

Multiple Regression

- $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \varepsilon$
- As our measuring gets better/cheaper, we have more variables
 - Earth sciences advance (climatology, remote sensing)
- Two key challenges
 - Too many variables to make sense of
 - Collinearity (correlation between variables)
 Less noticed in ANOVA because usually factors are "orthogonal" i.e. uncorrelated





Scaling in R

bd=read.table("/851/birds.csv",sep=
 ",",h=T)
names(bd)
bd2=bd[bd\$Rangesize>0,]
summary(bd2)
summary(lm((log10(TotalAbund+.1))~(
 log10(Mass))+(log10(Rangesize))
 ,data=bd2))
summary(lm(scale(log10(TotalAbund+.
 1))~scale(log10(Mass))+scale(lo
 g10(Rangesize)),data=bd2))

Handling Many Variables II -Nested models

- Model A is nested in model B if all the terms in A are also in B
- Example
 - Y=a+bx is nested in y=a+bx+bx² and in y=a+bx+cz
- Einstein says "make a model as simple as possible but not simpler than needed"
 - You can compare the R² (adjusted?) of nested models
 - Can also do ANOVA tests (F-statistics) for statistical significance (degrees of freedom adjust for increased parameters)



Comparing nested models II -Nested test In R: anova(m1,m2,m3,...) Reports progessive tests In general with nested models Test most "complex" term in largest model If not significant, throw it out and look at smaller model EG x^2 or interaction term invalides interpretation of x or A+B Don't use effect estimates, r² from model with higher order term in

Handling many variables III -Stepwise regression

- Have dozens of variables
- Let the computer pick the equation
- Forward and backward steps
 - Forward
 - start with intercept
 - do a regression with each variable added on
 - Pick one with greatest increase in adjusted $r^2 \mbox{ or } F$
 - Repeat
 - Backward
 - Similar but start with all variables in
 - Calculate with each variable dropped, keep best
 - Stepwise = mix forwards and backwards

Stepwise Caveat

- Stepwise is OK for exploration
- Freedman's paradox
 - If you have enough variables, even random data will have some strong relations
 - Stepwise regression pulls these out
- Whittingham bad! Murtaugh fine!
 - On website
- Beware of significance tests on regressions developed stepwise
- If you want to do a significance test, have an a priori hypothesis!



Handling many variables IV -Really radical technique

- Think!
- Form hypotheses use your knowledge to select variables a priori. Use biology.
- EG climate
 - What factors most important
 - Max Ann Temp, Growing season length correlated which matters more? Use one!
 - Plant model winter freezing, summer growth, summer drought
 - → Annual Min Temp, Growing Degree Days, Palmer Drought Severity Index









Assesment I - Common symptoms in regression

- From Belsley 1991
 - Coefficients have the wrong sign
 - Important predictors have high p-values
 - Deletion of one column or row causes the betas to change drastically



- A square pxp (p=#variables) matrix containing covariances
- $P = c^* exp(-(x-\mu)\Sigma^{-1}(x-\mu)/2)$
- Standardized matrix
 - each column converted to a z-score (subtract mean, divide by std dev)
 - Avoid units making some variables appear more important
- Standardized covariance matrix=correlation matrix





























