# Methods in Psycholinguistics
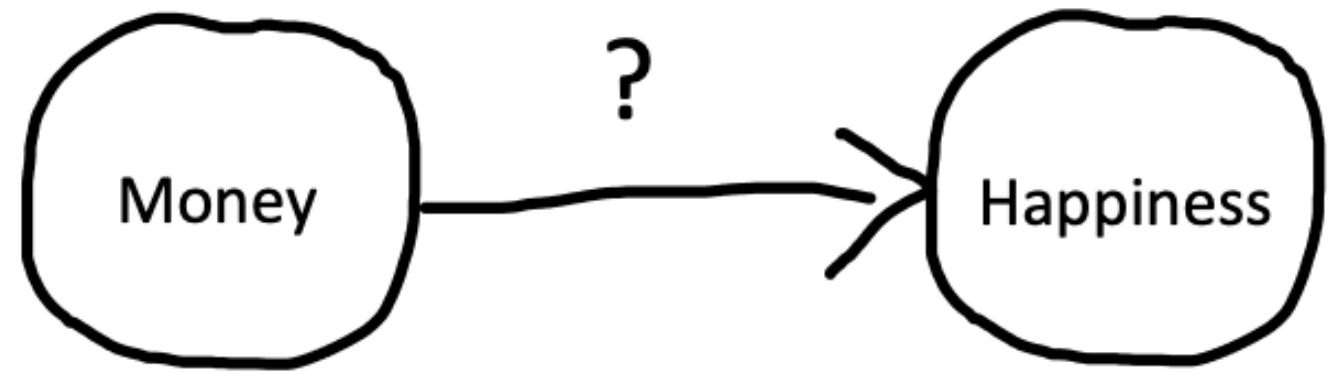# — Estimation, inference —
# — Linear regression —

Judith Degen
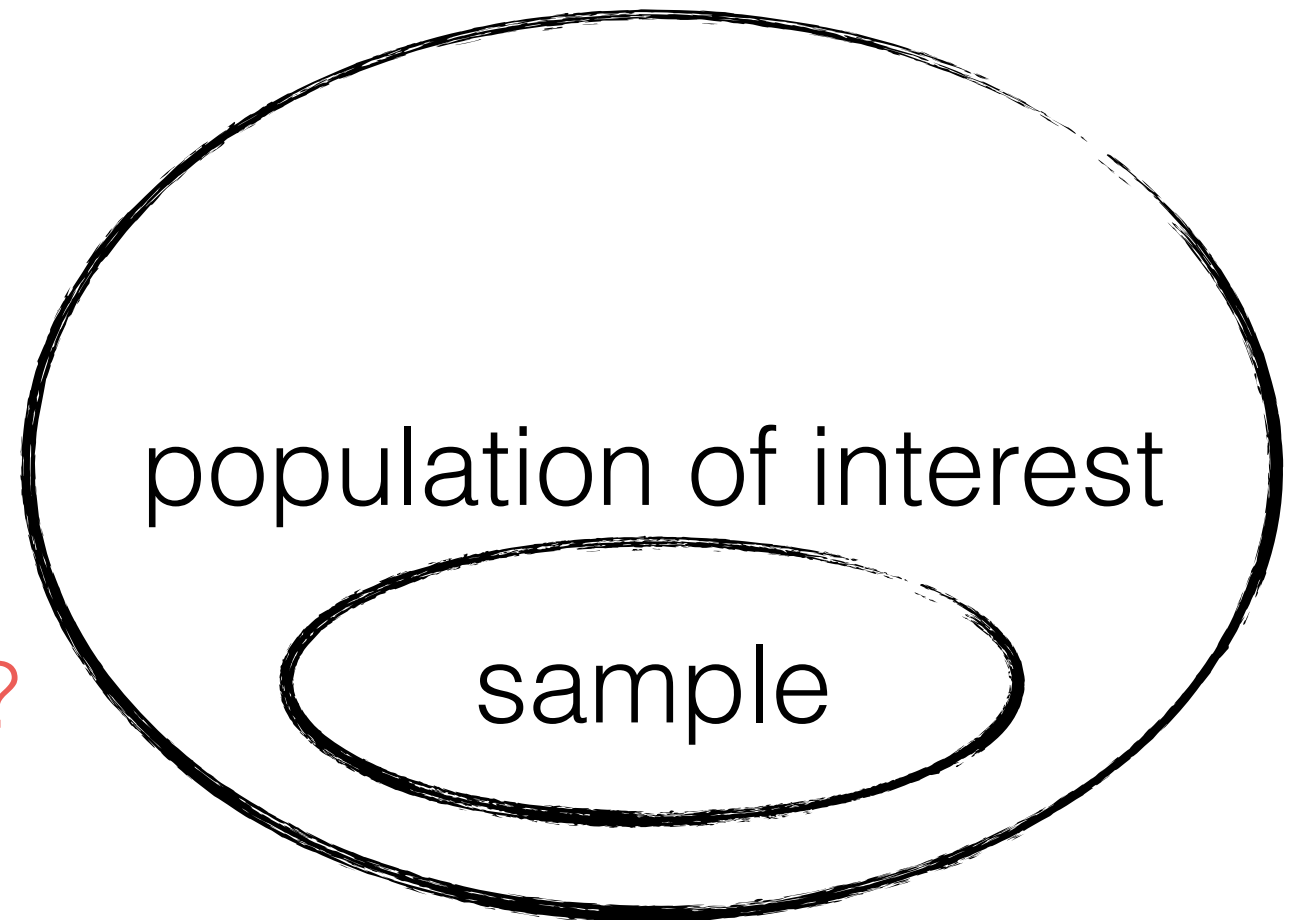
# Why do stats?

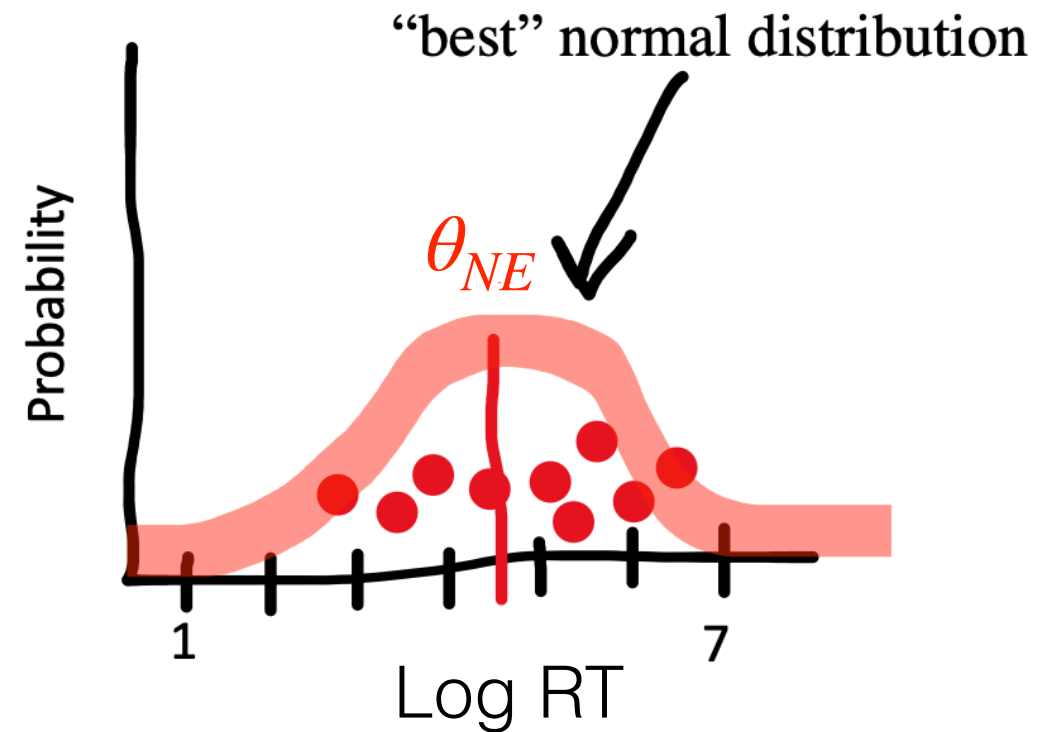# Estimation

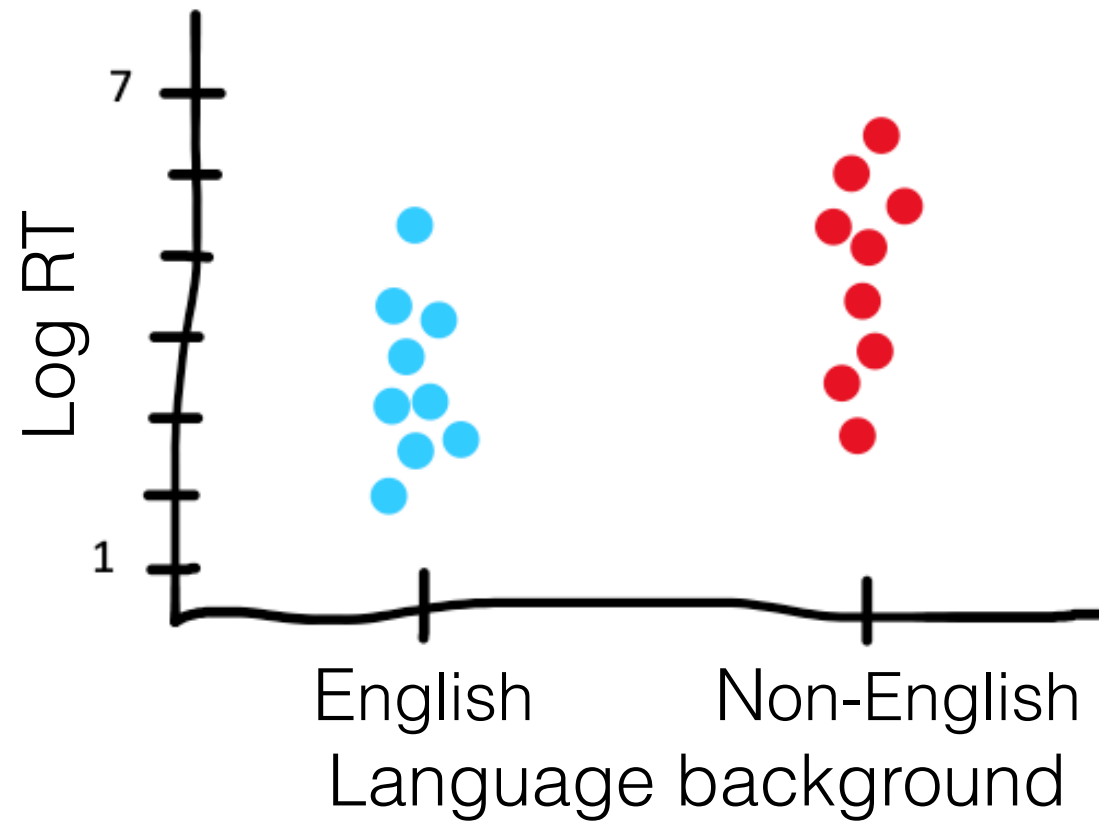How big is the effect of X on Y?



# Inference

How likely is a difference in Y between $x_1$ and $x_2$ to generalize to the population?

# Estimation



$\theta_{NE}$ is the mean of the best-fitting normal distribution

# Methods of estimation



Maximum Likelihood Estimation (MLE)

likelihood

$$p(data|\theta_M)$$

Bayesian Estimation

# Estimating treatment (condition) effect



MLE and Bayesian estimation return similar results for $\hat{\beta}$ with
- large amounts of data
- weak prior beliefs

# Inference

|  | **Frequentist** | **Bayesian** |
|---|---|---|
| **Hypothesis testing** | *p* value from null hypothesis significance test | Bayes factor |
| **Estimation with uncertainty** | estimate with confidence interval | posterior distribution with credible interval |

Null Hypothesis Significance Testing (NHST) — still very much the norm

Bayesian stats favor estimation mindset, but hypothesis-testing also possible with Bayes factors

# Bayes Factors

| BF range | Interpretation |
|----------|----------------|
| $< 1$ | Negative evidence (supports H0) |
| $1 - 3$ | Barely worth mentioning |
| $3 - 10$ | Substantial |
| $10 - 30$ | Strong |
| $30 - 100$ | Very strong |
| $> 100$ | Decisive |

Likelihood of data under hypothesis of non-zero difference

$$BF = \frac{p(data|H_1)}{p(data|H_0)}$$

Likelihood of data under null hypothesis of zero difference

# Linear regression

# What will we cover?

- introduction to Generalized Linear Models (GLMs) and Generalized Linear Mixed Models (GLMMs)
  - mathematical background
  - intuition / conceptualization
  - geometric interpretation
  - common issues & solutions for GLM/GLMMs

# What kind of data can you analyze with GLMs?

- continuous (nominal)
  response/reading times, slider ratings, speech onset times,…
  …………linear regression

- categorical (binary)
  truth value judgments,
  any binary choice prediction…
  ……….logistic regression

- ordered discrete (ordinal)
  Likert scale ratings…
  ……….ordinal regression

- unordered discrete
  any choice between more than two options
  …..multinomial regression

# Generalized Linear Models

Goal: model effects of predictors (independent variables) X on a response (dependent variable) Y

# Reviewing GLMs

Assumptions of the generalized linear model:

1. Predictors $X_i$ influence $Y$ through the mediation of a linear predictor $\eta$

2. $\eta$ is a linear combination of the $X_i$

$$\eta = \alpha + \beta_1 X_1 + \cdots + \beta_N X_N$$

3. $\eta$ determines predicted mean $\mu$ of $Y$

$$\eta = g(\mu) \qquad \text{(link function)}$$

4. There is some noise distribution $P$ around the predicted mean $\mu$ of $Y$:

$$P(Y = y; \mu)$$

# Linear regression

Linear regression is a kind of generalized linear model.

The predicted mean is simply the linear predictor:

$$\eta = l(\mu) = \mu$$

Noise is normally (=Gaussian) distributed around 0 with standard deviation $\sigma$ :

$$\epsilon \sim N(0, \sigma)$$

This results in the traditional linear regression equation:

Predicted mean $\mu = \eta$   Noise $\sim N(0, \sigma)$

$$Y = \alpha + \beta_1 X_1 + \cdots + \beta_n X_n + \epsilon$$

# An example: lexical decision

tpozt          *Word or non-word?*

house          *Word or non-word?*

Measure response times (RT)

Question: which factors predict RTs?

# The dataset

- lexical decisions from 79 concrete nouns, each seen by 21 participants  (1,659 observations)

- **Outcome/response:** log-transformed lexical decision times

- **Inputs:**

  - continuous: e.g. frequency

  - categorical: e.g., native language (English vs other)

# The basic model

Let's assume that frequency has a *linear* effect on average log RT, and trial-level noise is *normally distributed.*

If $x_i$ is frequency, this simple model is:

Given

Noise $\sim N(0, \sigma_\epsilon)$

$$RT_{ij} = \alpha + \beta x_{ij} + \overbrace{\epsilon_{ij}}$$

Inferences

E.g. "Does frequency affect RT?"—> is $\beta$ reliably non-zero?

Let's translate this into R

```
Coefficients:
                Estimate Std. Error t value Pr(>|t|)
(Intercept)    6.588778    0.022296 295.515   <2e-16 ***
Frequency     -0.042872    0.004533  -9.459   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' '

Residual standard error: 0.2353 on 1657 degrees of freedom
Multiple R-squared:  0.05123,   Adjusted R-squared:  0.05066
```

$$\text{Noise} \sim N(0, \sigma_\epsilon)$$

$$RT_{ij} = \boxed{\alpha} + \boxed{\beta} x_{ij} + \overbrace{\epsilon_{ij}}$$

*"There was a significant main effect of frequency such that more frequent words were responded to more quickly $(\beta = -0.04, SE = 0.004, t = -9.46, p < .0001)$."*

Why is $R^2$ so low even though frequency has tiny p-value?

```
Coefficients:
                              Estimate Std. Error t value Pr(>|t|)
(Intercept)                   6.318309   0.007435  849.78  <2e-16 ***
LanguageBackgroundNon-English 0.155821   0.011358   13.72  <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.2289 on 1657 degrees of freedom
Multiple R-squared:  0.102,    Adjusted R-squared:  0.1015
```
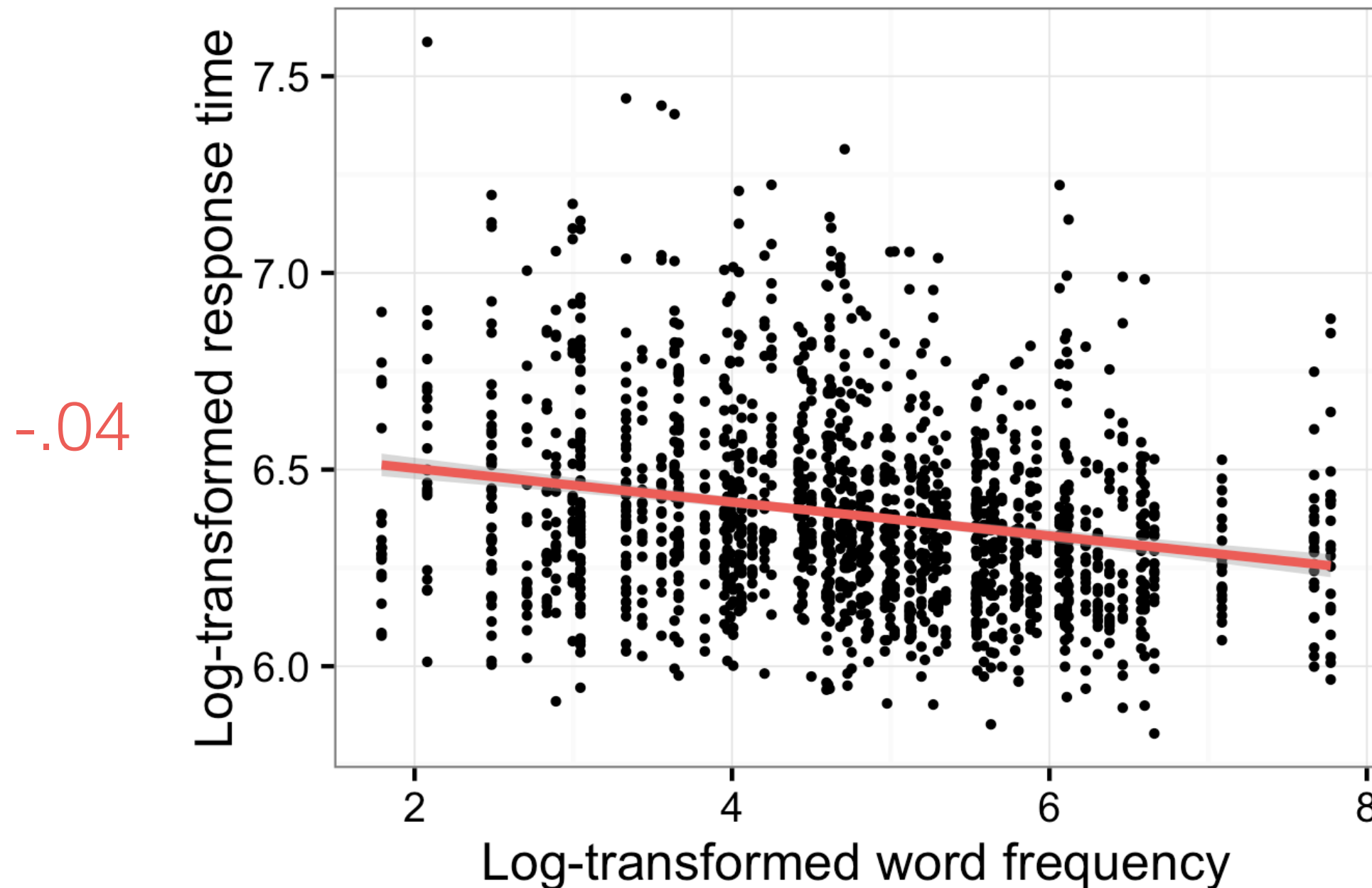
$$\text{Noise} \sim N(0, \sigma_\epsilon)$$

$$RT_{ij} = \boxed{\alpha} + \boxed{\beta} x_{ij} + \overbrace{\epsilon_{ij}}$$

*"There was a significant main effect of language background such that participants with a Non-English background responded more slowly ($\beta = 0.16, SE = 0.01, t = 13.72, p < .0001$)."*

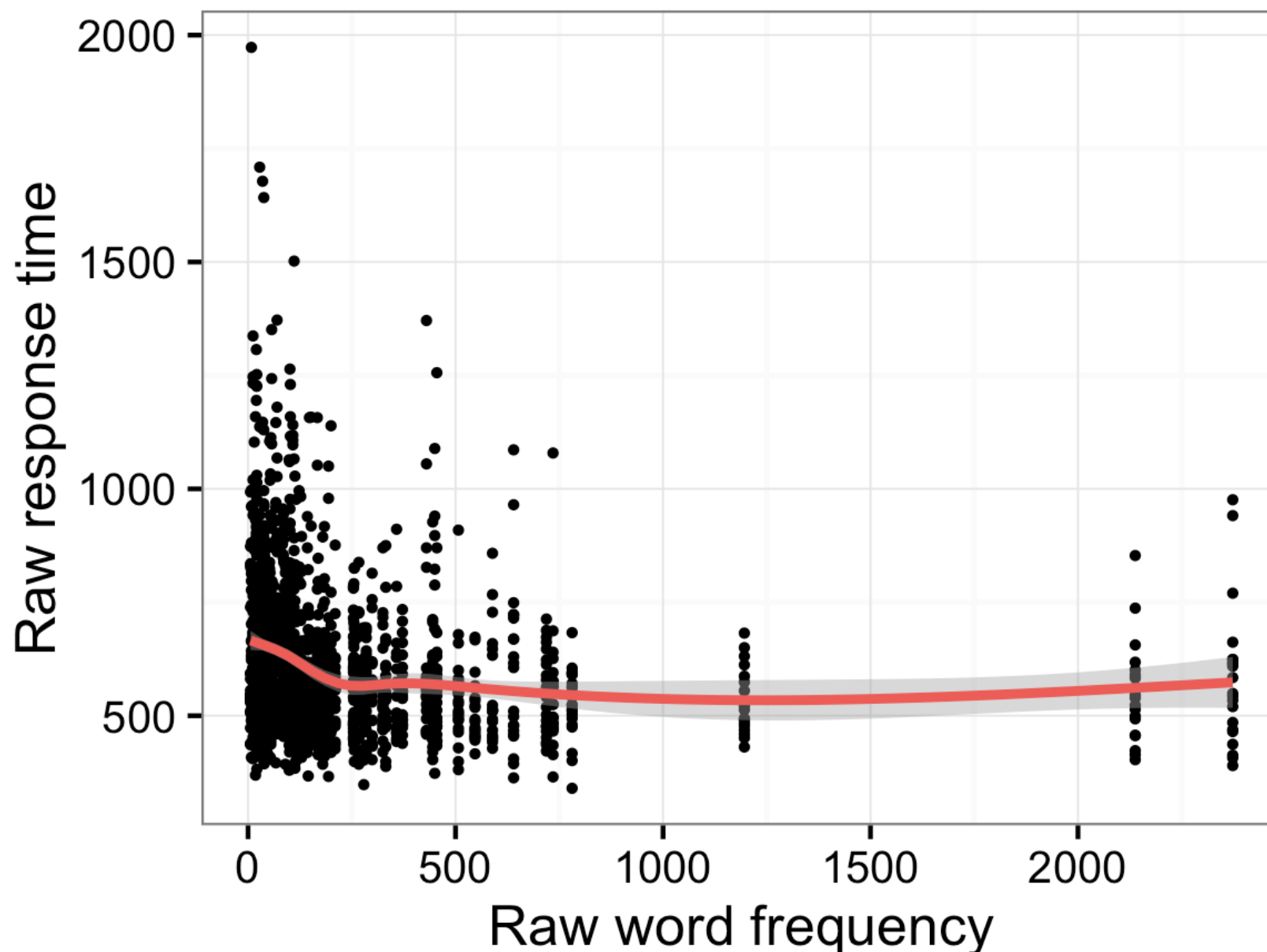Why is $R^2$ so low even though frequency has tiny p-value?

# Geometric intuitions



-.04

Geometric interpretation of linear regression: find slopes for predictors that minimize squared error

# Linearity assumption

Like ANOVA, the linear model assumes the outcome is linear in the *coefficients* (**linearity assumption**).



This doesn't mean that outcome and input *variables* need to be linearly related!

# Adding predictors (multiple regression)

Extend the simple model to include an additional predictor for **morphological family size** (number of words in the morphological family of the target word).

```
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  6.563853   0.026826 244.685  < 2e-16 ***
Frequency   -0.035310   0.006407  -5.511 4.13e-08 ***
FamilySize  -0.015655   0.009380  -1.669   0.0953 .
```

1. Is the interpretation of the output clear?
2. What is the interpretation of the intercept?
3. How much faster is the most frequent word expected to be read compared to the least frequent word?

# Categorical predictors

Extend the model to include a predictor for participants' **native language** (English vs other).

```
Coefficients:
                      Estimate Std. Error t value Pr(>|t|)
(Intercept)           6.497073   0.025784 251.977  < 2e-16 ***
Frequency            -0.035310   0.006054  -5.832 6.56e-09 ***
FamilySize           -0.015655   0.008863  -1.766   0.0775 .
NativeLanguageOther   0.155821   0.011025  14.133  < 2e-16 ***
```

The output is a linear combination of predictors, so categorical predictors need to be coded numerically —> Default in R: dummy/treatment coding (more later)

What is the "mean" that is being predicted in this model?

# Interactions

Interactions are products of predictors.
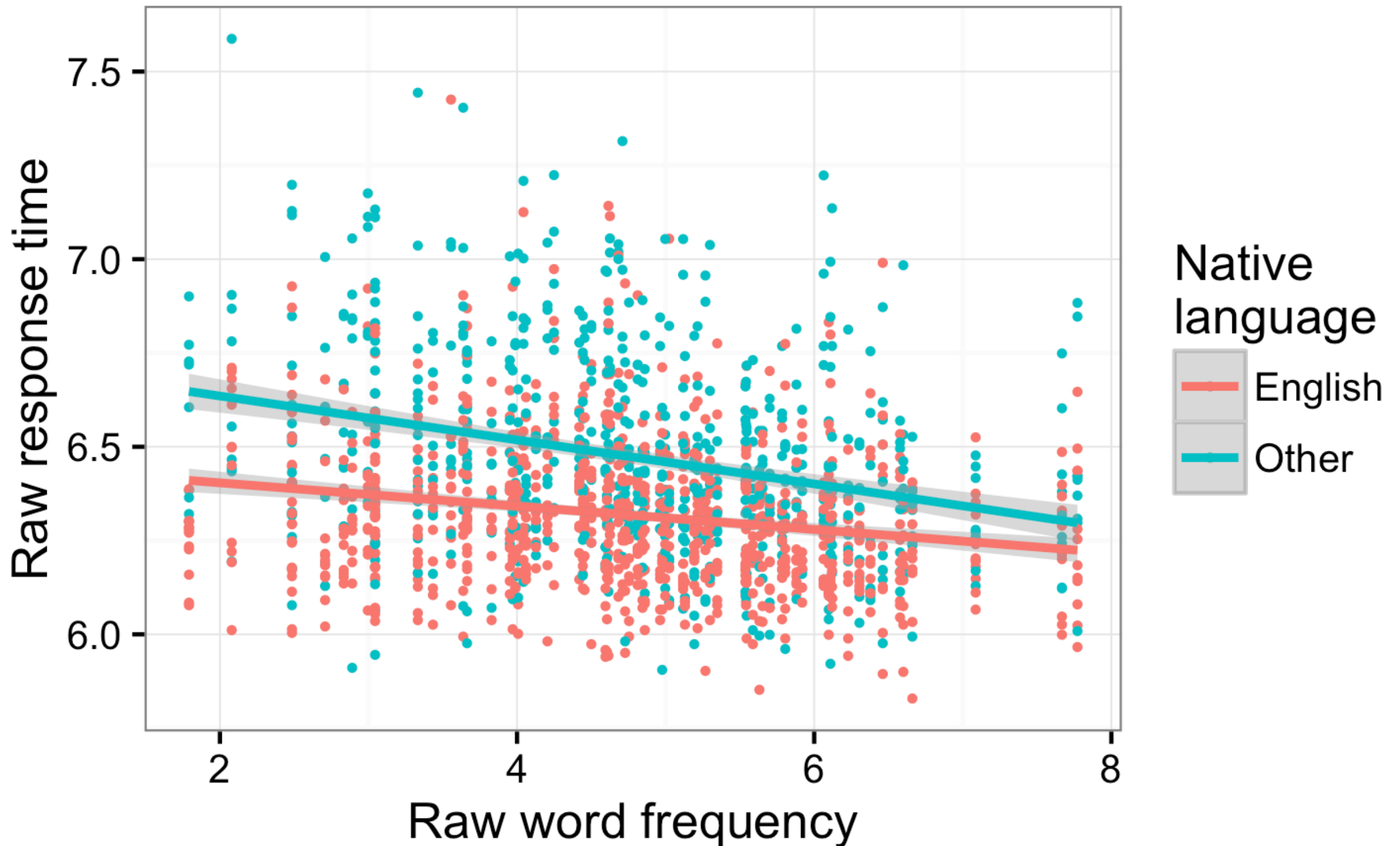Interpretation of significant interactions: the slope of one predictor differs for different values of the other predictor.

```
Coefficients:
                                Estimate Std. Error t value Pr(>|t|)
(Intercept)                     6.441135   0.031140 206.847  < 2e-16 ***
FamilySize                     -0.015655   0.008839  -1.771 0.076726 .
Frequency                      -0.023536   0.007079  -3.325 0.000905 ***
NativeLanguageOther             0.286343   0.042432   6.748 2.06e-11 ***
Frequency:NativeLanguageOther  -0.027472   0.008626  -3.185 0.001475 **
```

How should we interpret the interaction between frequency and native language?

# Plotting the interaction

# Determining parameters

How do we choose parameters (model coefficients) $\beta_i$ and $\sigma$?

**Find the best ones.** (see <u>Andrew Ng's videos</u>)

Two major approaches:
1. Maximum Likelihood Estimation (ML): pick parameter values that maximize the (log) probability of data, i.e., maximize $P(Y|\beta_i, \sigma)$
2. Bayesian inference: infer best model parameters via Bayes' rule, given a prior distribution over model parameters

$$P(\beta_i, \sigma|Y) = \frac{\overbrace{P(Y|\beta_i, \sigma)}^{\text{Likelihood}} \cdot \overbrace{P(\beta_i, \sigma)}^{\text{Prior}}}{P(Y)}$$
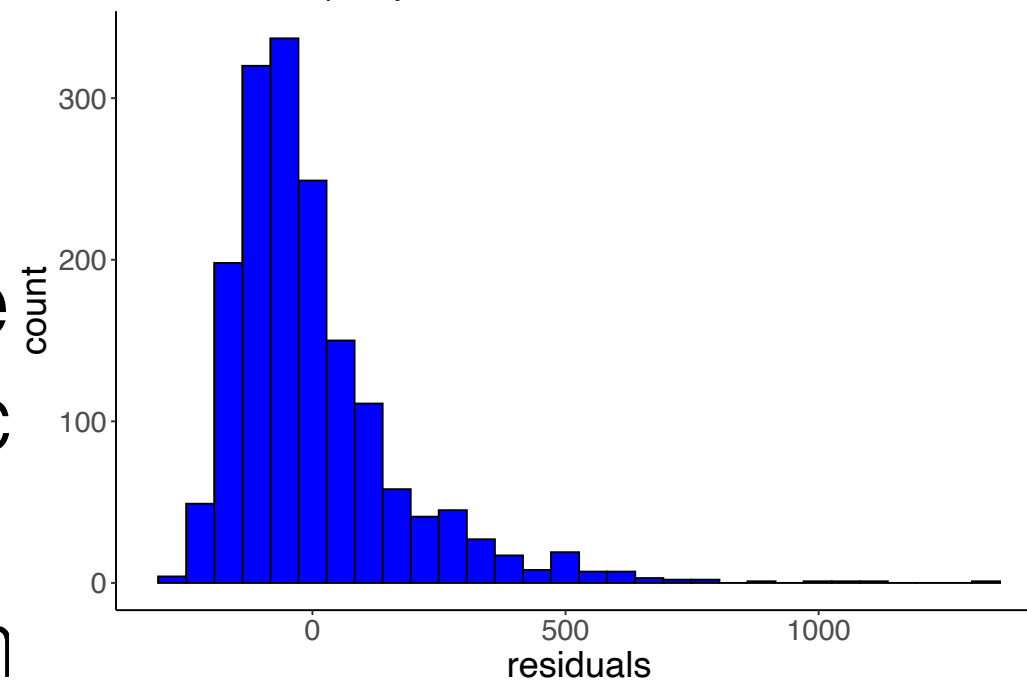
# Hypothesis testing in psycholinguistic research

- often, we make predictions not just about the **existence**, but also about the **direction** of the effect

- sometimes, we're also interested in effect **shapes** (e.g., non-linearities)

- unlike ANOVA, regression analyses test hypotheses about effect **direction**, **shape**, and **size** without requiring post-hoc analyses

  - if predictors are coded appropriately (more later)

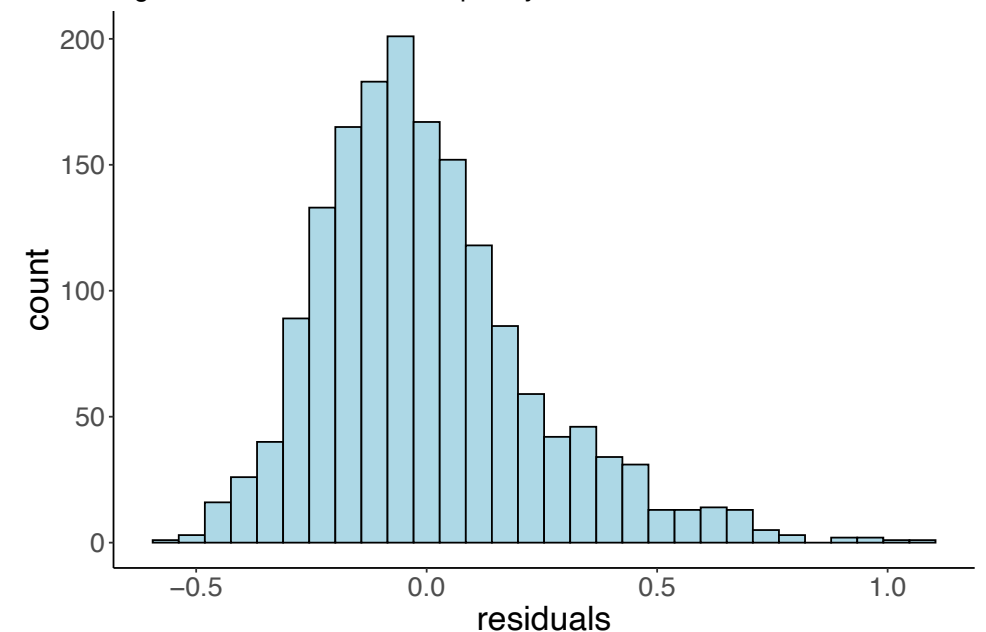  - if the model can be trusted (more later)

# Assumptions of linear regression

- homoscedasticity of residuals: e
same across all values of predic

- normality of residuals: error term
distributed

- independence of samples

**Distribution of residuals**
raw RT and Frequency

**Distribution of residuals**
log−transformed RT and Frequency

# Assumptions of linear regression

- homoscedasticity of residuals: error term is the same across all values of predictor

- normality of residuals: error term is normally distributed

- **independence of samples**

$$\sigma = \sqrt{\frac{\sum (x_i - \mu)^2}{N}}$$

**Standard deviation** of the population: average amount of variability in the data

$$\sigma_{\bar{x}} = \frac{\sigma_x}{\sqrt{n}}$$

**Standard error** of a sampling distribution: estimate of population standard deviation