

adv. data analysis

Midterm 2

generalized linear models

8.1-8.2

basics

- We want to use non-linear models while still having the benefits of linear regression...

$$\eta = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

we use a link function

$$\eta = g(\mu)$$

- most glm's are made up of responses that are members of exponential family dist.
- almost anything can be a link function, but canonical link functions exist...

family	canonical link	var. func.
normal	$\eta = \mu$	1
Poisson	$\eta = \log \mu$	μ
binomial	$\eta = \log(\mu/(1-\mu))$	$\mu(1-\mu)$
gamma	$\eta = \mu^{-1}$	μ^2
inv. gauss.	$\eta = \mu^{-2}$	μ^3

- The distributions also have dispersion parameters ϕ . The Poisson and binomial have $\phi=1$.
- To estimate β 's we use an algorithm...

② we need an initial estimate $\hat{\mu}^0 = \text{fitted values}$

① then we solve for adjusted dependent variable

$$z_i = \eta_i^0 + (y_i - \eta_i^0) \frac{\partial \eta_i}{\partial \mu_i}$$

② Find weights... $w_i = \left(\frac{\partial \eta_i}{\partial \mu_i}\right)^2 V(\hat{\mu}_i^0)$

③ Fit linear model to estimate $\hat{\beta}^1$ and use this as the next iteration of $\eta \rightarrow \eta^1$

④ repeat until convergence

- This method is good, but the std. errors are wrong. Find with $\text{Var}(\hat{\beta}) = (X^T W X)^{-1} \hat{\phi}$ w/ W diag. weight matrix X = model matrix

8.3 hypothesis tests

- We can measure goodness of fit by comparing our model to the null model (model w/ no possible var. and interactions). We use an equation to find deviance which can be compared to a χ^2 dist. (w/ df = obs. - param.). $\leq .05$ imply good fit.
- We can also compare two nested models (χ^2 larger, w smaller) by subtracting their deviances and comparing it to a χ^2 dist. $D_w - D_a \rightarrow df = \text{df}_a - \text{df}_w$. $\leq .05$ imply bigger model better

8.4 diagnostics

- We can't measure residuals the same way we do for linear models, since variance of response is not uniform for most glm's. So we use Pearson Residuals or Deviance Residuals.
- We can use diagnostics to check Leverage & Influence

Unusual Points

we look at a QQ Plot and Cook's D

we look at a half-normal plot and test the deviance to test the dependence (very similar conclusion as tables)

- We can also use Fisher Test on our table (fisher.test()) this will return an odds ratio that is a measure of association between variables.
- We can also visually judge dependence w/ Mosaic Plots (mosaicplot(..., color=T, main=NULL, las=1))

5.1

count reg. & poisson

- Lets take a closer look at one kind of glm \rightarrow Poisson!

$$Y_i \sim \text{Pois}(\mu_i)$$

$$\log \mu_i = \eta_i = X_i^T \beta$$

Pois. dist. are typically skewed right.



- Poisson links work really well for unbounded counts (0, 1, 2, ...). Counts w/ small prob. of success. Events where time between them is not exp. dist.
- For Poisson, sum is also Poisson and indep.
- Poisson is a type of count reg. model.

5.2-5.5

Poisson variations

- Dispersed Poisson \rightarrow modification of Poisson to allow more var. in resp. we let Y 's rate λ also be a RV...

$$Y \sim \text{Pois}(\lambda) \quad \lambda \sim \text{Gamma}(\alpha, \beta)$$

$$E[Y] = \mu \quad \text{Var}(Y) = \mu(1 + \phi) / \phi$$

in R \rightarrow glm(..., family=quasipoisson, ...) $\phi > 1 \rightarrow$ overdispersion $\phi < 1 \rightarrow$ underdispersion (models w/ overdispersion should be compared w/ F test rather than χ^2)

- Rate Models \rightarrow occurrences may depend on a size variable that determines # of opportunities for occurrences (ie # of goals scored depends on time played), so we use a rate model where covariate is fixed at 1.

$$\text{in R} \rightarrow \text{glm}(\dots \sim \text{offset}(\dots))$$

- Hurdle Model \rightarrow what if we have way more 0's in our data than a Poisson or similar model would predict? Well we have a couple ways to deal w/ that. One way is to imagine there is a hurdle that an observer has to get over, that most people never get over, then when they get it it's more likely they'll get more (ie first book published, first robbery)

$$P(Y=0) = f_1(0) \quad P(Y=j) = \frac{1-f_1(0)}{1-f_1(0)} f_2(j), j > 0$$

$$\text{in R} \rightarrow \text{hurdle}(y \sim x, \text{data})$$

- Zero Inflated Poisson Model \rightarrow same problem as above, data has lots of 0's. Idea of ZIP model is only a certain population will have >0. (compare w/ χ^2 dist.)

$$P(Y=0) = \phi + (1-\phi)f(0) \quad P(Y=j) = (1-\phi)f(j), j > 0$$

$$\text{in R} \rightarrow \text{zeroinfl}(y \sim x, \text{data})$$

- Negative Binomial \rightarrow series of indep. trials each w/ prob. of success = p . Z = # of trials until k th success. This is a neg. binomial dist. is we let $Y = Z - k$ $\rightarrow p = (1-p)^k$

$$P(Y=y) = \binom{y+k-1}{k-1} p^k (1-p)^y$$

we end up w/ a glm w/...

$$E[Y] = \mu = k/p$$

$$\text{Var}(Y) = k/p + k/p^2 = \mu + \mu^2/k$$

$$\eta = \log \mu = \log k - \log p$$

$$\text{in R} \rightarrow \text{glm}(\dots, \text{family} = \text{negative.binomial}(\dots))$$

$$\text{glm.nb}()$$

6.1-6.2

two-by-two tables

- Contingency Table \rightarrow a way to show cross classified data on 2 or more categorical variables. Variables can be nominal (no natural order) or ordinal (has a natural order). A great way to start observing data and dependence.

- A two-by-two table has two variables...

- We can judge indep. w/ a χ^2 test on the table.

- We can also fit a Poisson w/ count as the response and test the deviance to test the dependence (very similar conclusion as tables)

- We can also use Fisher Test on our table (fisher.test()) this will return an odds ratio that is a measure of association between variables.

- We can also visually judge dependence w/ Mosaic Plots (mosaicplot(..., color=T, main=NULL, las=1))

var1	outcome 1	outcome 2	total
var2	0.1	0.2	
	count	count	
	total		

6.4 matched pairs

- w/ matched pairs we build a contingency table where we observe one measure on two matched objects.
- We can also build a glm and then build a matched pairs contingency table of residuals.
- These matched pair contingency tables can be really useful for testing symmetry. We can use a quasi-symmetry model w/ a normal dist. with $\text{mcmc}()$ test w/ Chi dist. If the test results show there is evidence of lack of marginal homogeneity.

student	A	B	C	D
teacher				
	count	count	count	count

6.5 3-way contingency table

- What if we produce a two-way contingency table that shows/doesn't show a relationship we don't/ do expect to be there? Well maybe we need to condition upon a third variable, this is how we make a 3-way contingency table

- marginal association \rightarrow association between 2 var.

- conditional association \rightarrow association obs. within each category in the 3rd var.

- Simpson's paradox \rightarrow when the marginal association and the conditional association lead to different conclusions.

- We can also use the Mantel-Haenszel test which is designed to test indep. in 2x2 tables across K strata. Output is a approx. dist. on χ^2 dist. If stat. sign \rightarrow association between vars. (mantelhaen.test(ct, exact=T))

- We can also test association by comparing fit of 4 types of models...

- Mutual Independence \rightarrow This model will be best if all 3 var. are indep.

- Joint Independence \rightarrow This model will be best if 2 var. are dep. but indep. of the 3rd. in R \rightarrow glm(count ~ var1 + var2 + var3, family=poisson)

- Conditional Independence \rightarrow This model will be best if the first and second var. are indep. given the third. in R \rightarrow glm(count ~ var1 + var2 + var3, family=poisson)

- Uniform Association \rightarrow This model will be best if each of the variables have some association with both of the other vars. in R \rightarrow glm(count ~ var1 + var2 + var3, family=poisson)

10.1-10.3 basics

random effects

defining the parameter: fixed effect

- mixed effect \rightarrow when we have random effect and fixed effects in one model...

- simplest design of random effect model...

$$y_{ij} = \mu + \alpha_i + \epsilon_{ij}$$

Intraclass correlation $\rho = \frac{\sigma_\alpha^2}{\sigma_\alpha^2 + \sigma_\epsilon^2}$

$$\rho = 0 \rightarrow \text{no var}$$

$$\rho = 1 \rightarrow \text{lots of var}$$

variances were originally found using MSE and MSA, but this is mathematically very complicated and can return negative values, so now we use MLE or the less biased REML

$$y_{ijk} = \mu + \tau_i + \gamma_j + \epsilon_{ijk}$$

fixed random error

$$y_{ij} = \mu + \alpha_i + \epsilon_{ij}, \quad \alpha_i \sim N(0, \sigma_\alpha^2), \quad \epsilon_{ij} \sim N(0, \sigma_\epsilon^2)$$

$$\hat{\alpha} = DZ^T V^{-1} (y - X\beta)$$

- in R \rightarrow lmer4::lmer(response ~ 1 + (1|re), ...)
- if our model it doesn't make sense to try and estimate the RE, instead we estimate its func. ...
- to test two RE models against each other create a null model w/ component your testing @ test criterion ratio statistics (CRMs) of nested models against each other it should be approx. χ^2 dist unless parameters of null model are on boundary of param. space.
- in R \rightarrow RLRsim::exactLRT(mod, nullmod), for one mod \rightarrow RLRsim::exactRLRT(mod)
- For balanced data (same # of obs. in all levels) we can use sum of squares in ANOVA decomposition for hypothesis testing.
- For fixed effects we should test w/ F-test.
- We can also use parametric bootstrap to find p-values for LRT. To do this we generate data under null mod using fitted param. estimates. we repeat this many times to find p-values. in R \rightarrow confint(mod, method="boot")
- We can use AIC to compare many models.

10.4 prediction

- best linear unbiased predictors \rightarrow w/ operators we can derive BLUPs to predict new values.

- in R \rightarrow predict(mod, re.form = ~0)

$$\text{predict(mod, newdata = dF(pseudo = "a"))}$$

- Func. does not give us intervals or std. errors cause matrix is hard.

- For simple models we can bootstrap to get pred. intv.

$$\text{in R} \rightarrow \text{lme4::bootMer()}$$

10.5 diagnostics

- Assumptions for our mixed lin. model are very similar to linear models. Diagnostics...

- mixed effect models have more than one kind of fitted residual - dependent on estimated random effect

- random effects models are extra sensitive to outliers

- you can look at normality at group levels.

10.6 blocks as re's

- blocks \rightarrow properties of the experimental units, groupings, typically we are not interested in the blocks effects specifically, but we need to account for their effects. So we use RE's!

- in R \rightarrow lmer(response ~ fixed + (1|block))

- we can then test our mod to see if blocking is helpful or if we should stick to fixed.

