

Statistical Modeling

BANA 7042

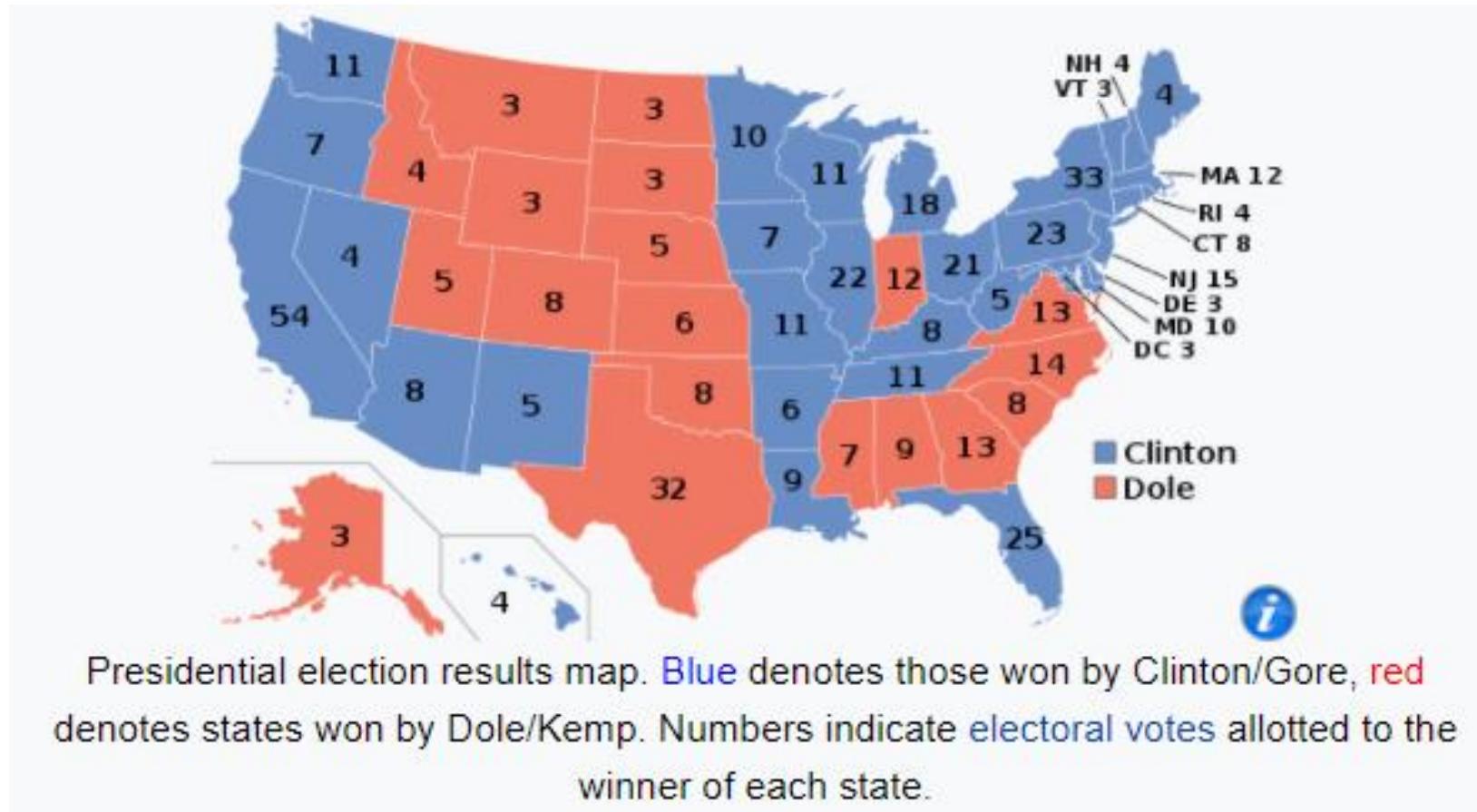
Lecture 5: Modeling multinomial data

Dr. Dungang Liu

1996 US Presidential Election







Presidential election results map. Blue denotes those won by Clinton/Gore, red denotes states won by Dole/Kemp. Numbers indicate electoral votes allotted to the winner of each state.

1996 American National Election Study

The data set “nes96” is a 10 variable subset of the 1996 American National Election Study. It contains information of 944 voters.

```
library("faraway")
data(nes96)
str(nes96)
?nes96
summary(nes96)
```

The data structure

```
> str(nes96)
'data.frame': 944 obs. of 10 variables:
 $ popul : int 0 190 31 83 640 110 100 31 180 2800 ...
 $ TVnews: int 7 1 7 4 7 3 7 1 7 0 ...
 $ selfLR: Ord.factor w/ 7 levels "extLib"<"Lib"<...: 7 3 2 3 5 3 5 5 5 4 3 ...
 $ ClinLR: Ord.factor w/ 7 levels "extLib"<"Lib"<...: 1 3 2 4 6 4 6 4 6 3 ...
 $ DoleLR: Ord.factor w/ 7 levels "extLib"<"Lib"<...: 6 5 6 5 4 6 4 5 3 7 ...
 $ PID    : Ord.factor w/ 7 levels "strDem"<"weakDem"<...: 7 2 2 2 1 2 2 5 4 1
...
$ age    : int 36 20 24 28 68 21 77 21 31 39 ...
$ educ   : Ord.factor w/ 7 levels "MS"<"HSdrop"<...: 3 4 6 6 6 4 4 4 4 3 ...
$ income: Ord.factor w/ 24 levels "$3Kminus"<"$3K-$5K"<...: 1 1 1 1 1 1 1 1 1 1
...
$ vote   : Factor w/ 2 levels "Clinton", "Dole": 2 1 1 1 1 1 1 1 1 1 ...
```

Variable dictionary

popul

population of respondent's location in 1000s of people

TVnews

days in the past week spent watching news on TV

selfLR

Left-Right self-placement of respondent: an ordered factor with levels extremely liberal, extLib < liberal, Lib < slightly liberal, sliLib < moderate, Mod < slightly conservative, sliCon < conservative, Con < extremely conservative, extCon

ClinLR

Left-Right placement of Bill Clinton (same scale as selfLR): an ordered factor with levels extLib < Lib < sliLib < Mod < sliCon < Con < extCon

DoleLR

Left-Right placement of Bob Dole (same scale as selfLR): an ordered factor with levels extLib < Lib < sliLib < Mod < sliCon < Con < extCon

PID

Party identification: an ordered factor with levels strong Democrat, strDem < weak Democrat, weakDem < independent Democrat, indDem < independent independentind < independent Republican, indRep < weak Republican, weakRep < strong Republican, strRep

Variable dictionary (cont'd)

age

Respondent's age in years

educ

Respondent's education: an ordered factor with levels 8 years or less, MS < high school dropout, HSdrop < high school diploma or GED, HS < some College, Coll < Community or junior College degree, CCdeg < BA degree, BAdeg < postgraduate degree, MAdeg

income

Respondent's family income: an ordered factor with levels \$3Kminus < \$3K-\$5K < \$5K-\$7K < \$7K-\$9K < \$9K-\$10K < \$10K-\$11K < \$11K-\$12K < \$12K-\$13K < \$13K-\$14K < \$14K-\$15K < \$15K-\$17K < \$17K-\$20K < \$20K-\$22K < \$22K-\$25K < \$25K-\$30K < \$30K-\$35K < \$35K-\$40K < \$40K-\$45K < \$45K-\$50K < \$50K-\$60K < \$60K-\$75K < \$75K-\$90K < \$90K-\$105K < \$105Kplus

vote

Expected vote in 1996 presidential election: a factor with levels Clinton and Dole

Source

Sapiro, Virginia, Steven J. Rosenstone, Donald R. Kinder, Warren E. Miller, and the National Election Studies. AMERICAN NATIONAL ELECTION STUDIES, 1992-1997: COMBINED FILE [Computer file]. 2nd ICPSR version. Ann Arbor, MI: University of Michigan, Center for Political Studies [producer], 1999. Ann Arbor, MI: Inter-university Consortium for Political and Social Research [distributor], 1999.

Summary statistics

```
> summary(nes96)
```

popul	TVnews	selfLR	ClinLR	DoleLR
Min. : 0.0	Min. :0.000	extLib: 16	extLib:109	extLib: 13
1st Qu.: 1.0	1st Qu.:1.000	Lib :103	Lib :317	Lib : 31
Median : 22.0	Median :3.000	sliLib:147	sliLib:236	sliLib: 43
Mean : 306.4	Mean :3.728	Mod :256	Mod :160	Mod : 87
3rd Qu.: 110.0	3rd Qu.:7.000	sliCon:170	sliCon: 67	sliCon:195
Max. :7300.0	Max. :7.000	Con :218	Con : 36	Con :460
		extCon: 34	extCon: 19	extCon:115
PID	age	educ	income	vote
strDem :200	Min. :19.00	MS : 13	\$60K-\$75K:103	Clinton:551
weakDem:180	1st Qu.:34.00	HSdrop: 52	\$50K-\$60K:100	Dole :393
indDem :108	Median :44.00	HS :248	\$30K-\$35K: 70	
indind : 37	Mean :47.04	Coll :187	\$25K-\$30K: 68	
indRep : 94	3rd Qu.:58.00	CCdeg : 90	\$105Kplus: 68	
weakRep:150	Max. :91.00	BAdeg :227	\$35K-\$40K: 62	
strRep :175		MAdeg :127	(Other) :473	

What question could be answered?

- What are the underlying factors that would influence an individual's
 - Vote? --- A binary variable
 - Party Identification? --- A multinomial variable
- Are there any common factors?
- Generally speaking, Democrats voted for Clinton, and Republican voted for Dole. --- Is this true?
- Let us focus on 5 variables: "vote", "PID", "income", "educ", and "age".

Initial data processing

- Aggregate the levels of “PID” such that we have 3 categories: “Democrat”, “Independent”, and “Republican”.

```
party<-nes96$PID  
levels(party)<-c("Democrat", "Democrat", "Independent", "Independent", "Independent",  
"Independent", "Independent", "Republican", "Republican")
```

- Convert “income” to a numeric variable by taking the midpoint of each income range.

```
inca<-  
c(1.5, 4, 6, 8, 9.5, 10.5, 11.5, 12.5, 13.5, 14.5, 16, 18.5, 21, 23.5, 27.5, 32.5, 37.5, 42.5, 47.5,  
55, 67.5, 82.5, 97.5, 115)  
  
income<-inca[unclass(nes96$income) ]
```

Summary of the processed data

```
rnes96<-data.frame(party, income, education=nes96$educ, age=nes96$age)  
summary(rnes96)
```

	party	income	education	age
Democrat	:380	Min. : 1.50	MS : 13	Min. :19.00
Independent	:239	1st Qu.: 23.50	HSdrop: 52	1st Qu.:34.00
Republican	:325	Median : 37.50	HS :248	Median :44.00
		Mean : 46.58	Coll :187	Mean :47.04
		3rd Qu.: 67.50	CCdeg : 90	3rd Qu.:58.00
		Max. :115.00	BAdeg :227	Max. :91.00
			MAdeg :127	

Visualization: Party versus Education

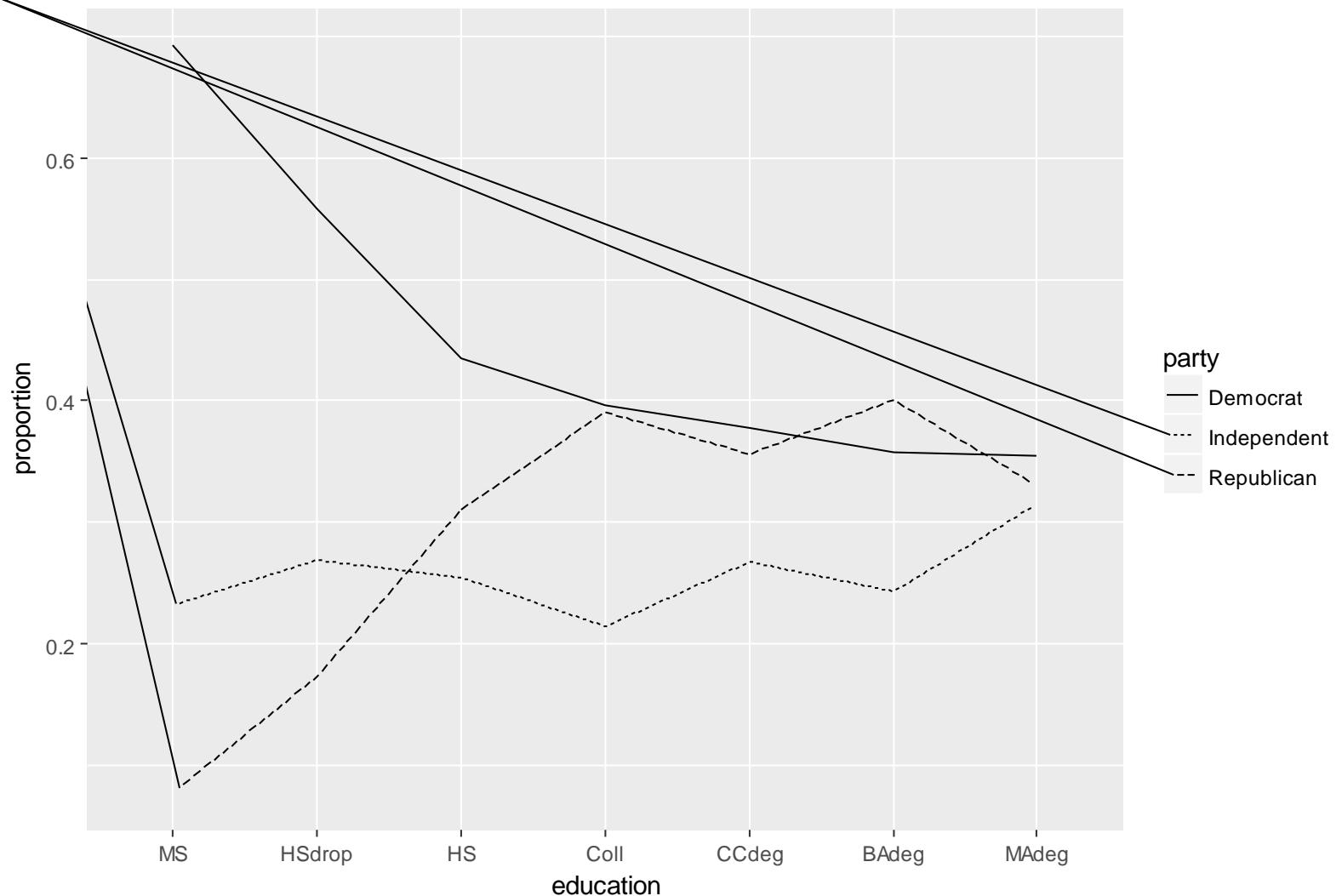
```
install.packages("dplyr")
```

```
library("dplyr")
```

```
egp<-group_by(rnes96, education, party) %>%  
  summarise(count=n()) %>% group_by(education) %>%  
  mutate(etotal=sum(count), proportion=count/etotal)
```

```
library("ggplot2")
```

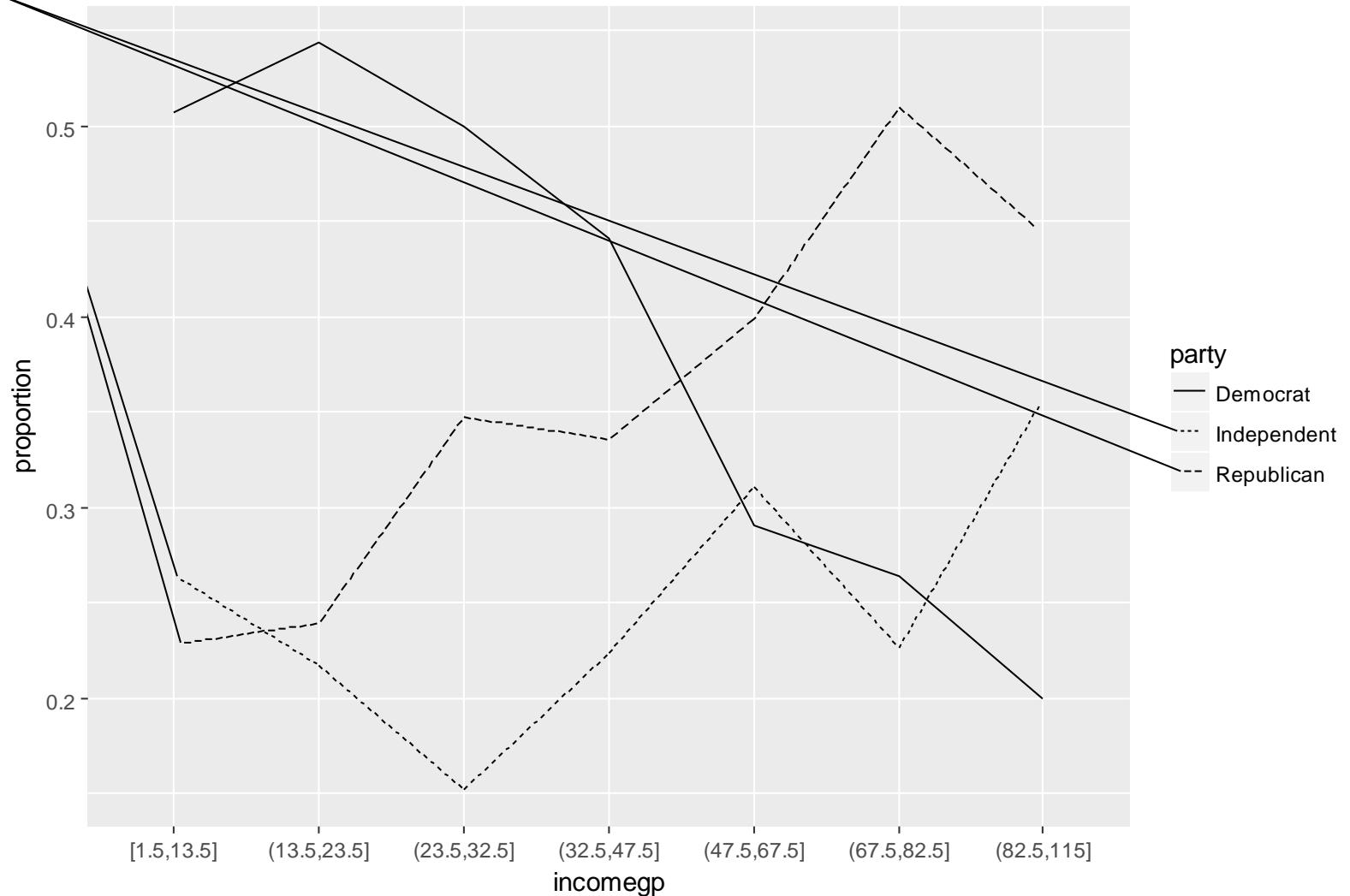
```
ggplot(egp, aes(x=education, y=proportion,  
group=party, linetype=party)) + geom_line()
```



Visualization: Party versus Income

```
igp<-mutate(rnes96, incomegp=cut_number(income,7)) %>%
group_by(incomegp, party) %>% summarise(count=n()) %>%
group_by(incomegp) %>% mutate(etotal=sum(count),
proportion=count/etotal)
```

```
ggplot(igp, aes(x=incomegp, y=proportion, group=party,
linetype=party)) + geom_line()
```



Are the trends real or by chance?

- As the education level increases, the proportion of democrats decreases and that of republicans increases.
- As the income increases, the proportion of democrats decreases and that of republicans increases.
- Can we build a model to examine the plausible association?

Modeling ordinal data

From “binary” to “ordinal multinomial”

- Unlike the variable “vote” which is binary, “party” has three categories and the categories have an “order”.
- How can we extend the idea of modeling “vote” to that of modeling “party”?
- Suppose we use a logistic regression to model “vote”. How do we interpret the model using a **latent variable**?

The latent variable for a binary variable

- Suppose there is a latent variable Z for the binary variable “vote”.

$$Z \sim \text{logistic}(\mu = -\beta x, \sigma^2 = 1)$$

$$\text{“vote”} = Y = \begin{cases} 0 & \text{if } Z \leq c, \\ 1 & \text{otherwise.} \end{cases}$$

$$\begin{aligned}\Pr\{Y = 0\} &= \Pr\{Z \leq c\} = \Pr\{\text{logistic}(\mu, 1) \leq c\} \\ &= \Pr\{\text{logistic}(0, 1) \leq c - \mu\} = F(c + \beta x)\end{aligned}$$

- This is one way to establish a logistic regression model.

The latent variable for an ordinal variable

- Suppose there is a latent variable Z for the ordinal variable “party”.

$$Z \sim \text{logistic}(\mu = -\beta x, \sigma^2 = 1)$$

$$\text{“party”} = Y = \begin{cases} 0 & \text{if } Z \leq c_1, \\ 1 & \text{if } Z \leq c_2, \\ 2 & \text{otherwise.} \end{cases}$$

$$\begin{aligned} \Pr\{Y \leq 0\} &= \Pr\{Z \leq c_1\} = \Pr\{\text{logistic}(\mu, 1) \leq c_1\} \\ &= \Pr\{\text{logistic}(0, 1) \leq c_1 - \mu\} = F(c_1 + \beta x) \end{aligned}$$

$$\begin{aligned} \Pr\{Y \leq 1\} &= \Pr\{Z \leq c_2\} = \Pr\{\text{logistic}(\mu, 1) \leq c_2\} \\ &= \Pr\{\text{logistic}(0, 1) \leq c_2 - \mu\} = F(c_2 + \beta x) \end{aligned}$$

Fit a proportional odds model

```
library(VGAM)  
nmod<-vglm(party ~ age + education + income,  
family=cumulative(parallel=TRUE), rnes96)  
summary(nmod)
```

```
vglm(formula = party ~ age + education + income, family = cumulative(parallel = TRUE),  
      data = rnes96)
```

Pearson residuals:

	Min	1Q	Median	3Q	Max
logit(P[Y<=1])	-2.059	-0.9564	-0.3698	1.0152	1.974
logit(P[Y<=2])	-2.804	-1.1156	0.3450	0.9138	1.707

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept):1	0.644871	0.244989	2.632	0.00848 **
(Intercept):2	1.737347	0.250638	6.932	4.16e-12 ***
age	-0.005775	0.003878	-1.489	0.13649
education.L	-0.724109	0.392387	-1.845	0.06498 .
education.Q	0.781386	0.358405	2.180	0.02924 *
education.C	-0.040190	0.297745	-0.135	0.89263
education^4	0.019942	0.236013	0.084	0.93266
education^5	0.079399	0.193456	0.410	0.68149
education^6	0.061111	0.157335	0.388	0.69771
income	-0.012739	0.002186	-5.826	5.68e-09 ***

Number of linear predictors: 2

Names of linear predictors: logit(P[Y<=1]), logit(P[Y<=2])

Residual deviance: 1984.211 on 1878 degrees of freedom

Log-likelihood: -992.1056 on 1878 degrees of freedom

Exponentiated coefficients:

age	education.L	education.Q	education.C	education^4	education^5	education^6	income
0.9942419	0.4847561	2.1844985	0.9606074	1.0201419	1.0826362	1.0630167	0.9873423

Transform “education” to years

```
edu.years<-c(9,10.5,12,13,14,16,18)
education<-edu.years[unclass(nes96$educ) ]
rnes96$education<-education

nmod<-vglm(party ~ age + education + income,
family=cumulative(parallel=TRUE), rnes96)
summary(nmod)
```

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept):1	0.506585	0.441597	1.147	0.251312
(Intercept):2	1.589747	0.444444	3.577	0.000348 ***
age	-0.002902	0.003752	-0.773	0.439242
education	-0.012316	0.028803	-0.428	0.668931
income	-0.012877	0.002181	-5.904	3.55e-09 ***

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

Number of linear predictors: 2

Names of linear predictors: logit(P[Y<=1]), logit(P[Y<=2])

Residual deviance: 1994.642 on 1883 degrees of freedom

Log-likelihood: -997.3211 on 1883 degrees of freedom

Number of iterations: 4

No Hauck-Donner effect found in any of the estimates

Exponentiated coefficients:

age	education	income
0.9971018	0.9877591	0.9872058

Proportional odds model

- The proportional odds model for “party”

$$\Pr\{Y \leq 1\} = F(\alpha_1 + \beta x) \text{ and } \Pr\{Y \leq 2\} = F(\alpha_2 + \beta x)$$

- Where F is the Inverse of the logit function.
- Why do we call it “odds” model?
- Why do we call it “proportional” model?

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)	
(Intercept) : 1	0.506585	0.441597	1.147	0.251312	
(Intercept) : 2	1.589747	0.444444	3.577	0.000348	***
age	-0.002902	0.003752	-0.773	0.439242	
education	-0.012316	0.028803	-0.428	0.668931	
income	-0.012877	0.002181	-5.904	3.55e-09	***

Exponentiated coefficients:

age	education	income
0.9971018	0.9877591	0.9872058

Non-proportional odds model

- The proportional odds model for “party”

$$\Pr\{Y \leq 1\} = F(\alpha_1 + \beta_1 x) \text{ and } \Pr\{Y \leq 2\} = F(\alpha_2 + \beta_2 x)$$

where F is the Inverse of the logit function.

- The beta coefficients rely on the categories.

```
nmod.2<-vglm(party ~ age + education + income,  
family=cumulative(parallel=FALSE), rnes96)  
summary(nmod.2)
```

```
vglm(formula = party ~ age + education + income, family = cumulative(parallel = FALSE),  
      data = rnes96)
```

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept):1	0.527897	0.490515	1.076	0.28183
(Intercept):2	1.511425	0.503245	3.003	0.00267 **
age:1	-0.002214	0.004101	-0.540	0.58931
age:2	-0.003784	0.004323	-0.875	0.38144
education:1	-0.007138	0.032075	-0.223	0.82389
education:2	-0.013991	0.032699	-0.428	0.66875
income:1	-0.016077	0.002550	-6.304	2.90e-10 ***
income:2	-0.010231	0.002408	-4.249	2.15e-05 ***

Names of linear predictors: logit(P[Y<=1]), logit(P[Y<=2])

Residual deviance: 1986.59 on 1880 degrees of freedom

Log-likelihood: -993.2951 on 1880 degrees of freedom

Exponentiated coefficients:

age:1	age:2	education:1	education:2	income:1	income:2
0.9977883	0.9962231	0.9928875	0.9861065	0.9840517	0.9898212

Other link functions

- We can use different link function in the proportional model for “party”

$$\Pr\{Y \leq 1\} = F(\alpha_1 + \beta x) \text{ and } \Pr\{Y \leq 2\} = F(\alpha_2 + \beta x)$$

- The function F could be the CDF of
 - A logistic distribution (Proportional Odds Model)
 - A normal distribution (Ordered Probit Model)
 - An extreme value distribution (Proportional Hazards Model)

Probit and complementary log-log links

```
nmod.probit<-vglm(party ~ age + education + income,  
family=cumulative(parallel=TRUE, link=probit), rnes96)  
summary(nmod.probit)
```

```
nmod.cloglog<-vglm(party ~ age + education + income,  
family=cumulative(parallel=TRUE, link=cloglog), rnes96)  
summary(nmod.cloglog)
```

```
vglm(formula = party ~ age + education + income, family = cumulative(parallel = TRUE,  
link = probit), data = rnes96)
```

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept):1	0.306698	0.271604	1.129	0.258809
(Intercept):2	0.976247	0.272625	3.581	0.000342 ***
age	-0.001896	0.002305	-0.823	0.410731
education	-0.006770	0.017717	-0.382	0.702398
income	-0.008059	0.001331	-6.055	1.4e-09 ***

Names of linear predictors: probit(P[Y<=1]), probit(P[Y<=2])

Residual deviance: 1994.127 on 1883 degrees of freedom

Log-likelihood: -997.0636 on 1883 degrees of freedom

Exponentiated coefficients:

age	education	income
0.9981054	0.9932533	0.9919731

```
vglm(formula = party ~ age + education + income, family = cumulative(parallel = TRUE,  
link = cloglog), data = rnes96)
```

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept):1	-0.049780	0.296676	-0.168	0.8667
(Intercept):2	0.693885	0.296105	2.343	0.0191 *
age	-0.002080	0.002462	-0.845	0.3980
education	-0.010662	0.019396	-0.550	0.5825
income	-0.008054	0.001511	-5.331	9.77e-08 ***

Names of linear predictors: cloglog(P[Y<=1]), cloglog(P[Y<=2])

Residual deviance: 2002.988 on 1883 degrees of freedom

Log-likelihood: -1001.494 on 1883 degrees of freedom

Exponentiated coefficients:

age	education	income
0.9979217	0.9893947	0.9919786

Modeling multinomial data

Multinomial data without an order

- We have modeled “party” as an ordered multinomial response.
 - Democrat
 - Independent
 - Republican
- What if we treat “party” as a multinomial variable without an order?
- How can we build a model using “education”, “age”, and “income” as explanatory variables?

Review of logistic regression models for a binary response

- Suppose the binary response (**two categories!**)

$$Y \sim \text{Bernoulli}(p)$$

- We will make two additional assumptions on top of this assumption:
 1. Y relies on X only through its mean, e.g $Y \sim \text{Bernoulli}(p(X))$
 2. The logit transformation of the parameter p

$$\text{logit}(p) = \log\left(\frac{p}{1-p}\right) = \beta_1 + \beta_2 X$$

- How can we extend this modeling idea to **multinomial** data?

Modeling multinomial data

- Suppose Y follows a multinomial distribution

$Y = A$ with prob p_1 , $Y = B$ with p_2 , $Y = C$ with p_3

- We will make two additional assumptions on top of this assumption:
 1. Y relies on X only through the three probabilities $p_1(X), p_2(X), p_3(X)$
 2. The log odds

$$\log\left(\frac{p_2}{p_1}\right) = \alpha_2 + \beta_2 X \quad \text{and} \quad \log\left(\frac{p_3}{p_1}\right) = \alpha_3 + \beta_3 X$$

Notes

- A constraint: $p_1 + p_2 + p_3 = 1$
- We use p_1 as the baseline probability. It does not matter which category we use as the baseline.
- The probability

$$p_j = \frac{\exp \eta_j}{1 + \sum_{j=2}^J \exp \eta_j}$$

$$\eta_1 = 0 \text{ and } \eta_j = \alpha_j + \beta_j x$$

```
library("faraway")
data(nes96)

### Process the data
party<-nes96$PID
levels(party)<-c("Democrat", "Democrat", "Independent",
"Independent", "Independent", "Republican", "Republican")
inca<-
c(1.5, 4, 6, 8, 9.5, 10.5, 11.5, 12.5, 13.5, 14.5, 16, 18.5, 21, 23.5, 27.5, 32
.5, 37.5, 42.5, 47.5, 55, 67.5, 82.5, 97.5, 115)
income<-inca[unclass(nes96$income) ]
rnes96<-data.frame(party, income, education=nes96$educ,
age=nes96$age)
edu.years<-c(9, 10.5, 12, 13, 14, 16, 18)
education<-edu.years[unclass(nes96$educ) ]
rnes96$education<-education
```

```
install.packages("nnet")
library("nnet")
mmod<-multinom(party~age+education+income, rnes96)
summary(mmod)
```

How to interpret the coefficients?

```
> summary(mmod)
Call:
multinom(formula = party ~ age + education + income, data = rnes96)
```

Coefficients:

	(Intercept)	age	education	income
Independent	-1.177840	0.000238586	-0.0007666165	0.01614117
Republican	-1.332179	0.004328175	0.0132647908	0.01747883

Std. Errors:

	(Intercept)	age	education	income
Independent	0.6095759	0.005181846	0.03986776	0.003089633
Republican	0.5622177	0.004736439	0.03660254	0.002875076

Residual Deviance: 1984.366

AIC: 2000.366

Model selection and comparison

```
### Model selection based on AIC
```

```
mmodi<-step(mmod)
```

```
### Model comparison based on the significance
```

```
deviance(mmodi)-deviance(mmod)
```

```
mmod$edf-mmodi$edf
```

```
pchisq(deviance(mmodi)-deviance(mmod), mmod$edf-mmodi$edf, lower=F)
```

```

> mmodi<-step(mmod)
Start: AIC=2000.37
party ~ age + education + income
      Df     AIC
- education 6 1996.539
- age        6 1997.325
<none>      8 2000.366
- income     6 2042.822

Step: AIC=1996.54
party ~ age + income
      Df     AIC
- age      4 1993.424
<none>    6 1996.539
- income   4 2048.850

Step: AIC=1993.42
party ~ income
      Df     AIC
<none>    4 1993.424
- income   2 2045.272

```

Model comparison

```
### Difference between deviances  
> deviance(mmodi)-deviance(mmod)  
[1] 1.058172
```

```
### Difference between degrees of freedom  
> mmod$edf-mmodi$edf  
[1] 4
```

```
### Chi^2 test  
> pchisq(deviance(mmodi)-deviance(mmod), mmod$edf-  
mmodi$edf, lower=F)  
[1] 0.9008507
```

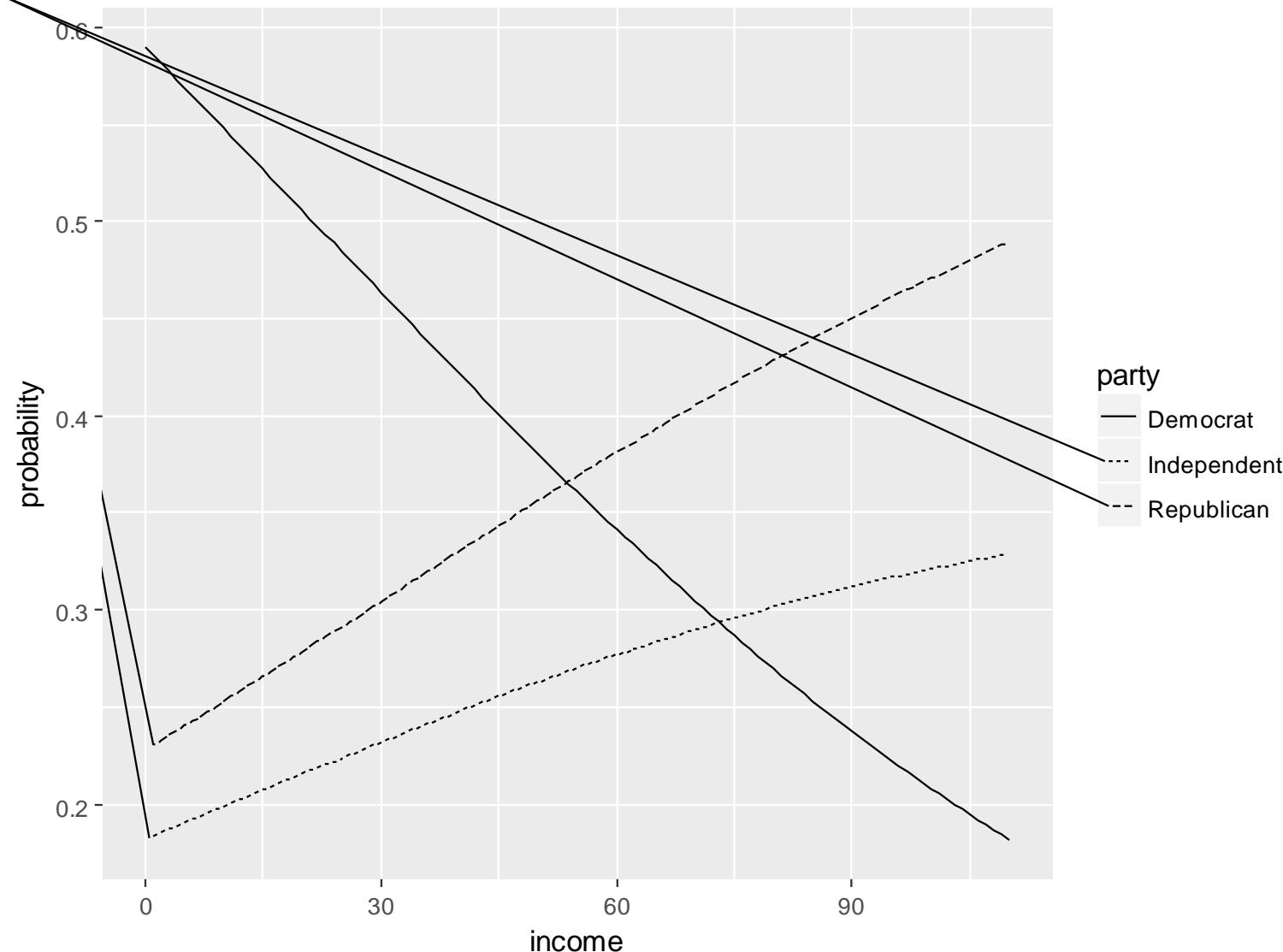
Prediction

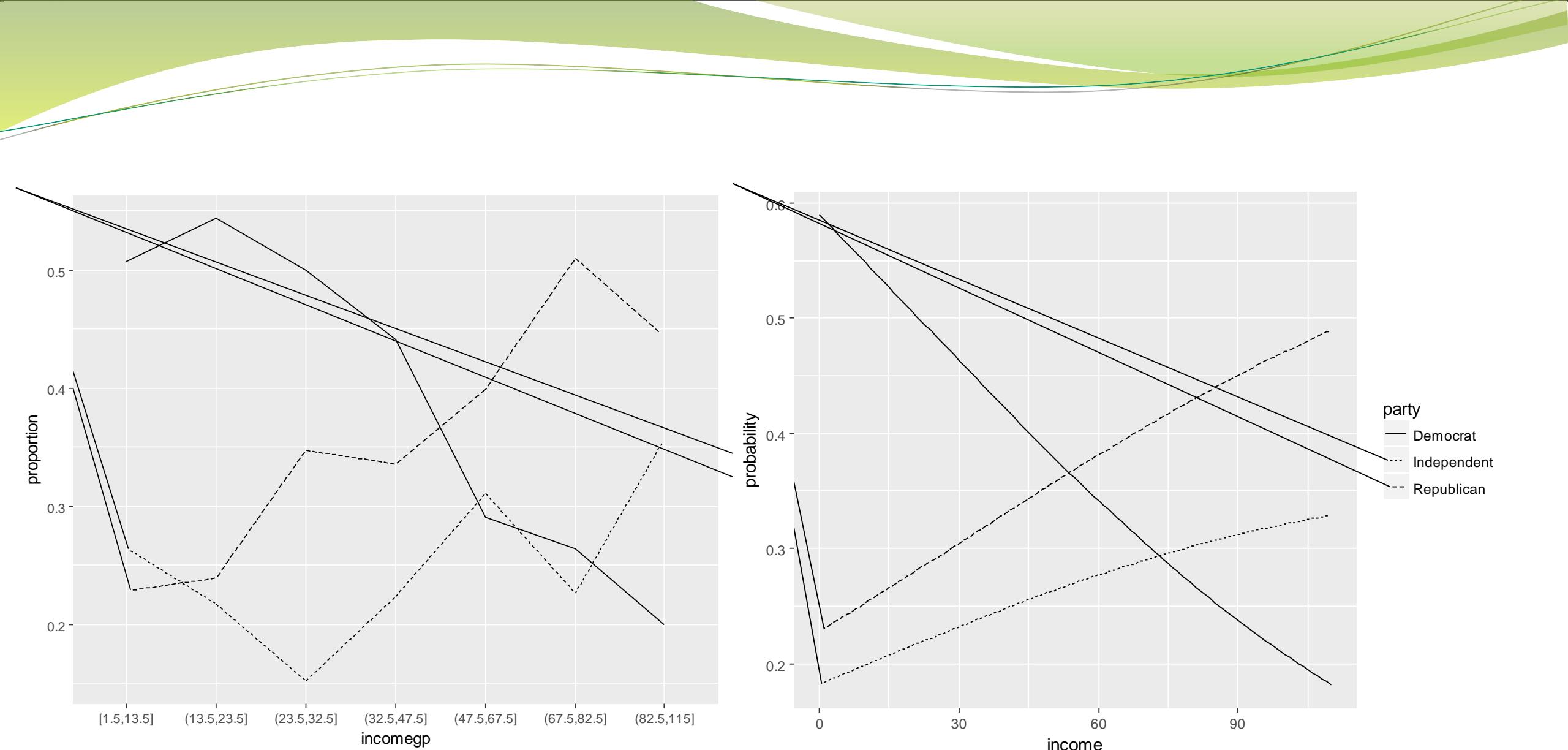
```
### Predict the outcome  
predict(mmodi,data.frame(income=30))  
predict(mmodi,data.frame(income=90))  
  
### Give the probabilities of all possible outcomes  
predict(mmodi,data.frame(income=30),type="probs")  
predict(mmodi,data.frame(income=90),type="probs")
```

```
> predict(mmodi,data.frame(income=30))  
[1] Democrat  
Levels: Democrat Independent Republican  
> predict(mmodi,data.frame(income=90))  
[1] Republican  
Levels: Democrat Independent Republican  
  
> predict(mmodi,data.frame(income=30),type="probs")  
          Democrat Independent Republican  
0.4635683   0.2319741   0.3044576  
> predict(mmodi,data.frame(income=90),type="probs")  
          Democrat Independent Republican  
0.2375755   0.3121136   0.4503109
```

Plot the fitted curve

```
incomes<-0:110 ### Annual income 0-110K  
preds<-  
  data.frame(income=incomes,predict(mmodi,data.frame(i  
ncome=incomes),type="probs"))  
  
install.packages("tidyverse")  
library("tidyverse")  
lpred<-gather(preds,party,probability,-income)  
ggplot(lpred,aes(x=income,y=probability,group=party,li  
netype=party))+geom_line()
```





Prediction table

```
> xtabs (~predict(mmodi) + rnes96$party)
```

		rnes96\$party			
		predict(mmodi)	Democrat	Independent	Republican
Democrat	Democrat	284	123	166	
	Independent	0	0	0	
Republican	Republican	96	116	159	