

# Ordinal and categorical variables

Ben Bolker

October 29, 2018



Licensed under the Creative Commons attribution-noncommercial license (<http://creativecommons.org/licenses/by-nc/3.0/>). Please share & remix noncommercially, mentioning its origin.

```
library(ggplot2)
theme_set(theme_bw())
library(scales)    ## squish
library(gridExtra) ## grid.arrange()
library(nnet)      ## multinom()
library(plyr)
library(reshape2)
library(faraway)  ## data
library(RColorBrewer) ## nice colours
```

## Ordered predictors

(Not the primary topic but feel like I ought to mention it.)

*Ordered* factors are the case where there is a natural ordering to the responses.

This is (confusingly) different from the usual unordered-factor case, where the order of the levels is still used (1) to determine the order of the categories for high-level plotting and (2) to determine contrasts (which level is the baseline).

Options for dealing with ordered (or otherwise messy) predictors:

- assume linearity (equal differences in predicted values between successive levels); convert the factor back to numeric
- use `contr.sdif` from the MASS package
- use `ordered` instead of `factor`
- use `cut`, `cut_number`, `cut_interval` to convert continuous predictors to factors

Don't snoop!

Ordered factors: contrasts

```
ff <- function(n) {
  cc <- zapsmall(contr.poly(n))
  ## polynomials are scaled so that sum(c^2)=1; prettify
```

```

sign(cc)*MASS::fractions(cc^2)
}
ff(3)

##      .L      .Q
## [1,] -1/2  1/6
## [2,]   0 -2/3
## [3,]  1/2  1/6

ff(5)

##      .L      .Q      .C      ^4
## [1,] -2/5  2/7 -1/10  1/70
## [2,] -1/10 -1/14  2/5 -8/35
## [3,]   0 -2/7   0 18/35
## [4,]  1/10 -1/14 -2/5 -8/35
## [5,]  2/5  2/7  1/10  1/70

```

No increase in parsimony over treatment contrasts, but improved interpretability. Linear, quadratic models are nested within the ordered-factor model.

### *Categorical responses*

We can either model these as *multinomial*, or as conditional Poisson (i.e., if we take a set of independent Poisson deviates  $x_i$  they are equivalent to a multinomial sample out of  $\sum_i x_i$  with  $p_i = \lambda_i / \sum \lambda_i$ ).

In either case we have to define

$$\mathcal{L} \propto \sum_i N_i \log p_i$$

Multinomial distributions are also conditionally *binomial* if we only want to consider one category vs. all the others ...

Here's a data set on US political preferences:

10 variable subset of the 1996 American National Election Study. Missing values and "don't know" responses have been listwise deleted. Respondents expressing a voting preference other than Clinton or Dole have been removed.

```

library(faraway)
data(nes96)
nn <- subset(nes96, select=c(PID, age, educ, income))
## summary(nn)

```

For simplicity, lump party identifications into three categories:

```
nn$party <- factor(sub("str|weak|ind"), "", nn$PID)
```

Get a numeric value for the average income in a category:

```
## income breakpoints
incbrks <- c(0,
             unique(readr::parse_number(nn$income)),
             125)
## take average of breakpoints
inc_avg <- (incbrks[-1]+incbrks[-length(incbrks)])/2
```

Name the vector:

```
names(inc_avg) <- levels(nn$income)
```

Now something like `inc_avg["$3K-$5K"]` would work ...

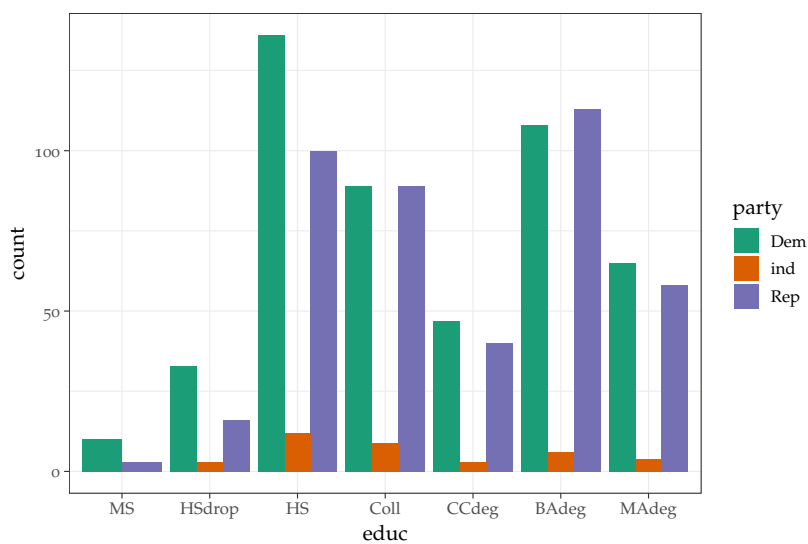
Numeric versions of variables:

```
nn <- transform(nn, nincome=inc_avg[nn$income],
               neduc=as.numeric(educ))
```

Categorical versions of variables:

```
cincome <- cut_number(nn$income, 7)
cage <- cut_number(nn$age, 7)
cdata <- with(nn, data.frame(party, educ, cincome, cage))
```

```
(ggplot(cdata, aes(x=educ, fill=party))
 +geom_bar(position="dodge")+
 scale_fill_brewer(palette="Dark2")
)
```



Rescale data, get proportions of parties by education and party:

```
tt <- with(nn, table(educ, party))
tot <- rowSums(tt)
tt <- sweep(tt, 1, tot, "/")
tt <- data.frame(tt, tot) ## automatically "melted"

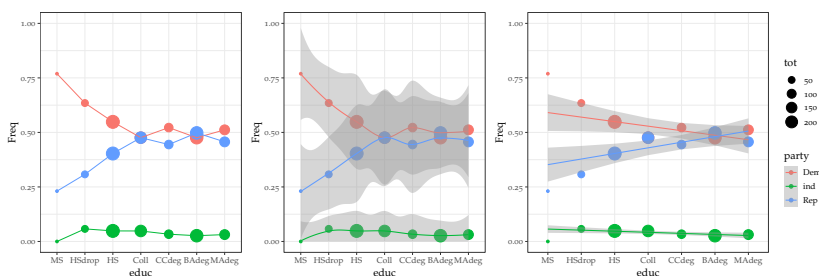
## Warning in data.frame(tt, tot): row names were found
## from a short variable and have been discarded

tt$neduc <- as.numeric(tt$educ)
```

Three ways to plot the results:

```
g1 <- ggplot(tt, aes(x=educ, y=Freq,
                    colour=party))+
  geom_point(aes(size=tot))+
  scale_y_continuous(limits=c(0, 1), oob=squish)
library(gridExtra)
g1A <- g1+geom_line(aes(group=party))+theme(legend.position="none")
g1B <- g1+geom_smooth(aes(x=as.numeric(educ)), method="loess")+
  theme(legend.position="none")
g1C <- g1 + geom_smooth(aes(group=party, weight=tot),
                       method="glm",
                       method.args=list(family=binomial))
```

```
grid.arrange(g1A, g1B, g1C, ncol=3, widths=unit(c(1, 1, 1.4)), units="null")
```



### Multinomial responses

Non-ordered categorical responses. We have to predict the effects of *each* predictor on *each* response.

```
library(nnet)
m1 <- multinom(party ~ age+educ+nincome, data=nn)
summary(m1)
```

What do the parameters mean? e.g. the first element of the intercept vector is the log-odds of the probability of being Independent vs. Democrat in the baseline level; the second is the log-odds of the probability of being Republican vs Democrat in the baseline level.

Test this:

```
z <- data.frame(party=c("Democrat", "Democrat", "Ind", "Republican"))
```

We take the coefficient (the intercept), compute the logistic function (`plogis`), and compute the fractional equivalent.

```
MASS::fractions(plogis(coef(multinom(party~1, data=z))))
## # weights:  6 (2 variable)
## initial value 4.394449
## final value 4.158883
## converged
##           (Intercept)
## Ind           1/3
## Republican 1/3
```

Both of the probabilities are  $1/3$ :

number of independents / [number of ind + number of dem] =  $1/3$

number of republicans / [number of R + number of D] =  $1/3$

Change the reference level to Independent:

```
z$party <- relevel(z$party, "Ind")
```

```
MASS::fractions(plogis(coef(multinom(party~1, data=z))))
## # weights:  6 (2 variable)
## initial value 4.394449
## final value 4.158883
## converged
##           (Intercept)
## Democrat    2/3
## Republican  1/2
```

number of D / [number of I + number of D] =  $2/3$

number of R / [number of R + number of I] =  $1/2$

Fit with numeric rather than ordinal predictors:

```
m2 <- multinom(party ~ age+neduc+nincome, nn)
```

Without education at all:

```
m3 <- update(m2, ~.-neduc)
```

What do the parameters mean??

```
summary(m2)

## Call:
## multinom(formula = party ~ age + neduc + nincome, data = nn)
##
## Coefficients:
##      (Intercept)          age          neduc          nincome
## ind   -2.560991  0.002804454 -0.21395608  0.01686278
## Rep   -1.164684  0.007441529  0.01217699  0.01302126
##
## Std. Errors:
##      (Intercept)          age          neduc          nincome
## ind   0.7862200  0.010845152  0.12194267  0.005887065
## Rep   0.3121893  0.004199209  0.04666894  0.002441064
##
## Residual Deviance: 1521.778
## AIC: 1537.778
```

To the extent that the non-intercept parameters are similar between groups, this suggests that we might be able to get away with a proportional-odds model (see below).

Finding best AIC (smallest AIC is best;  $< 2\Delta\text{AIC}$  is a small difference;  $> 10\Delta\text{AIC}$  is a big difference).

```
trace <- TRUE ## I don't know why, but this prevents an error
(dd <- drop1(m1)) ## test="Chisq" is ignored
```

Compared to best model:

```
delta_AIC <- dd$AIC-min(dd$AIC)
names(delta_AIC) <- rownames(dd)
round(delta_AIC,2)

## <none>    age    educ nincome
##   10.30   11.75    0.00   39.10
```

We can't get  $p$  values from `drop1`, but we can do likelihood ratio tests:

```
anova(m1,m2,m3) ## education: test categorical vs linear vs null model
```

```
## Likelihood ratio tests of Multinomial Models
```

```
##
```

```
## Response: party
```

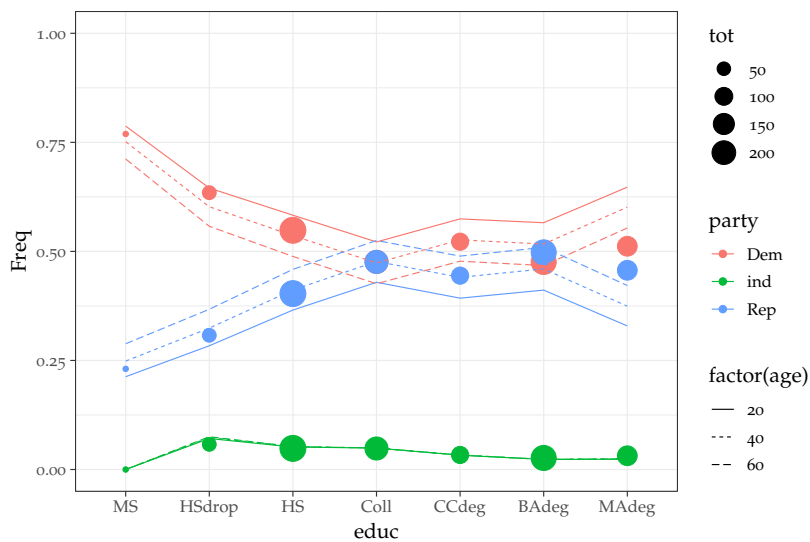
```
##           Model Resid. df Resid. Dev  Test   Df LR stat.
## 1      age + nincme      1882   1525.317
## 2 age + neduc + nincme      1880   1521.778 1 vs 2    2  3.539461
## 3 age + educ + nincme      1870   1511.612 2 vs 3   10 10.165237
##      Pr(Chi)
## 1
## 2 0.1703789
## 3 0.4261181
```

```
predict.multinom...
```

```
prepdata <- data.frame(nincome=mean(nn$nincome),
                      expand.grid(age=c(20,40,60),educ=levels(nn$educ)))
probs <- predict(m1,newdata=prepdata,type="probs")
```

```
prepdata <- data.frame(prepdata,probs)
predmelt <- rename(melt(prepdata,id.vars=1:3),
                  c(variable="party",value="Freq"))
```

```
g1 + geom_line(aes(group=interaction(party,age),
                                lty=factor(age)),data=predmelt)
```



What else can I do with a multinomial fit?

```

methods(class="multinom")

## [1] add1      anova      coef      confint   drop1
## [6] extractAIC logLik    model.frame predict   print
## [11] summary    vcov
## see '?methods' for accessing help and source code

```

(Sometimes there are starred functions, which are hidden inside packages: e.g. to look at them you would need `nnet:::drop1.multinom`.)

### *Ordinal responses*

Multiple categorical levels of response, but ordered.

*Proportional odds* (or *proportional probability*, depending on link function).

`polr` function from the MASS package; also the ordinal package.

```

library(MASS)
p1 <- polr(party ~ age+educ+nincome, nn)
drop1(p1, test="Chisq")

## Single term deletions
##
## Model:
## party ~ age + educ + nincome
##           Df    AIC      LRT Pr(>Chi)
## <none>      1538.8
## age         1 1542.0  5.2199  0.02233 *
## educ        6 1535.1  8.3304  0.21488
## nincome     1 1566.2 29.4579 5.715e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

p2 <- polr(party ~ age+neduc+nincome, nn)
drop1(p2, test="Chisq")

## Single term deletions
##
## Model:
## party ~ age + neduc + nincome
##           Df    AIC      LRT Pr(>Chi)
## <none>      1537.1
## age         1 1538.0  2.9736  0.08463 .
## neduc       1 1535.1  0.0484  0.82593
## nincome     1 1564.3 29.2493 6.364e-08 ***
## ---

```



```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Note correlation among parameters:

```
round(cov2cor(vcov(p2)),2)

##
## Re-fitting to get Hessian
##      age neduc nincome Dem|ind ind|Rep
## age      1.00  0.13   0.03   0.74   0.74
## neduc    0.13  1.00  -0.38   0.63   0.63
## nincome  0.03 -0.38   1.00   0.12   0.12
## Dem|ind  0.74  0.63   0.12   1.00   1.00
## ind|Rep  0.74  0.63   0.12   1.00   1.00
```

Or using the ordinal package (more flexible/newer):

```
library(ordinal)
p3 <- clm(party ~ age+educ+nincome, data=nn)
coef(p1)

##      age      educ.L      educ.Q      educ.C      educ^4
## 0.009522628 0.573552066 -0.742893138 0.069254713 -0.044684004
##      educ^5      educ^6      nincome
## -0.081227547 -0.138260492 0.012412721

coef(p3)

##      Dem|ind      ind|Rep      age      educ.L      educ.Q
## 1.268256953 1.433490808 0.009522676 0.573548506 -0.742897303
##      educ.C      educ^4      educ^5      educ^6      nincome
## 0.069250944 -0.044695028 -0.081200313 -0.138256927 0.012412867
```

Comparing log-likelihoods and AICs between multinomial and proportional-odds models:

```
logLik(m1)

## 'log Lik.' -755.8062 (df=18)

logLik(p1)

## 'log Lik.' -759.3974 (df=10)

AIC(m1)

## [1] 1547.612
```

```

AIC(p1)
## [1] 1538.795

library(bbmle) ## prettier AIC tables

## Loading required package: stats4
##
## Attaching package: 'bbmle'
## The following object is masked from 'package:ordinal':
##
## slice

AICtab(m1,p1)

##      dAIC df
## p1  0.0 10
## m1  8.8 18

```

Alternative test of non-proportionality (for individual predictor variables):

```

p4 <- update(p3, nominal= ~age)
anova(p3, p4)

## Likelihood ratio tests of cumulative link models:
##
##      formula:                nominal: link: threshold:
## p3 party ~ age + educ + nincome ~1      logit flexible
## p4 party ~ age + educ + nincome ~age     logit flexible
##
##      no.par   AIC  logLik LR.stat df Pr(>Chisq)
## p3         10 1538.8 -759.40
## p4         11 1540.8 -759.39 0.0189 1      0.8906

```