

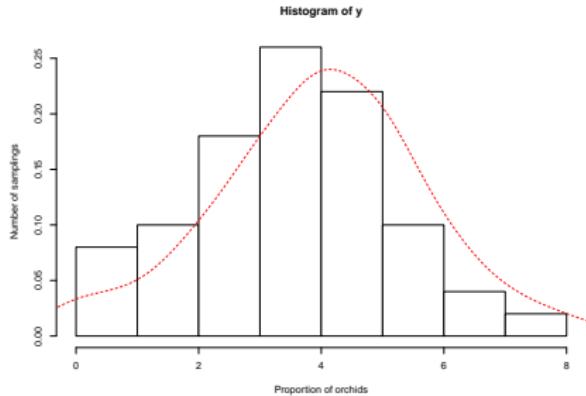
# Introduction to Linear Models

Bert van der Veen

Department of Mathematical Sciences, NTNU

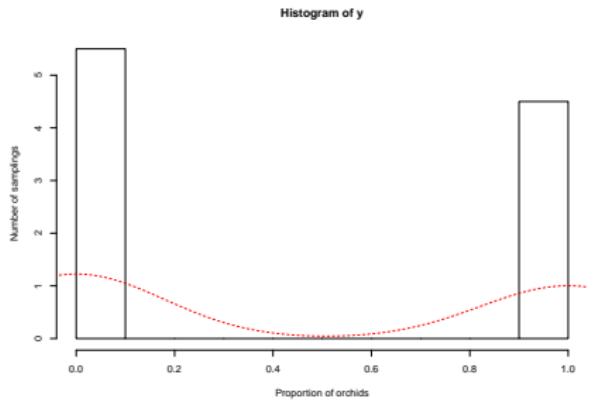
## The orchids example

```
set.seed(12345) # For reproducibility
n.times <- 50;n.picks=10;p.orchid <- 0.4
y <- rbinom(n.times, size = n.picks, prob = p.orchid)
```



## The orchids example

```
set.seed(12345) # For reproducibility
n.times <- 100;n.picks=1;p.orchid <- 0.4
y <- rbinom(n.times, size = n.picks, prob = p.orchid)
```



## The orchids example

The Binomial distribution is approximately normal for large  $n_{picks}$  and  $\pi$  away from 0 and 1

## Normality

**Real ecological observations  
are rarely normally distributed.**

But it is a nice starting point when learning GLMs.

We can use it (e.g.,) when the mean is far enough from zero.



# The normal distribution

---

$$f(y_i; \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{(y_i - \mu)^2}{2\sigma^2}\right\} \quad (1)$$

## The normal distribution (2)

Likelihood:

$$f(y_i; \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{(y_i - \mu)^2}{2\sigma^2}\right\} \quad (2)$$

log-Likelihood:

$$\log\{f(y_i; \mu, \sigma)\} = -\frac{1}{2} \log(\sigma^2 2\pi) - \frac{(y_i - \mu)^2}{2\sigma^2} \quad (3)$$

Two parameters:  $\mu$  and  $\sigma$

- ▶  $\mu$  is the mean; the middle of the distribution
- ▶  $\sigma$  is the standard deviation; it controls the width

Not only used for data, also the basis of many statistics (e.g., asymptotic sampling distributions)

## Estimating $\mu$

---

Same process as before: calculate gradient and find estimator

$$\frac{\partial \log\{\mathcal{L}(\mathbf{y}; \hat{\mu}, \sigma)\}}{\partial \mu} = \frac{1}{2\sigma^2} \left( 2 \sum_{i=1}^n y_i - 2n\mu \right) \quad (4)$$

Giving..

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n y_i \quad (5)$$

It is a linear function of  $y_i$  so is normally distributed.

## Estimating $\sigma^2$

---

Same process as before. MLE is biased so gets a small correction.

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \hat{\mu})^2 \quad (6)$$

Is it a quadratic function of  $y_i$  so is  $\chi^2$ -distributed.

## Uncertainty of $\hat{\mu}$

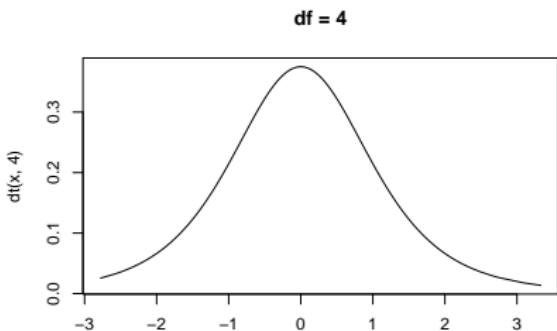
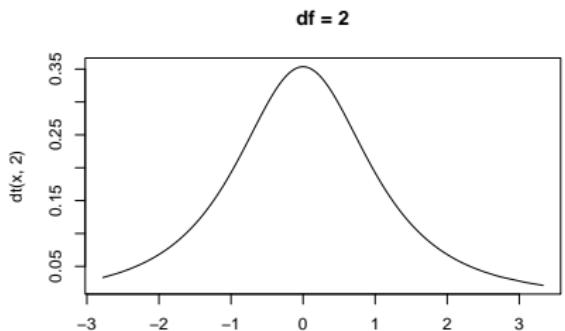
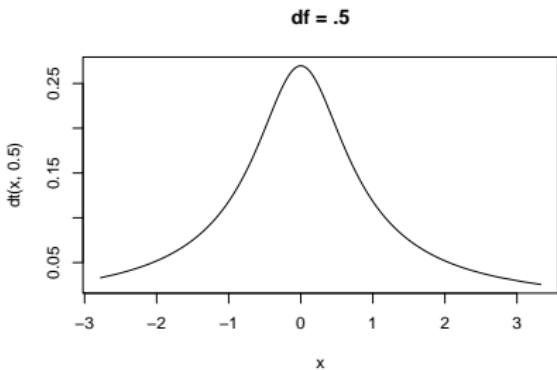
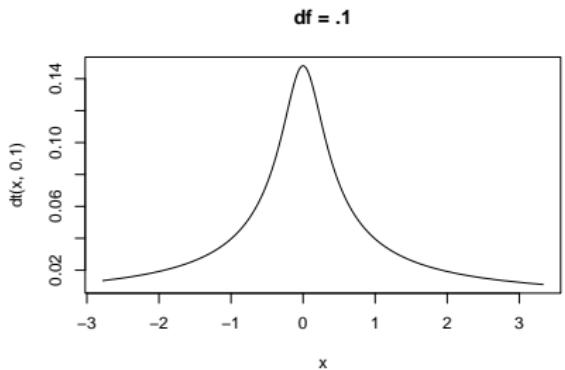
---

$$\begin{aligned}\text{var}(\hat{\mu}) &= \mathbb{E}(\hat{\mu}^2) - \mathbb{E}(\hat{\mu})\mathbb{E}(\hat{\mu}) \\ &= \frac{1}{n^2} \sum_{i=1}^n \mathbb{E}(y_i^2) - \mathbb{E}(\hat{\mu})\mathbb{E}(\hat{\mu}) \\ &= \frac{1}{n}\sigma^2\end{aligned}\tag{7}$$

- ▶ Depends on  $n$  (small  $n$ , large uncertainty)
- ▶ Depends on  $\sigma^2$
- ▶ Which is estimated by  $\hat{\sigma}^2$
- ▶ But that estimate also has uncertainty

So, we use the  $t$ -distribution to represent that additional uncertainty.

# The t-distribution



# The t-test

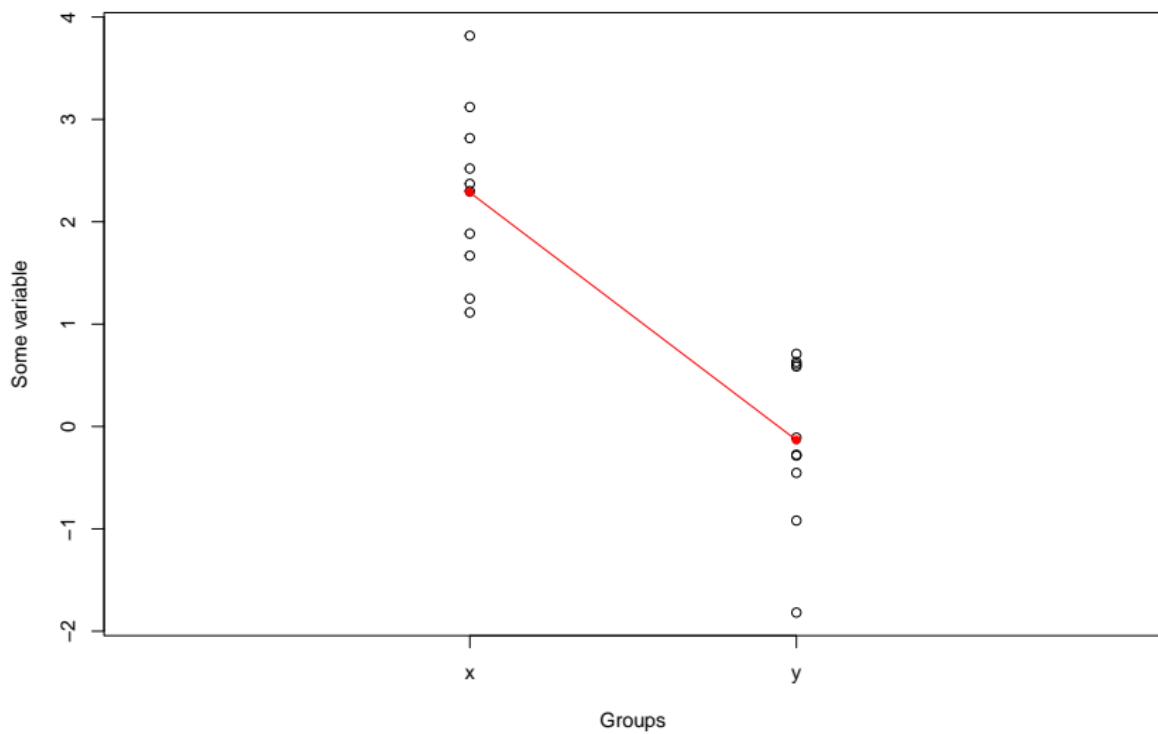
---

- ▶ The t-test is type of linear regression
- ▶ With one covariate and two groups

```
set.seed(12345)
y <- rnorm(10)
x <- rnorm(10, mean = 2)
t.test(x, y)
```

```
##
##  Welch Two Sample t-test
##
## data: x and y
## t = 6.5336, df = 17.979, p-value = 3.872e-06
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  1.641042 3.196802
## sample estimates:
## mean of x mean of y
## 2.2859778 -0.1329441
```

## t-test visualized



## t-test as linear regression

---

```
data <- data.frame(y = c(x, y),
                     var = c(rep(c("x","y"),each=10)))
lm(y~0+var, data = data)
```

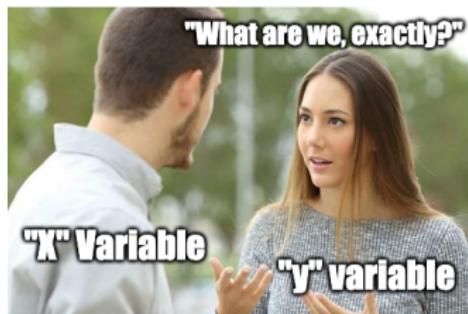
```
##  
## Call:  
## lm(formula = y ~ 0 + var, data = data)  
##  
## Coefficients:  
##       varx      vary  
##  2.2860   -0.1329
```

## What is a linear regression?

---

Models with a continuous **response variable** as a function of one or more **explanatory variable**. Variables are connected by linear equations.

- ▶  $y_i$ : the **response variable**, can only be numerical
- ▶  $x_i$ : the **explanatory variable**, can be categorical (0,1) or numerical



$$y_i = \alpha + x_i \beta + \epsilon_i \sim \mathcal{N}(0, \sigma^2) \quad (8)$$

## Synonyms

---

- ▶ Covariate
- ▶ Predictor (variable)
- ▶ Explanatory variable
- ▶ Independent variable

They all refer to  $x_i$ .

## What is the goal of regression?

---

We measure data  $y_i$  and want to infer its with  $x_i$

Steps:

- 1) We decide on a model
- 2) We estimate the parameters
- 3) We check if it is a valid and good model
- 4) We draw our conclusion (with uncertainty)

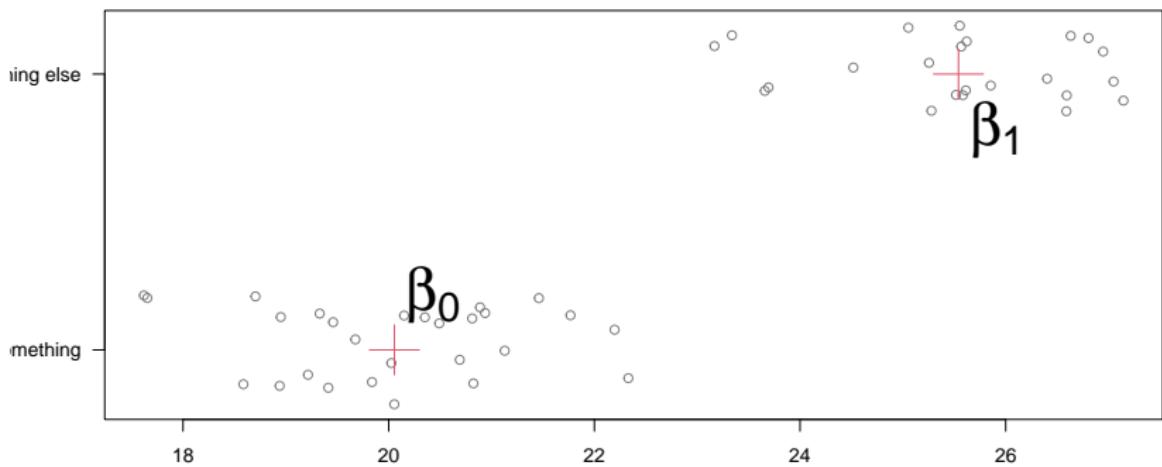
## Examples of linear models: categorical $x_i$

---

$$\mu_i = \begin{cases} \beta_0 & \text{if } X_i = 0 \\ \beta_1 & \text{if } X_i = 1 \end{cases} \quad y_i = (1 - x_i)\beta_0 + x_i\beta_1 + \epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

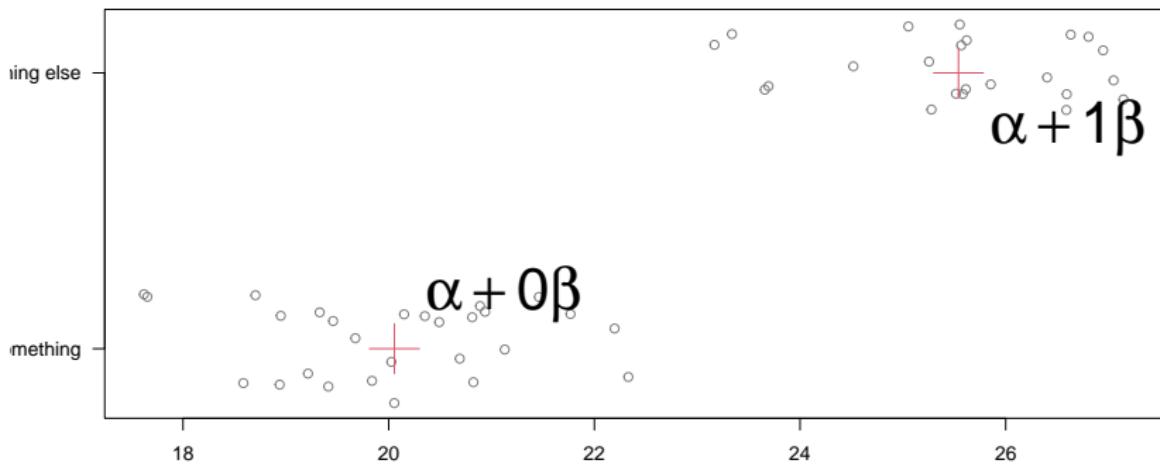
$$\mu_i = \begin{cases} \alpha & \text{if } X_i = 0 \\ \alpha + \beta & \text{if } X_i = 1 \end{cases} \quad y_i = \alpha + x_i\beta_1 + \epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

## Examples of linear models: categorical $x_i$



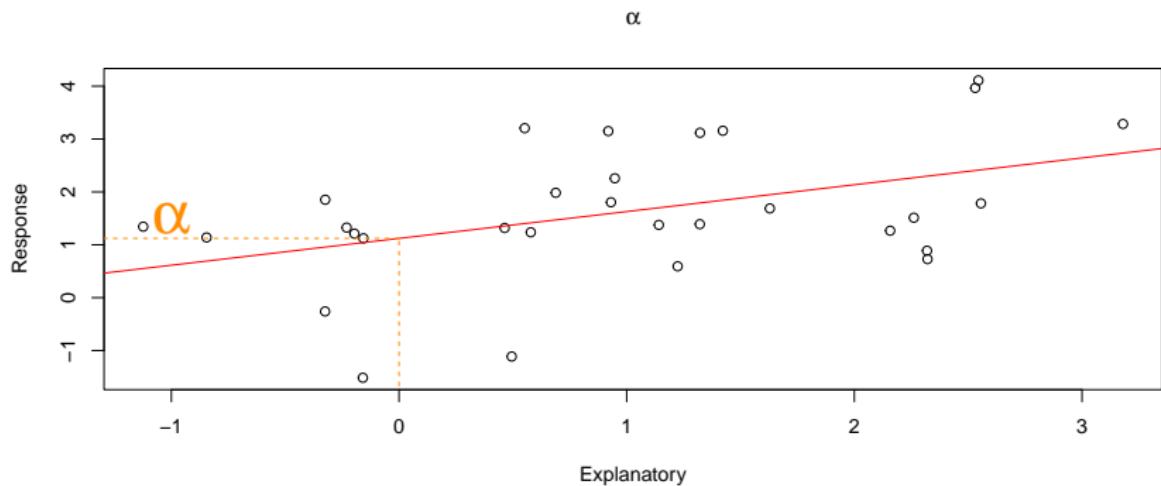
- ▶  $\beta_0$  is the group 1 mean
- ▶  $\beta_1$  is the group 2 mean

## Examples of linear models: categorical $x_i$



- ▶  $\alpha$  is the mean of the first group
- ▶  $\beta$  is the deviation from the mean of the first group

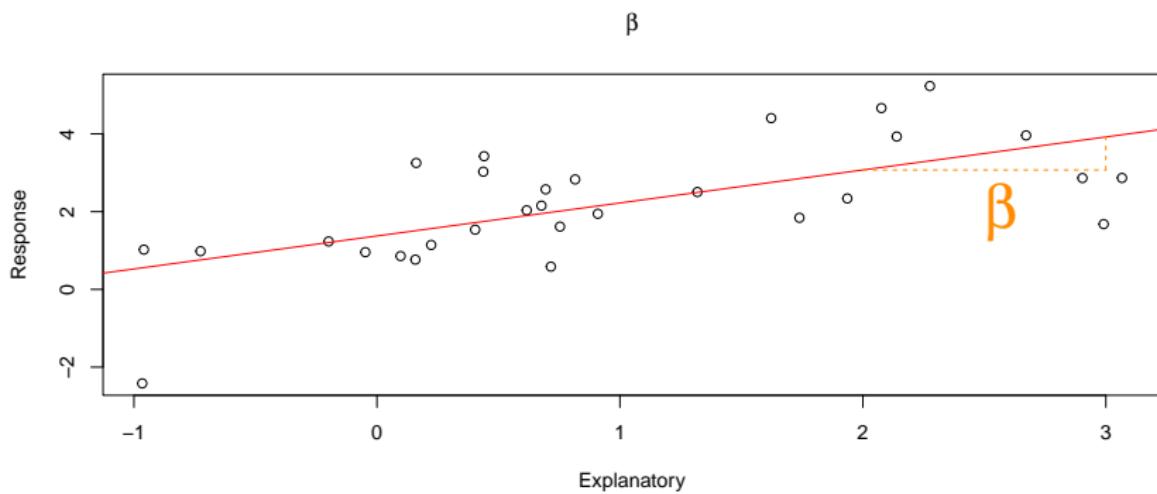
## Examples of linear models: continuous $x_i$



$$y_i = \alpha + x_i\beta + \epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

►  $\alpha$ : the intercept is the value of  $y_i$  where  $x_i = 0$

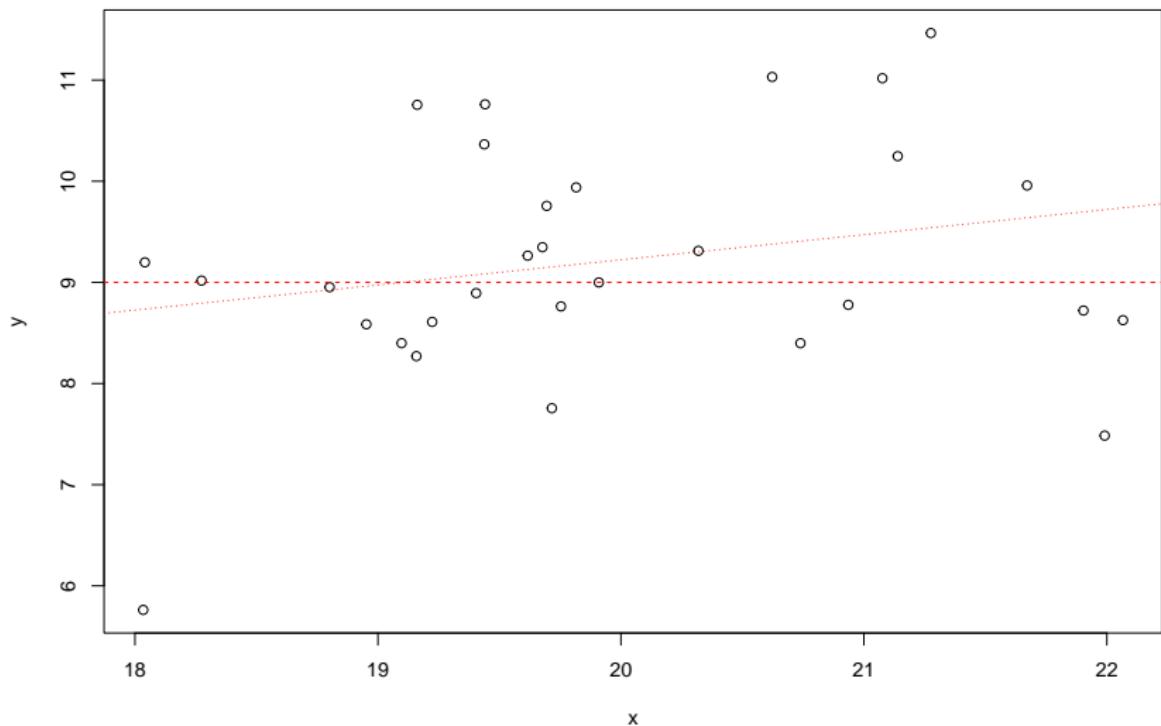
## Examples of linear models: continuous $x_i$



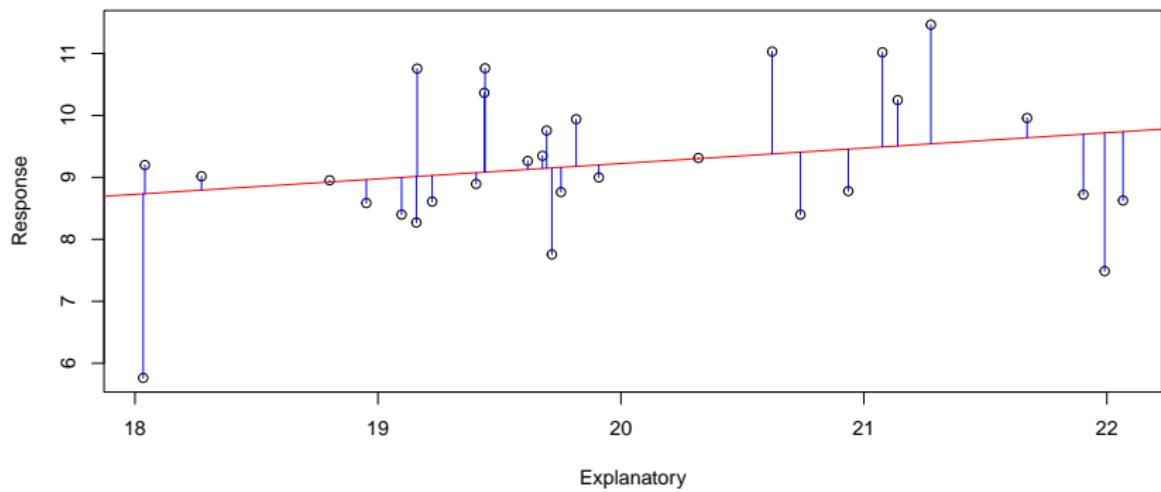
$$y_i = \alpha + x_i \beta + \epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

- ▶  $\alpha$ : the intercept is the value of  $y_i$  where  $x_i = 0$
- ▶  $\beta$ : the slope is the change in  $y_i$  for a unit increase in  $x_i$

## What is the best line?



## How good is the line?



Distance from model (line) to data: “error”  $\epsilon_i$

## Least squares estimation

---

Minimize the sum of squared residuals:

$$RSS = \sum_{i=1}^n \epsilon_i^2 \tag{9}$$

which is the same to maximizing the normal likelihood!

- ▶  $\hat{\alpha} = \frac{1}{N}(\sum y_i - \hat{\beta} \sum x_i)$
- ▶  $\hat{\beta} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$
- ▶  $\hat{\sigma}^2 = \frac{1}{N-1} \sum (y_i - (\hat{\alpha} + x_i \hat{\beta}))^2$
- ▶  $\hat{\mu}_i = \hat{\alpha} + x_i \hat{\beta}$

## Our model

---

$$y_i = \alpha + \beta x_i + \epsilon_i \sim \mathcal{N}(0, \sigma^2) \quad (10)$$

- ▶  $y_i$ : our data
- ▶  $\alpha, \beta$  describe our line
- ▶  $\epsilon_i$  quantifies distance to the model

## Assumptions

---

We make some critical assumptions here

- 1) the relationship between  $y_i$  and  $x_i$
- 2) distribution of the errors

## Other assumptions for the Errors

---

$$\epsilon_i \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2) \quad (11)$$

- ▶ are normally distributed
- ▶ have constant variance (“Homoscedasticity”)
- ▶ are independent
- ▶ no outliers

## Assumptions for the errors

---

**all the errors together tell us how good the line is**

- ▶ this is the same as finding the line by maximum likelihood estimation 😊

$$y_i \sim \mathcal{N}(\alpha + x_i\beta, \sigma^2) \quad (12)$$

## Summary simple linear models

---

- ▶ includes t-test, anova (analysis of variance) and regression
- ▶ all use the same mathy bits (and model)
- ▶ **interpretation** depends on the type of **variable**
- ▶ GLMs take the same form

## Summary

---

- ▶ Normal-distribution and t-distribution
- ▶ Simple linear regression: one covariate
- ▶ Least squares estimation
- ▶ Difference in LMs with categorical or continuous explanatory variables
  - ▶ Categorical: intercept/mean parameter
  - ▶ Continuous: slope parameter
- ▶ Fortunately we have the `lm()` function in R!
- ▶ More on assumptions checking tomorrow

## Summary (2)

---

t-test: linear model with categorical covariate of two groups

ANOVA: linear model with categorical covariate of multiple groups

ANCOVA: linear model with categorical and continuous covariate  
(interaction)