# Stat 412/512

## WRAPPING UP INFERENCE

Jan 28 2015

Charlotte Wickham

stat512.cwick.co.nz

# Announcement

DA#1 posted

Don't need to do regression diagnostics (i.e. residual plots)

Read report description even if it seems familiar from ST511

Submit a report that contains no R code or raw R output.  Also submit an R code file.

# Two models

Full model:

μ{ log Energy | log Mass, Type}

$= \beta_0 + \beta_1$ log Mass $+ \beta_2$ *bird* $+ \beta_3$ *ebat*

Reduced model:

μ{ log Energy | log Mass, Type}

$= \beta_0 + \beta_1$ log Mass

If the reduced model is the truth:
  then $\beta_2$ and $\beta_3$ should be estimated close to zero
  both models should fit about the same
  the residuals in both models should be about the same size

```
> anova(fit_eq, fit_bats)
Analysis of Variance Table

Model 1: log(Energy) ~ log(Mass)
Model 2: log(Energy) ~ log(Mass) + Type
  Res.Df    RSS Df Sum of Sq      F Pr(>F)
1     18 0.58289
2     16 0.55332  2  0.029574 0.4276 0.6593
```

There is no evidence that the mean log energy differs for birds, echolocating bats and non-echolocating bats after accounting for body mass (extra sum of squares F-test, p-value = 0.66).

# Another example

We relied on a parallel lines regression to answer our question of interest, we might also want to test this is reasonable.

Fit separate lines model (check assumptions look good)

Test whether interaction terms can be dropped.
Full model:

$\mu\{$ log Energy | log Mass, Type$\}$

**6 parameters**

$\quad = \quad \beta_0 + \beta_1$ log Mass $+ \beta_2$ *bird* $+ \beta_3$ *ebat* $+$

$\qquad \beta_4$ log Mass x *bird* $+ \beta_5$ log Mass x *ebat*

Reduced model:

**4 parameters**

$\mu\{$ log Energy | log Mass, Type$\}$

$\quad = \quad \beta_0 + \beta_1$ log Mass $+ \beta_2$ *bird* $+ \beta_3$ *ebat*

$\beta_4 = \beta_5 = 0$

```
> anova(fit_bats, fit_sep)    ← post code
Analysis of Variance Table
```

Model 1: log(Energy) ~ log(Mass) + Type
Model 2: log(Energy) ~ log(Mass) + Type +
   log(Mass):Type

| Res.Df | RSS | Df | Sum of Sq | F | Pr(>F) |
|---|---|---|---|---|---|
| 1 | 16 | 0.55332 | | | |
| 2 | 14 | 0.50487 | 2 | 0.04845 | 0.6718 | 0.5265 |

**Your Turn:** Write a summary of this result.

There is no evidence,
the effect of log mass on mean log energy depends on animal
type.                                *no evidence of*
                                     *an interaction*

There is no evidence that the relationship between mean log
energy and log body mass differs for birds, echolocating bats
and non-echolocating bats (extra sum of squares F-test, p-

# Extra SS F-test

**Null hypothesis:**

The parameters in the full model are constrained.

Reduced model is correct.

**Alternative hypothesis:**

The parameters in the full model are unconstrained.

*prefer the full model*

A small p-value gives us evidence against the reduced model.

# Overall regression F-test

**Null hypothesis:**

$\mu\{ Y \mid X \} = \beta_0$ ← *All β's apart from $\beta_0$ are zero*

↑ many

constant mean

**Alternative hypothesis:**

The parameters in the full model are unconstrained. *At least one β (other than intecept) is not zero*

If we reject this null, then not all parameters are zero, (this is not the same as all parameters are non-zero)

For the bats:

**Null**: $\mu\{ \log \text{Energy} \mid \log \text{Mass, Type}\} = \beta_0$

**Alternative:** $\mu\{ \log \text{Energy} \mid \log \text{Mass, Type}\}$

$= \beta_0 + \beta_1 \log \text{Mass} + \beta_2 \, bird + \beta_3 \, ebat$

```
> summary(fit_bats)

Call:
lm(formula = log(Energy) ~ log(Mass) + Type, data = case1002)

Residuals:
    Min      1Q  Median      3Q     Max
-0.23224 -0.12199 -0.03637  0.12574  0.34457

Coefficients:
                               Estimate Std. Error t value Pr(>|t|)
(Intercept)                    -1.57636    0.28724  -5.488 4.96e-05 ***
log(Mass)                       0.81496    0.04454  18.297 3.76e-12 ***
Typenon-echolocating birds      0.10226    0.11418   0.896    0.384
Typeecholocating bats           0.07866    0.20268   0.388    0.703
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.186 on 16 degrees of freedom
Multiple R-squared: 0.9815,    Adjusted R-squared: 0.9781
F-statistic: 283.6 on 3 and 16 DF,  p-value: 4.464e-14
```

overal
regression F-test

A extra sum of squares F-test, with the reduced model:

$$\mu\{ \log \text{Energy} \mid \log \text{Mass}, \text{Type}\} = \beta_0$$

I.e. Null: $\beta_1 = \beta_2 = \beta_3 = 0$

# Meadowfoam case study

Intensity could be treated as continuous variable:

$\mu\{$ *flowers* | *Intensity*, *early*$\} =$

$$\beta_0 + \beta_1 early + \beta_2 Intensity \quad \longleftarrow$$

## Or as a categorical variable:

$\mu\{$ *flowers* | *Intensity*, *early*$\} =$

$$\beta_0 + \beta_1 early + \beta_2 L300 + \beta_3 L450 +$$
$$+ \beta_4 L600 + \beta_5 L750 + \beta_6 L900$$

$$\mu\{ \textit{flowers} \mid \textit{Intensity, early}\} =$$
$$\beta_0 + \beta_1 \textit{early} + \beta_2 \textit{Intensity}$$

# 3 parameters



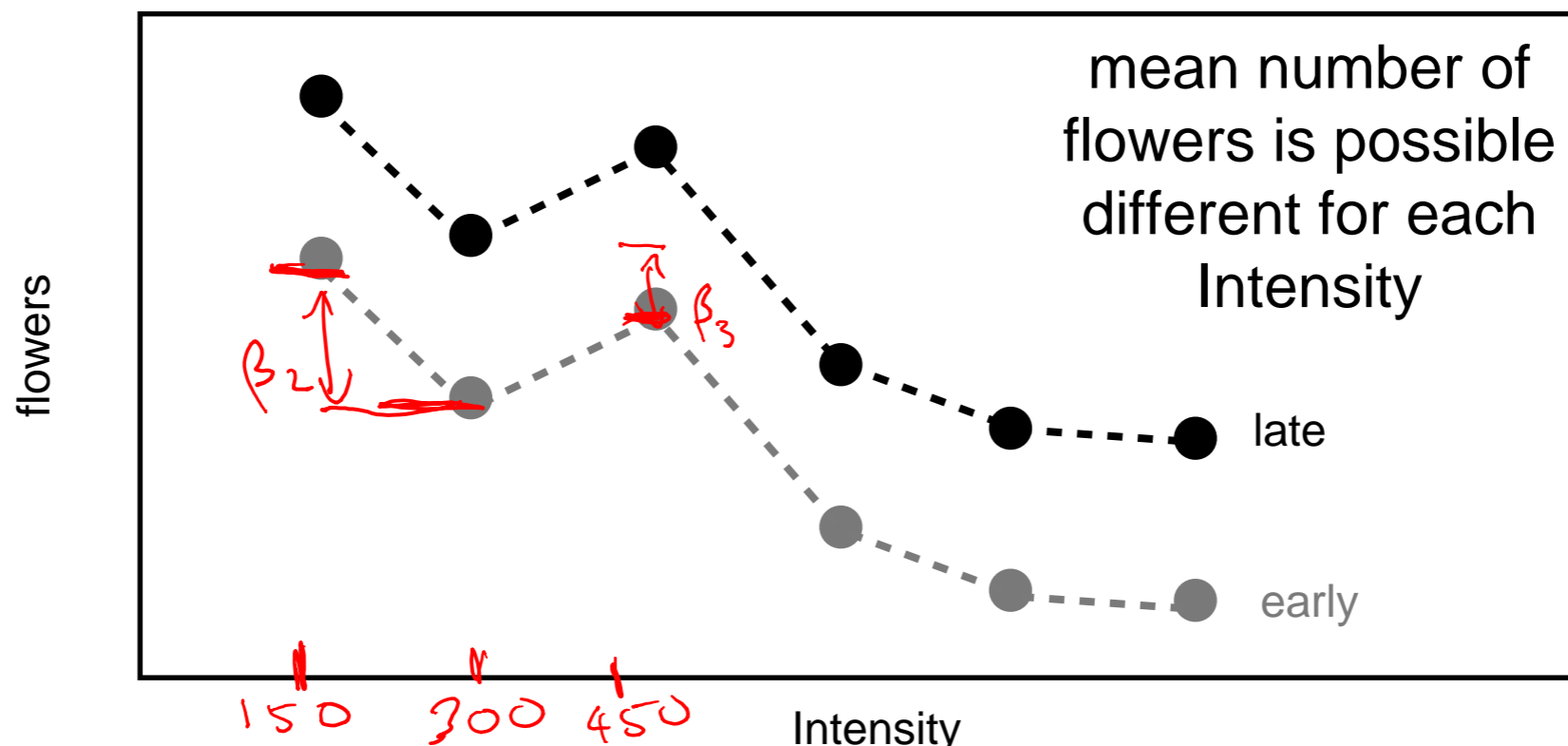mean number of flowers is a straight line function of Intensity

late

early

Intensity

straight
$\beta_2$ slope

$$\mu\{ \textit{flowers} \mid \textit{Intensity, early}\} =$$
$$\beta_0 + \beta_1 \textit{early} + \beta_2 L300 + \beta_3 L450 +$$
$$+ \beta_4 L600 + \beta_5 L750 + \beta_6 L900$$

# 7 parameters



mean number of flowers is possible different for each Intensity

late

early

$\beta_2$

(extra
parameters)

$\beta_2$

$\beta_3$

150   300  450

Intensity
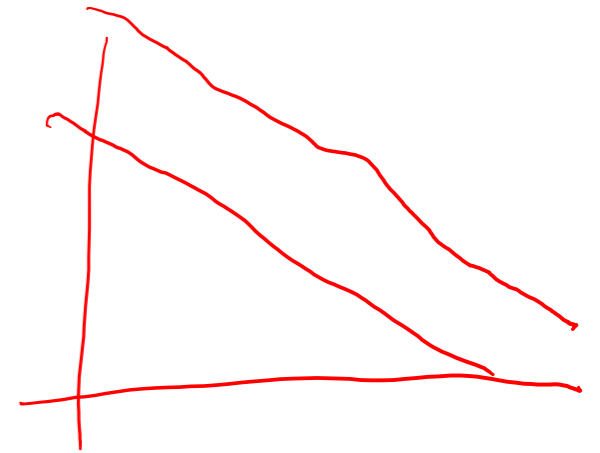
# In general

A model with a continuous variable is a constrained case of a model with the same variable represented as categories.

An extra sum of squares F-test can be used to compare them.

```
> fit_cont <- lm(Flowers ~ Intens + Time,
          data = case0901)
> fit_ind <- lm(Flowers ~ factor(Intens) + Time,
          data = case0901)

> anova(fit_cont, fit_ind)
Analysis of Variance Table

Model 1: Flowers ~ Intens + Time
Model 2: Flowers ~ factor(Intens) + Time
  Res.Df    RSS Df Sum of Sq      F Pr(>F)
1     21 871.24
2     17 767.47  4    103.76 0.5746 0.6848
```
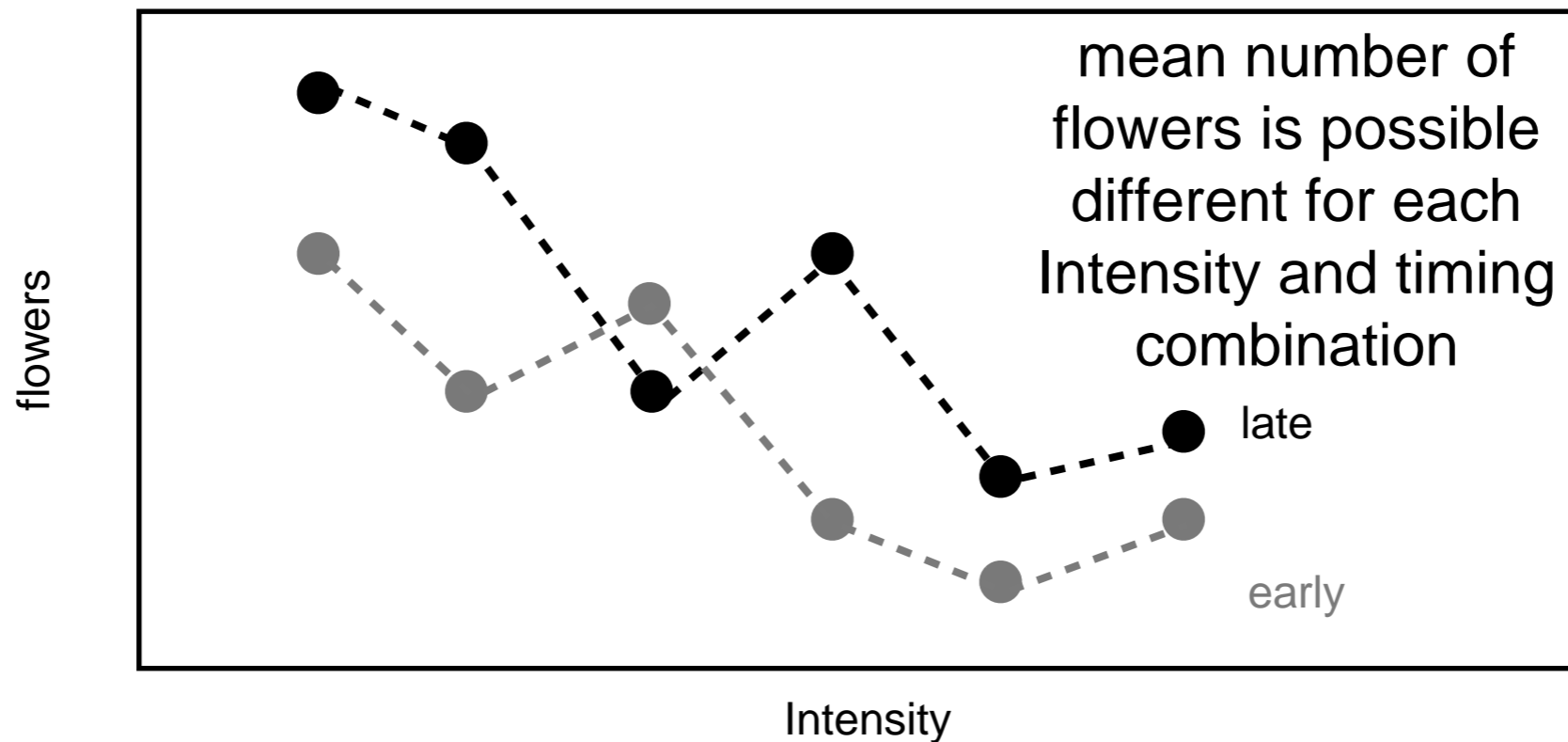
**Your Turn:** Write a summary of this result.

There is no evidence against the mean number of Flowers being a straight line function of Intensity (extra SS F-test ... p= 0.68)

# HW#2

$\mu\{$ *flowers | Intensity, early*$\} =$

$\beta_0 + \beta_1 early + \beta_2 L300 + \beta_3 L450 + \beta_4 L600 + \beta_5 L750 + \beta_6 L900 +$

$\beta_7 L300 xearly + \beta_8 L450 xearly + \beta_9 L600 xearly + \beta_{10} L750 xearly + \beta_{11} L900 xearly$

## 12 parameters



mean number of flowers is possible different for each Intensity and timing combination

late

early

flowers

Intensity

The assumptions for the F-test, are that:

**the full model is appropriate**

and the usual regression assumptions:

- constant spread
- the response is normally distributed around the mean
- observations are independent

Before doing the F-test you need to check these!

*except once in DA #1*

A small p-value gives us evidence against the reduced model, **assuming** the full model is true.
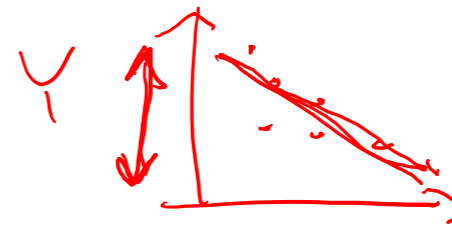
If the full model is inappropriate:

the response is non-linear
you left out important terms, etc

the F-test tells you nothing.

# $R^2$

R-squared tells you the proportion of variance in the response explained by the explanatories.

```
> summary(fit_intensonly)    µ{ flowers | Intensity, early} =  intensity
....
Residual standard error: 8.94 on 22 degrees of freedom
Multiple R-squared: 0.5947,   Adjusted R-squared: 0.5763
F-statistic: 32.28 on 1 and 22 DF,  p-value: 1.03e-05
```

```
> summary(fit_cont)    µ{ flowers | Intensity, early} =  intensity + TIME
...
Residual standard error: 6.441 on 21 degrees of freedom
Multiple R-squared: 0.7992,   Adjusted R-squared:  0.78
F-statistic: 41.78 on 2 and 21 DF,  p-value: 4.786e-08
```

The linear relationship with intensity explains 59% of the variability in the mean number of flowers per stem.

The additive effect of early explains an additional 20% of the variability in the mean number of flowers per stem.

# But $R^2$ always gets bigger

The more variables you add to the model, the bigger $R^2$ gets.

If you add as many variables as observations, then $R^2 = 1$.

**Adjusted R-squared**, tries to adjust for this.  If the adjusted $R^2$ increases then the additional variable explained more variance than expected by chance.

# The principle of parsimony

The simplest explanation is the best.

In statistics, translates to:

If two models fit the data equally well, use the simpler one.

In practice, leads to an **acceptance of the null** model in an F-test. (!)

not everyone agrees with this