

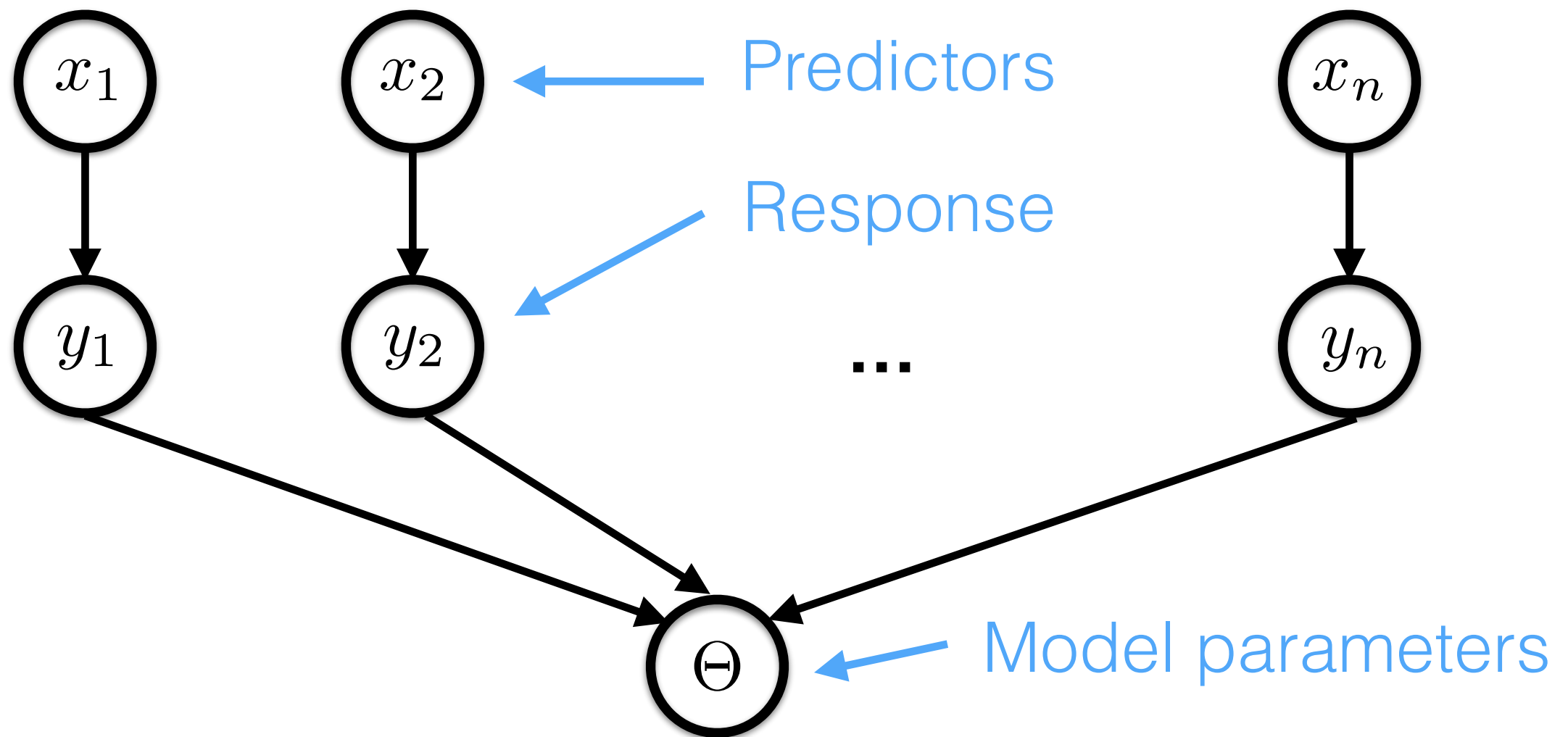
# Rutgers Linguistics Workshop on Mixed Effects Models — Linear regression —

09/16/2016

Judith Degen  
Stanford University

# 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  $X_i$

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

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

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

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

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

# Linear regression

Linear regression (which underlies ANOVA) 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:

$$Y = \text{Predicted mean } \mu = \eta \quad \text{Noise } \sim N(0, \sigma)$$
$$Y = \alpha + \beta_1 X_1 + \cdots + \beta_n X_n + \epsilon$$

# An example: lexical decision

Baayen, Feldman, & Schreuder (2006)

tpozt

*Word or non-word?*

house

*Word or non-word?*

Measure response times (RT)

Question: which factors predict RTs?

Let's analyze...  
open RStudio!

# 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

$$RT_{ij} = \alpha + \beta x_{ij} + \underbrace{\epsilon_{ij}}_{\text{Noise} \sim N(0, \sigma_\epsilon)}$$

Inferences

The diagram illustrates the flow of information in a linear model. The word "Given" is at the top, with three green arrows pointing down to  $RT_{ij}$  (green box),  $\alpha$  (blue box), and  $x_{ij}$  (green box). The word "Inferences" is at the bottom, with three arrows pointing up: a blue arrow to  $\alpha$ , a red arrow to  $\beta$  (red box), and a purple arrow to  $\sigma_\epsilon$  (purple box). The equation  $RT_{ij} = \alpha + \beta x_{ij} + \epsilon_{ij}$  is centered, with  $\epsilon_{ij}$  underlined and labeled "Noise  $\sim N(0, \sigma_\epsilon)$ ".

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	0.0222296	295.515	<2e-16 ***
Frequency	-0.042872	0.004533	-9.459	<2e-16 ***

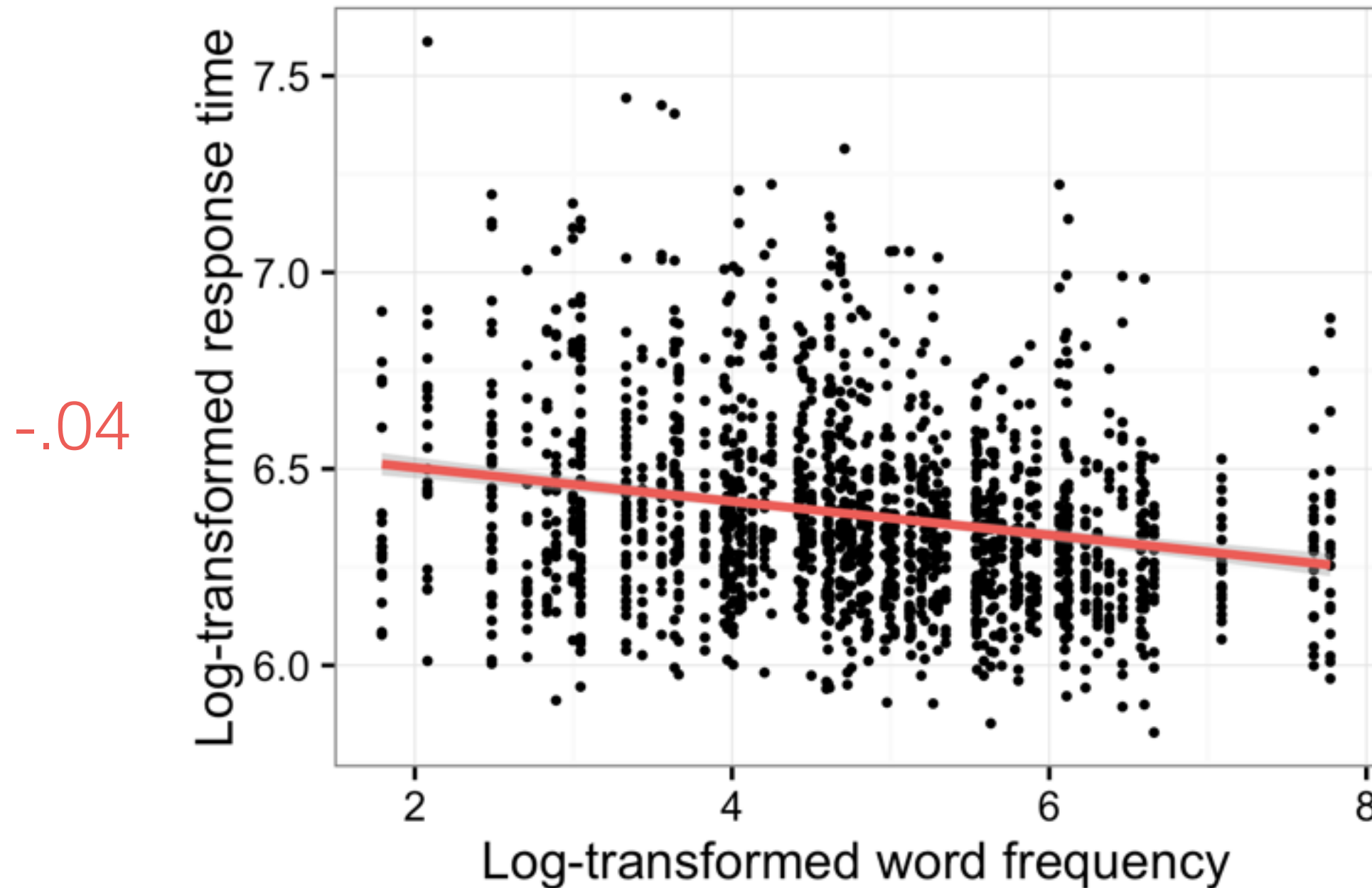
---  
Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1  
Residual standard error: 0.2353 on 1657 degrees of freedom  
Multiple R-squared: 0.05123, Adjusted R-squared: 0.05066

$$RT_{ij} = \alpha + \beta x_{ij} + \underbrace{\epsilon_{ij}}_{\text{Noise} \sim N(0, \sigma_\epsilon)}$$

*“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?

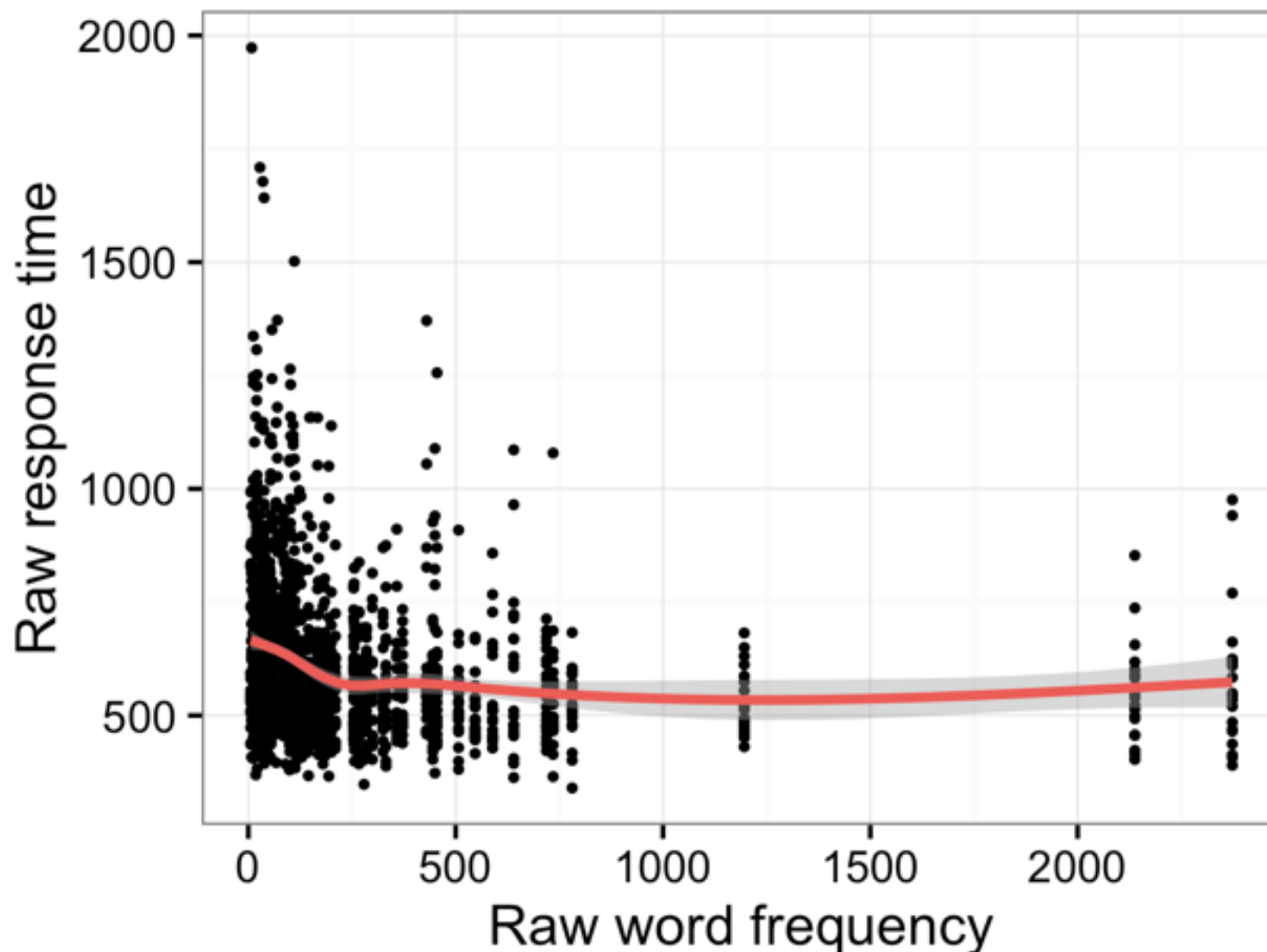
# Geometric intuitions



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 tomorrow)

What is the “mean” that is being predicted in this model?



# Interactions

Interactions are products of predictors.

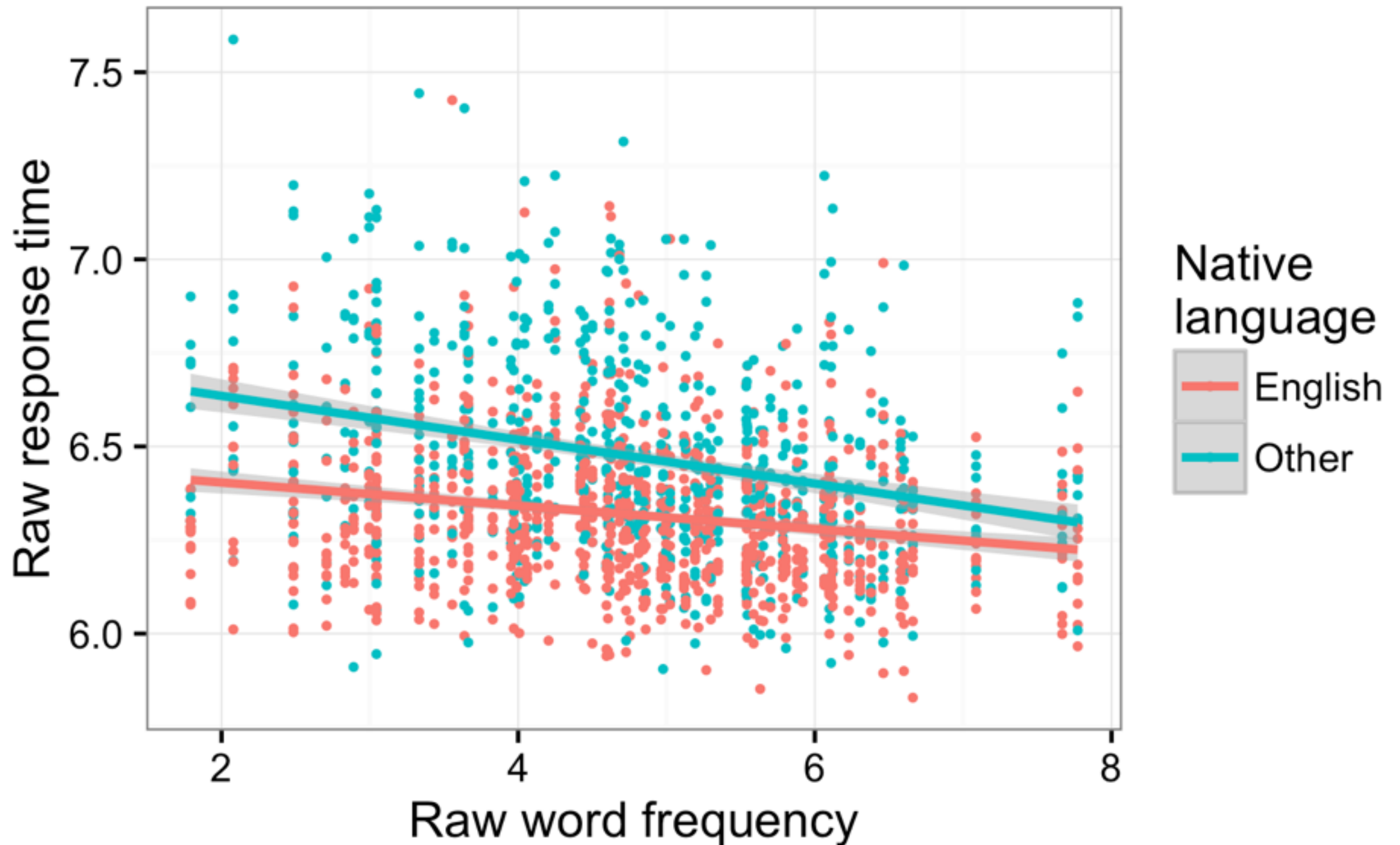
Interpretation of significant interactions: the slope of one predictor differs for different values of the other predictor.

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





# Linear regression vs. ANOVA

- **shared**

- linearity assumption (though investigation of non-linearities easily possible in regression)
- assumption of normality
- assumption of independence (of noise)
- ANOVA is basically linear regression with only categorical predictors

- **different**

- Generalized Linear Model
- consistent and transparent way of treating continuous and categorical predictors
- regression encourages a priori explicit coding of hypothesis (reducing post-hoc tests)

# 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 (see tomorrow)
  - if the model can be trusted (see tomorrow)

# Determining parameters

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

**Find the best ones.** (cf Andrew Ng's videos)

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

$$P(\beta_i, \sigma|Y) = \frac{\overbrace{P(Y|\beta_i, \sigma)}^{\text{Likelihood}} \cdot \overbrace{P(\beta_i, \sigma)}^{\text{Prior}}}{P(Y)}$$