

ADV. DATA ANALYSIS

MIDTERM 1

OLD

Lets start w/ reviewing basic R and data analysis steps we learned in STOR 455 ...

DATA ANALYSIS

- Anytime we get a new data set we begin by looking at initial graphical images, 5 # summaries, skewness & outliers, and correlations.
- We might want to clean some data or make new variables

LINEAR MODEL

$$y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_{p-1} X_{p-1} + \epsilon$$

response parameters predictors error
intercept

- We find the least squares estimate of the parameters $\hat{\beta}$, in order to minimize ϵ .
- The difference between \hat{y} (actual response) and predicted response \hat{y} are the residuals.
- Variance of the error σ^2
- We can put in qualitative predictors w/ binary dummy vars.
- We can also add interaction terms $\rightarrow \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2$ which allows us to explore how the influence of x_1 on the model changes based on if $x_2 = 0$ or 1.

INTERPRETATION

- Intercept** \rightarrow "given all predictors (X 's) are 0, β_0 will be the response"
- Coeffs** \rightarrow "given all other predictors are held constant we predict a β_n increase in response for every 1 increase in X_n ."
- Hypothesis Testing** \rightarrow To determine sign. of predictors in the model. (assume errors are iid.) Run F test and look at p-values.
Z test \rightarrow comparing 2 sample means, $\frac{\bar{x} - \mu}{SE(\bar{x})}$
t test \rightarrow used with more unknowns (typically we use for testing sign. of coeffs. $(\beta - \beta_0) / SE(\beta)$ w/ df = n-1.
- Confidence Intervals** \rightarrow More useful in judging the effect of a predictor than a p-value is. 95% space around predictor.
- Judge outliers and influential points w/ Cook statistic. Examine leverage w/ half-norm plot. leverage \neq influence

MODEL SELECTION

- We can transform variables w/ things like logs if it seems there is not a linear relationship between a var. and response.
- But this can change interpretations!
- We mostly judge the model as a whole based on adjusted R^2 or AIC. We can remove variables that don't contribute.
Lower AIC = better model fit
increase in R^2 = improve variability explained
- Use Hypothesis testing to test 2 models against each other w/ an ANOVA test. Mostly we use this w/ nested models.
F test is stat. sign. then larger model is better than nested model.
- We can also have R run through a larger version of the model and pick out the best variables based on AIC (backward selection)
- Sometimes a normal linear model just isn't right, so we need to start looking at a way to expand on the linear model...

code ...
summary()
df\$newvar =
plot(newvar, df)
cor(data[, all new])

lm(formula, data)
summary(lm)

confint(mod)
cooks.distance()
halfnorm(
influence(
mod)\$hat

AICc()
anova(w, z,
test =)
step()

NEW

Time to expand on the original linear model w/ logistic regression! We will see it works much better for binary data...

LOGISTIC REGRESSION

- We notice w/ a normal lme w/ a binary response (y) that the linear regression will give us pred. outside [0, 1], what can we do to get a better model...

Logistic Regression...

$$\eta_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$$

η_i is a link function $\eta_i = g(\beta_i)$
 $\rightarrow P(y_i = 1) = p_i$
this looks like linear regression but $\neq y$

code ...
logit()
illogit()

glm(, family = binomial)

$$\eta = \log(p/(1-p)) \text{ aka } p = \frac{e^\eta}{1+e^\eta}$$

- We find the parameters (β_i) w/ the method of maximum likelihood ... $\ell(\beta)$
- But of course now we can't use our linear interpretations for this model

INTERPRETATION

- The most natural way to interpret logistic regression is odds \rightarrow an unbounded scale of probability...

3 to 1 on odds \rightarrow for every \$3 bet win \$1
 $3 \text{ to } 1, 3/1 = 0.75$
 $p = \frac{3}{3+1} = 75\%$ chance

3 to 1 against odds \rightarrow for every \$1 bet win \$3
 $1 \text{ to } 3, 1/3 = 0.25$
 $p = \frac{1}{1+3} = 25\%$ chance

- To translate odds to prob. and vice versa...
probability $= p = \frac{p}{1+p}$ odds $= o = \frac{p}{1-p}$
- You can see why this translates well to log regression...

$$\log(\text{odds}) = \log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots$$
$$\text{odds} = e^{\beta_0} e^{\beta_1 x_1} e^{\beta_2 x_2}$$

- Coef interp.** \rightarrow "a unit increase in x_i w/ others held constant, changes log odds of success by β_i and changes odds of success by a factor of e^{β_i} ."
- Relative risk** $\rightarrow \frac{\exp(\hat{\beta} \hat{x}_1)}{1 + \exp(\hat{\beta} \hat{x}_1)} \bigg/ \frac{\exp(\hat{\beta} \hat{x}_2)}{1 + \exp(\hat{\beta} \hat{x}_2)} = \text{illogit}(\hat{\beta} \hat{x}_1) / \text{illogit}(\hat{\beta} \hat{x}_2)$
 \hat{x} is the vector of values for the predictors, change 1 value between \hat{x}_1 and \hat{x}_2 holding the rest fixed.
"The case where $x_i = 1$ is $rr\%$ more likely to see success than the case where $x_i = 0$."
- Hypothesis Tests** \rightarrow use Chi sq. tests instead of F tests, otherwise works the same.
- Confidence Intervals** \rightarrow find w/ $\beta_i \pm 1.96 \cdot SE(\hat{\beta}_i)$ \rightarrow note: find conf.int first then convert from log odds
z* for .95

MODEL SELECTION

- We can't judge the model based on residuals, the raw residuals always seem to form 2 lines (one near 0 and one near 1) because y can only take 2 values. Instead we use Deviance and Deviance Residuals.
Null Dev. \leftarrow based on just intercept
Dev. of model \leftarrow smaller deviance = better fit
off likelihood ratio test
 \rightarrow a transformed res. we bin and get these as a test of fit.
- Can still use anova (w/ chi sq test) and AIC.
- Again we can use step() / backward elimination to choose variables.
- Other Measures of Goodness of Fit \rightarrow
Hosmer-Lemeshow Stat \rightarrow uses conf.int. to measure variation, if p-value $\leq .05$ then not a good fit.
Missclassification \rightarrow specificity and sensitivity. Test how many correct and incorrect in the model based on outcome in data. can conceal variation.

lib(gamlss)
 \rightarrow logitglm(y,
predict())