1 and 1582 DF,  p-value: < 2.2e-16
```

▶ A 1% increase in GDP is associated with a decrease of 0.93 children in TFR

# Example

```
summary(lm(log_tfr ~ log_gdp, data = country_ind))
```

```
##
## Call:
## lm(formula = log_tfr ~ log_gdp, data = country_ind)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.94552 -0.21097  0.00451  0.18577  1.15063
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  3.756075   0.056155   66.89   <2e-16 ***
## log_gdp     -0.306566   0.006081  -50.41   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.2878 on 1582 degrees of freedom
## Multiple R-squared:  0.6163, Adjusted R-squared:  0.6161
## F-statistic:  2541 on 1 and 1582 DF,  p-value: < 2.2e-16
```

▶ A 1% increase in GDP is associated with a 0.31% decrease in TFR

# Summary

- ▶ Often we may want to transform dependent or independent variables to make relationships more linear
- ▶ Log transforms are by far the most common
- ▶ This is because many variables are naturally log-normally distributed, e.g. income and GDP

Intro to multiple linear regression

# Motivation

- So far we have used regression to model the relationship between two variables
  - $Y_i$: dependent variable, response variable, outcome
  - $X_i$: independent variable, covariate, explanatory variable, predictor
- But for many problems, it's likely that the outcome of interest is associated with several different explanatory variables of interest
  - Time taken to build lego tower depends on number of blocks and number of distractions
  - Child's height depends on height of both parents
  - Income depends on age, education, occupation...
- We can extend the SLR set-up to include more that one independent variable/explanatory variable

# Multiple Linear Regression

How does the expected value (i.e. population mean) of a dependent variable differ across different levels of multiple independent variables?

We will go through how to estimate this model (with two independent variables)

- ▶ very similar to SLR process
- ▶ extends to more than 2 variables

# Example

- $Y_i$ is the dependent variable or response variable
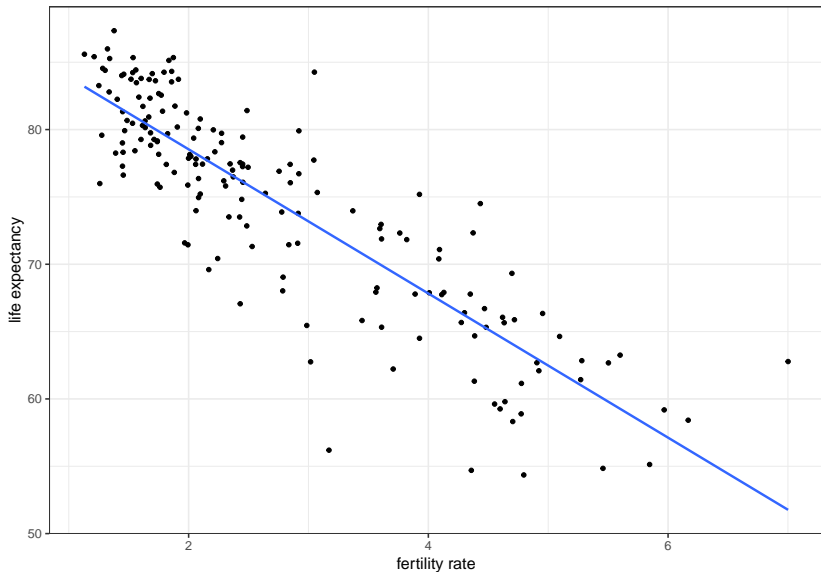- $X_{i1}$ and $X_{i2}$ are the independent variables, explanatory variables or predictors

Example:

- $\{Y_1, Y_2, \ldots, Y_{176}\}$ is life expectancy by country in 2017
- $\{X_{1,1}, X_{2,1}, \ldots, X_{176,1}\}$ is TFR by country in 2017
- $\{X_{1,2}, X_{2,2}, \ldots, X_{176,2}\}$ is child mortality by country in 2017
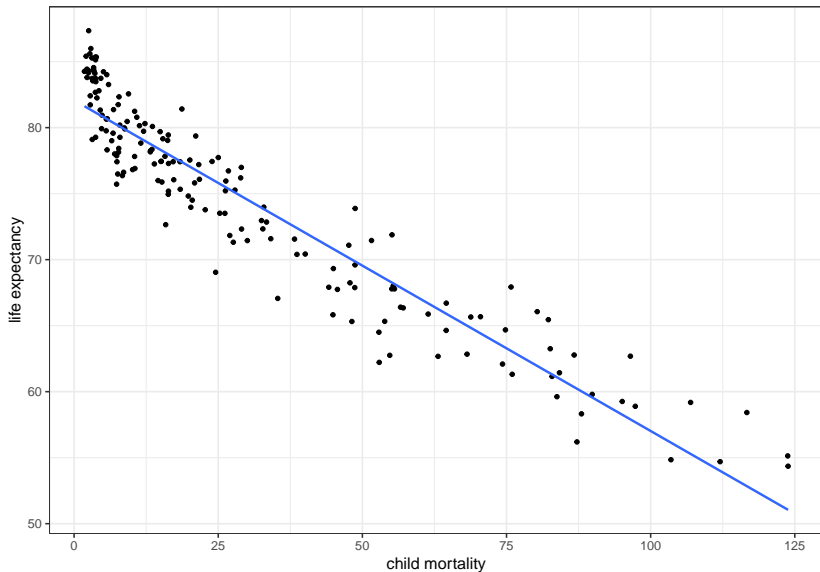
Research question:

- How does life expectancy differ across different levels of fertility and child mortality
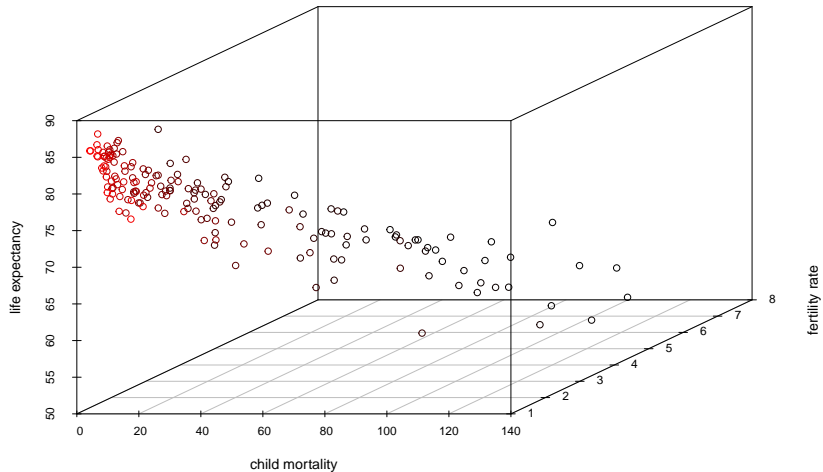- In other words, is life expectancy associated with fertility and child mortality, and if so, how?

# Scatter plot of fertility and life expectancy

# Scatter plot of child mortality and life expectancy

# Scatter plot of all variables

# MLR model

With two covariates, the MLR model is

$$Y_i = E(Y_i \mid X_{i1}, X_{i2}) + \varepsilon_i$$
$$= \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \varepsilon_i$$

Specifically, the most basic MLR model is a simple linear function of $X_{i1}$ and $X_{i2}$, and three parameters, $\beta_0$, $\beta_1$ and $\beta_2$.

# Interpretation

The MLR model: $E\left(Y_i \mid X_{i1}, X_{i2}\right) = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2}$

- What is $\beta_0$?

$$E\left(Y_i \mid X_{i1} = 0, X_{i2} = 0\right) = \beta_0 + \beta_1(0) + \beta_2(0)$$
$$= \beta_0$$

- $\beta_0$ is the is the expected value, or population mean, of $Y_i$ given both $X_{i1}$ and $X_{i2}$ equal zero.

# Interpretation

The MLR model: $E\left(Y_i \mid X_{i1}, X_{i2}\right) = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2}$

▶ What is $\beta_1$?

$$
\begin{aligned}
E\left(Y_i \mid X_{i1} = x_1 + 1, X_{i2} = x_2\right) &- E\left(Y_i \mid X_{i1} = x_1, X_{i2} = x_2\right) \\
&= \left(\beta_0 + \beta_1\left(x_1 + 1\right) + \beta_2 x_2\right) - \left(\beta_0 + \beta_1 x_1 + \beta_2 x_2\right) \\
&= \left(\beta_0 + \beta_1 x_1 + \beta_1 + \beta_2 x_2\right) - \left(\beta_0 + \beta_1 x_1 + \beta_2 x_2\right) \\
&= \beta_1
\end{aligned}
$$

▶ $\beta_1$ is the change in the expected value, or population mean, of $Y_i$ associated with a one unit increase in $X_{i1}$, **holding $X_{i2}$ constant at any value**

Same idea for $\beta_2$.

# Interpretation

The MLR model: $E\left(Y_i \mid X_{i1}, X_{i2}\right) = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2}$

- In general $\beta_1\left(x_1^* - x_1\right)$ is the change in the expected value of $Y_i$ associated with a $\left(x_1^* - x_1\right)$ change in $X_{i1}$, holding $X_{i2}$ constant
- $\beta_2\left(x_2^* - x_2\right)$ is the change in the expected value of $Y_i$ associated with a $\left(x_2^* - x_2\right)$ change in $X_{i2}$, holding $X_{i1}$ constant

## OLS Estimation

Back to example

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \varepsilon_i$$

where $Y_i$ is life expectancy, $X_{i1}$ is fertility rate and $X_{i2}$ is child mortality.

The estimates are

$$\hat{E}(Y_i \mid X_{i1}, X_{i2}) = \hat{\beta}_0 + \hat{\beta}_1 X_{i1} + \hat{\beta}_2 X_{i2}$$
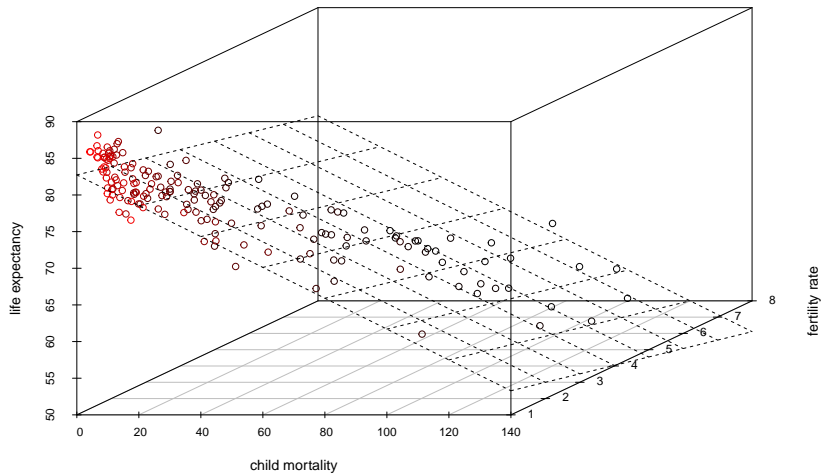$$= 83.8 - 1.07 X_{i1} - 0.21 X_{i2}$$

How should we interpret?
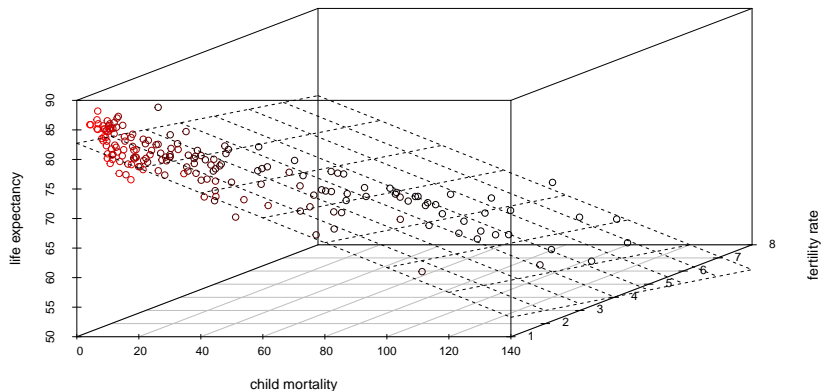
# Interpretation

The estimates are

$$\hat{E}\left(Y_i \mid X_{i1}, X_{i2}\right) = \widehat{\beta}_0 + \widehat{\beta}_1 X_{i1} + \widehat{\beta}_2 X_{i2}$$
$$= 83.8 - 1.07 X_{i1} - 0.21 X_{i2}$$

# OLS Estimation

# OLS Interpretation

For example, Japan:

$$Y_i = \hat{E}\left(Y_i \mid X_{i1} = 1.4, X_{i2} = 2.5\right) + \varepsilon_i$$
$$= 83.8 - 1.07 \times 1.4 - 0.21 \times 2.5 + 5.5$$
$$= 81.8 + 5.5$$

# MLR in R

▶ Can estimate MLR exactly the same way as SLR, just add additional variables with a + in the formula in `lm`

▶ Residuals, fitted values etc extracted in the same way

```
mod <- lm(life_expectancy~tfr+child_mort, data = country_ind_2017)
summary(mod)
```

```
##
## Call:
## lm(formula = life_expectancy ~ tfr + child_mort, data = country_ind_2017)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -6.7103 -1.6787 -0.1197  1.7379  5.7605
##
## Coefficients:
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept) 83.80622    0.56028 149.578  < 2e-16 ***
## tfr         -1.07102    0.30293  -3.536 0.000522 ***
## child_mort  -0.21031    0.01301 -16.171  < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.527 on 173 degrees of freedom
## Multiple R-squared:  0.9015, Adjusted R-squared:  0.9004
## F-statistic: 792.1 on 2 and 173 DF,  p-value: < 2.2e-16
```

# Next week

- Estimation
- Categorical covariates
- Hypothesis testing
- Interaction