

Week 5

## ANOVA with multiple variables/factors

outcome is still quantitative

you have multiple explanatory variables

goal:

to identify contributions of different variables

A/B test - Obama Campaigning

Relating score to rating (sinemas)

$$s_i = b_0 + b_1 \mathbb{I}(R_{a_i} = "PG") +$$

$$b_2 \mathbb{I}(R_{a_i} = "PG-13") +$$

$$b_3 \mathbb{I}(R_{a_i} = "R") +$$

$\epsilon$

$$\mathbb{I}(R_{a_i} = "PG") - \text{logical value } (R_{a_i} = "PG")$$

"indicator factors"

average values:

$$b_0 = \text{avg for } G$$

$$b_0 + b_2 = \text{avg for PG-13}$$

$$b_0 + b_1 = \text{avg for PG}$$

$$b_0 + b_3 = \text{avg for R}$$

## ANOVA in R

aov  
instead of lm

$$\text{aovObj} = \text{aov}(\text{movies\$score} \sim \text{movies\$rating})$$

aovObj\$coeff

(Intercept)	\$rating PG	\$rating PG-13	\$rating R
67.65	-12.59	-11.81	-12.02

↑  
avg for G

## Adding a second factor

$$S_i = b_0 + b_1 \mathbb{I}(R_{ai} = \text{"PG"}) + b_2 \mathbb{I}(R_{ai} = \text{"PG-13"}) + \dots + y_1 \mathbb{I}(G_i = \text{"action"}) + y_2 \mathbb{I}(G_i = \text{"animated"}) + \dots + e$$

↑  
genre

↑  
lots of genres

↑  
only 2 variables!

$aov2 = aov(\text{score} \sim \text{rating} + \text{genre})$

$\text{summary}(aov2)$

F value



variation explained  
by factor

for 1st (rating) -  
only by itself

for 2nd (genre) -  
only what wasn't  
explained by rating.  
(after taking into account  
rating)

↑  
order  
↓  
matters!



$aov2 = aov(\text{score} \sim \text{genre} + \text{rating}) \leftarrow \text{better}$

• you can add quantitative variables

→ so we may add box office

$aov3 = aov(\text{score} \sim \text{genre} + \text{rating} + \text{box office})$



only 1 degree of freedom

Language:

unit - one observation

treatment - applied to units

factors - controlled by experimenters

replicates - multiple (independent) units with the same factors / treatments

Useful resources

(wiki pages: Experimental Design  
ANOVA,  
A/B Testing)

Binary outcomes

like:

- Dead / Alive
- Win / Loss
- Success / Failure
- etc

} binary outcomes  
or  
0/1 outcomes

Linear regression - may be not the best

## Ravens data

win num      win      score      opponent. score  
(W/L)

Linear model:

$$RW_i = b_0 + b_1 \cdot S_i + e_i$$

$RW_i$  1 if W, 0 if not

$S_i$  - number of points they scored

$b_0$  - probability of a Ravens win <sup>with</sup> (0 score)

$b_1$  - increase in probability by each 1 point

$e_i$  - error

$\ln(\text{win num} \sim \text{raven. score})$

not good

(sometimes  $Pr > 1$  which is impossible)

etc

Binary Outcomes:  $RW_i$   
0/1

Probability (0, 1)

$$\Pr(RW_i \mid RS_i, b_0, b_1)$$

Odds (0,  $\infty$ )

$$\frac{\Pr(RW_i \mid RS_i, b_0, b_1)}{1 - \Pr(RW_i \mid RS_i, b_0, b_1)}$$

Log odds -  $\log(\text{odd})$   
( $-\infty, \infty$ )

### Linear vs Logistic Regression

Linear

Logistic

instead of modelling  $RW_i$   
we model the probability based on  
odds

$$\Pr(RW_i | RS_i, b_0, b_1) = \frac{e^{b_0 + b_1 RS_i}}{1 + e^{b_0 + b_1 RS_i}}$$

$$\underbrace{\log \left( \frac{\Pr(RW_i | RS_i, b_0, b_1)}{1 - \Pr(RW_i | RS_i, b_0, b_1)} \right)}_{\text{log odds}} = \underbrace{b_0 + b_1 RS_i}_{\text{Ravens score}}$$

↑  
any number  
 $(-\infty; +\infty)$   
not  $[0, 1]$

interpretation

$b_0$  - log odds of Win if they score 0 points

$b_1$  - log odds ratio of win probability for each point scored (compared to 0 points)

$\exp(b_1)$  - odds ratio of win probability for each point scored (compared to 0 points)

in R:

glm command

$\text{logRegRavens} = \text{glm}(\underbrace{\text{winNum}}_{\text{outcome}} \sim \underbrace{\text{score}}_{\text{covariates}}, \text{family} = \text{"binomial"})$   
↑  
logistic regression

coefficients are interpreted differently

↙ it gets log!  
 $\exp(\text{logRegRavens}\$coeff)$

↳ should see if  
score is bigger than 1  
(more chances to win)

$\exp(\text{confint}(\text{logRegRavens}))$

for confidence intervals



## Anova for logistic regression

anova(logRegParams, test = "Chisq")

↑  
analysis of Deviance Table

Simpson's Paradox - take a look at w.k.i.

## Interpreting Odds Ratios

- not probabilities
- odds ratio of 1 = no difference in odds  
↓
- log odds ratio of 0 = no difference in odds
- $0.5 < \text{odds ratio} < 2$  - "moderate effect"

• Relative Risks  $\frac{\Pr(RW_i | RS_i = 10)}{\Pr(RW_i | RS_i = 0)}$

often easier to interpret

Wiki  
↓

Odds Ratio

- for small probabilities  $RR \approx OR$ , but they are not the same

## Count Outcomes

- many data take form of counts
  - calls to a call center
  - number of flu cases in area
  - numbers of cars that cross a bridge
- data may also be in the form of rates
  - percent of students passing <sup>a</sup> the test
  - percent of hits to a website from a country

Linear regression is an option

Poisson distribution can be used to model this data

set.seed(3433)

~~plot~~ rpois(100, lambda = 100)



rate

(an avg number of calls  
coming to a call center)

Spread is bigger for bigger lambdas  
and it controls both mean and var

$\text{mean}(\text{pois.d}) \approx \text{var}(\text{pois.d})$   
very close

Web site traffic:

possible to fit regression here  
↓

$$NH_i = b_0 + b_1 JD_i + e_i$$

$NH_i$  - number of hits

$JD_i$  - day of the year (Julian day  
(from 01.01.70))

$b_0$  - number of hits on Julian day 0

$b_1$  - increase in number of hits per unit day

$e_i$  - error

ln(visits ~ julian)

Linear vs Poisson Regression (log-linear)

↓

$$\log(E(NH_i | JD_i, b_0, b_1)) = b_0 + b_1 JD_i$$

or

$$E[NH_i | JD_i, b_0, b_1] = e^{b_0 + b_1 JD_i} =$$

$$= e^{b_0} \cdot e^{b_1 JD_i}$$

↑

If  $JD_i$  is increased by 1,  $E[NH_i | JD_i, b_0, b_1]$   
is multiplied by  $\exp(b_1)$

$glm(visits \sim julian, family = "poisson")$

↑  
poisson regression

To model rates

$$\log(E[NHSS | JD_i, b_0, b_1]) = \log(NH_0) + b_0 + b_1 JD_i$$

↑  
(number of hots  
from a specific  
website)

more information

Wiki on Poisson Regression

## Model Checking and Model Selection

Basic assumptions for linear regressions

- Variance is constant
- trend is linear
- no big outliers
- no biases

- what to do if variance grows?
  - see if other variable explains the growth
  - sandwich library  
library(sandwich)  
 $lm1 = \text{data } lm(\text{data1} \sim \text{data2})$   
 $vcovHC(lm1)$

- the trend is not linear.

- use Poisson regression
- use data transformation (log, etc)
- use linear regression

+  $vcovHC$

- missing covariate

- check covariates carefully -  
use exploratory analysis
- report unexplained pattern

- outliers

how much influence do they have?

- fit regression with the outliers and then without
- if the two slopes are very different, then the outliers have a major effect

↓

more caution needed

so if you know for sure they are mistakes -  
remove and document them

If they are real - consider reporting how  
sensitive your estimate is to the  
outlier

Consider using a robust linear model fit

like `rlm {MASS}`

(downweights outliers)

Model Checking

- Default Plots

after fitting  $lm$ , try to see it's plot

~~the~~  $\text{plot}(lm)$

↑

Residuals vs Fitted

- Deviance

Commonly reported measure

might tell you that the model is wrong

-  $R^2$

may be a bad summary

## Model Selection

Usually you have a lot of variables -  
you have to do some sort of filtering

How to choose correct variables?

- have domain-specific knowledge (prev. experience)
- exploratory analysis  
(make plots, make plots of residuals, coloring them by different variables)
- Statistical selection
  - Step-wise (add/remove one var at a time)
  - AIC / BIC
  - etc

↑  
may bias your inference, so don't overdo  
statistical selection



## Error measures

- $R^2$  (not always good enough)
- Adjusted  $R^2$  - takes into account the number of estimated parameters
- AIC Information criteria
- BIC

## Model Selection - Step

$\text{lm1} = \text{lm}(\text{score} \sim \cdot, \text{data} = \text{movies})$

↙ all the terms

aic Formula =  $\text{step}(\text{lm1})$

↑  
recomputes, finds better orders,  
adds, deletes, etc.

↓  
 $\text{score} \sim \text{box. office} + \text{running. time}$

## Regsubsets

library(leaps)

regSub = regsubsets(score ~ ., data = movies)

↑  
calculates BIC score  
for all possible subsets

plot(regSub) ↓

Goal: minimize BIC

## Notes And Resources

- exploratory / visual analysis - is a key
- automatic selection produces an answer,  
but it may bias inference (outfit - fits yours,  
but won't fit other sample)
- you may think of separating the sample  
into 2 groups

use 1st to estimate

2nd to do the inference

- goal - not to get "causal" model.

LARS package

Elements of Machine Learning -  
(book from the TOCREAD list)

model selection  
choosing variables