# Lecture 25 EEP118

## Limited Dependent Variable

1. Logit Model Parameters/ Marginal Effects

   Parameters

   Marginal effects for continuous x

   Marginal effects for discrete x

2. Estimation Maximum Likelihood

3. Tests. Goodness of Fit. Likelihood Ratio Test

   The chi square distribution

Guest speaker: Law school

Study all of chapter 17.1

Posted all remaining DA and solutions, Practice final also

## Limited Depedent Variable Y

The basic context of this set of lectures is when Y is not continuous

Y=0 or 1, Y is binary. YES/NO

Use a Data set on Women labor force participation

Source: MROZ.RAW in Wooldridge. T.A. Mroz (1987), "The Sensitivity of an Empirical Model of Married Women's Hours of Work to Economic and Statistical Assumptions," Econometrica 55, 765-799.

Y= 1 or 0 column called inlf (short for in labor force)

Obs: N=753

inlf byte %9.0g inlf=1 if in labor force, 1975, inlf=0 otherwise

age byte %9.0g woman's age in years

educ byte %9.0g years of schooling

~~~~~~~~ %9.0g # kids < 6 years

kidsge6 byte %9.0g # kids 6-18

nwifeinc float %9.0g (faminc - wage*hours)/1000

hushrs int %9.0g hours worked by husband, 1975

husage byte %9.0g husband's age

huseduc byte %9.0g husband's years of schooling

huswage float %9.0g husband's hourly wage, 1975

city byte %9.0g =1 if live in SMSA

In [2]:
```r
# Load the 'pacman' package
library(pacman)
#packages to use load them now using the pacman "manager"
p_load(dplyr, haven, readr)
#Another great feature of p_load(): if you try to load a package that is not
p_load(ggplot2)


pacman::p_load(lfe, lmtest, haven, sandwich, tidyverse)
# lfe for running fixed effects regression
# lmtest for displaying robust SE in output table
# haven for loading in dta files
# sandwich for producing robust Var-Cov matrix
# tidyverse for manipulating data and producing plots


#The big difference with Stata that appears here is lm() by default
#doesn't compute robust SE - we have to use additional packages/functions
#to compute it. felm does allow for multi-way clustering by default though
#which is nice.

#I added an alternate version of the first plots to show that we can
#change the color of the points according to whether the prediction
#is in [0,1] or outside of it. You can also specify factor(inlf) for
#the latter plots of actual vs. predicted to only have the values 0 or 1 on
#the x-axis.


pacman::p_load(lfe, lmtest, margins, haven, sandwich, tidyverse)
# lfe for running fixed effects regression
# lmtest for displaying robust SE in output table
# haven for loading in dta files
# sandwich for producing robust Var-Cov matrix
# tidyverse for manipulating data and producing plots
#install.packages(sandwich)
#install.packages(lfe)
#install.packages(lmtest)
#install.packages(tidyverse)
library(sandwich)
```

```r
library(lmtest)
library(tidyverse)

#install for margins
#to get marginal effects from Logit
install.packages("mfx")
library(mfx)


# alternate plot theme for ggplot
theme_ed <- theme(
  legend.position = "bottom",
  panel.background = element_rect(fill = NA),
  # panel.border = element_rect(fill = NA, color = "grey75"),
  axis.ticks = element_line(color = "grey95", size = 0.3),
  panel.grid.major = element_line(color = "grey95", size = 0.3),
  panel.grid.minor = element_line(color = "grey95", size = 0.3),
  legend.key = element_blank())
```

In [3]:
```r
#load data
mydata<- read_dta("Lecture24MROZ.DTA")
#Summary stats inlf age educ kidslt6 kidsge6 nwifeinc hushrs husage huseduc
summary(mydata)
```

```
       inlf                hours            kidslt6            kidsge6
 Min.    :0.0000    Min.    :    0.0    Min.    :0.0000    Min.    :0.000
 1st Qu.:0.0000    1st Qu.:    0.0    1st Qu.:0.0000    1st Qu.:0.000
 Median :1.0000    Median :  288.0    Median :0.0000    Median :1.000
 Mean   :0.5684    Mean   :  740.6    Mean   :0.2377    Mean   :1.353
 3rd Qu.:1.0000    3rd Qu.: 1516.0    3rd Qu.:0.0000    3rd Qu.:2.000
 Max.   :1.0000    Max.   : 4950.0    Max.   :3.0000    Max.   :8.000

       age                educ              wage              repwage
 Min.   :30.00    Min.    : 5.00    Min.   : 0.1282    Min.   :0.00
 1st Qu.:36.00    1st Qu.:12.00    1st Qu.: 2.2626    1st Qu.:0.00
 Median :43.00    Median :12.00    Median : 3.4819    Median :0.00
 Mean   :42.54    Mean   :12.29    Mean   : 4.1777    Mean   :1.85
 3rd Qu.:49.00    3rd Qu.:13.00    3rd Qu.: 4.9708    3rd Qu.:3.58
 Max.   :60.00    Max.    :17.00    Max.   :25.0000    Max.   :9.98
                                    NA's    :325
      hushrs             husage            huseduc            huswage
 Min.   : 175    Min.    :30.00    Min.    : 3.00    Min.    : 0.4121
 1st Qu.:1928    1st Qu.:38.00    1st Qu.:11.00    1st Qu.: 4.7883
 Median :2164    Median :46.00    Median :12.00    Median : 6.9758
 Mean   :2267    Mean    :45.12    Mean    :12.49    Mean    : 7.4822
 3rd Qu.:2553    3rd Qu.:52.00    3rd Qu.:15.00    3rd Qu.: 9.1667
 Max.   :5010    Max.    :60.00    Max.    :17.00    Max.    :40.5090

      faminc              mtr             motheduc            fatheduc
 Min.   : 1500    Min.    :0.4415    Min.    : 0.000    Min.    : 0.000
 1st Qu.:15428    1st Qu.:0.6215    1st Qu.: 7.000    1st Qu.: 7.000
 Median :20880    Median :0.6915    Median :10.000    Median : 7.000
 Mean   :23081    Mean    :0.6789    Mean    : 9.251    Mean    : 8.809
 3rd Qu.:28200    3rd Qu.:0.7215    3rd Qu.:12.000    3rd Qu.:12.000
 Max.   :96000    Max.    :0.9415    Max.    :17.000    Max.    :17.000

      unem               city             exper             nwifeinc
 Min.   : 3.000    Min.    :0.0000    Min.    : 0.00    Min.    :-0.02906
 1st Qu.: 7.500    1st Qu.:0.0000    1st Qu.: 4.00    1st Qu.:13.02504
 Median : 7.500    Median :1.0000    Median : 9.00    Median :17.70000
 Mean   : 8.624    Mean    :0.6428    Mean    :10.63    Mean    :20.12896
 3rd Qu.:11.000    3rd Qu.:1.0000    3rd Qu.:15.00    3rd Qu.:24.46600
 Max.   :14.000    Max.    :1.0000    Max.    :45.00    Max.    :96.00000

      lwage             expersq
 Min.   :-2.0542    Min.    :    0
 1st Qu.: 0.8165    1st Qu.:   16
 Median : 1.2476    Median :   81
 Mean   : 1.1902    Mean    :  178
 3rd Qu.: 1.6036    3rd Qu.:  225
 Max.   : 3.2189    Max.    : 2025
 NA's    :325
```

# Fixing Problem 2, make sure predictions are between 0 and 1

Solution Problem 1 –

use a functional for for the probability as a function G( ) of the xs that stays between 0 and 1

e.g., the Logit Model!

the ratio of exponents in the logit below is always between 0 and 1

$$\text{Prob } [Y=1 \mid x] = G(\beta_0 + \beta_1\, x_1 + \beta_2\, x_2 + \cdots + \beta_3\, x_k)$$

Get a G that stays between 0 and 1, and **the Logit is**

**This ratio of exponentials is always between 0 and 1 no matter the betas and xs**

$$\text{Prob } [Y=1 \mid x] = G(\beta_0 + \ldots + \beta_k\, X_k) = \frac{e^{\beta_0 + \beta_1\, x_1 + \beta_2\, x_2 + \cdots + \beta_3\, x_k}}{1 + e^{\beta_0 + \beta_1\, x_1 + \beta_2\, x_2 + \cdots + \beta_3\, x_k}}$$

$$\text{Prob}[Y=1 \mid X] = \Lambda\,(\beta_0 + \beta_1\, x_1 + \beta_2\, x_2 + \cdots + \beta_3\, x_k)$$

24

In [4]:
```
###### Fixing Problem 2 so that predicted Y hats are less than 1 and greater
# In R, use the glm(formula, data, family = binomial(link = "logit")) functi

logit <- glm(inlf ~ nwifeinc + educ + age + kidslt6 + kidsge6, mydata, famil
summary(logit)
```

```
Call:
glm(formula = inlf ~ nwifeinc + educ + age + kidslt6 + kidsge6,
    family = binomial(link = "logit"), data = mydata)

Coefficients:
             Estimate Std. Error z value Pr(>|z|)
(Intercept)  0.722993   0.788698   0.917    0.359
nwifeinc    -0.034891   0.007884  -4.426 9.62e-06 ***
educ         0.257965   0.040744   6.331 2.43e-10 ***
age         -0.057553   0.012737  -4.519 6.23e-06 ***
kidslt6     -1.484437   0.198013  -7.497 6.55e-14 ***
kidsge6     -0.066363   0.067856  -0.978    0.328
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 1029.75  on 752  degrees of freedom
Residual deviance:  908.37  on 747  degrees of freedom
AIC: 920.37

Number of Fisher Scoring iterations: 4
```

Cannot easily interpret parameters here,

next class estimate implied marginal effects given the above estimated Logit parameters

Loading [MathJax]/extensions/Safe.js

Parameters not very meaningful here. (they enter two exponentials to get Phat)

What we want is if say education changes by one, how does the Prob(y=1) change?

# Logit Model Marginal Effects

For a <mark>continuous variable $x_1$ education</mark> for example:

Given that $P(y = 1) = \dfrac{e^{\beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k}}{1 + e^{\beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k}} = \dfrac{e^z}{1 + e^z}$ where

$$z = \beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k$$

Then $\quad \dfrac{\partial P(y=1)}{\partial x_1} = \dfrac{\partial P(y=1)}{\partial z}\dfrac{\partial z}{\partial x_1} = \dfrac{e^z}{(1+e^z)^2}\dfrac{\partial z}{\partial x_1}$

$$<=> \quad \frac{\partial P(y = 1)}{\partial x_1} = \frac{e^z}{(1 + e^z)^2}\beta_1$$
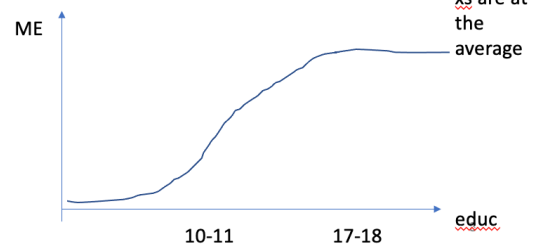
## Logit Model Marginal Effects (ME)

For a continuous variable $x_1$ education for example:

For education, given the estimates above, if education changes by one, then the ME on the Prob(y=1) is given by

$$\frac{e^{\widehat{\beta_0} + \widehat{\beta_1} educ + \cdots + \widehat{\beta_k} x_k}}{\left(1 + e^{\widehat{\beta_0} + \widehat{\beta_1} educ + \cdots + \widehat{\beta_k} x_k}\right)^2}\,\widehat{\beta_1} = \frac{e^{0.72544 + 0.2576\ educ + \cdots - 0.0351\ nwifeinc}}{(1 + e^{0.72544 + 0.2576\ educ + \cdots - 0.0351\ nwifeinc})^2}\ \mathbf{0.2576}$$

Where we substitute the estimated beta hats.

<mark>Note that the ME depends on the starting point of educ and also on all the other x's.</mark>

All other xs are at the average

ME

educ

10-11    17-18

# Logit Model Marginal Effects (ME)

### For a continuous variable $x_1$ education for example:

$$\frac{e^{\widehat{\beta_0}+\widehat{\beta_1}educ+\cdots+\widehat{\beta_k}x_k}}{\left(1+e^{\widehat{\beta_0}+\widehat{\beta_1}educ+\cdots+\widehat{\beta_k}x_k}\right)^2}\ \widehat{\beta_1}=\frac{e^{0.72544+0.2576\ educ+\cdots-0.0351\ nwifeinc}}{\left(1+e^{0.72544+0.2576\ educ+\cdots-0.0351\ nwifeinc}\right)^2}\ 0.2576$$

Where we substitute the estimated beta hats.

==Note that the ME depends on the starting point of educ==

==and also on all the other x's.==



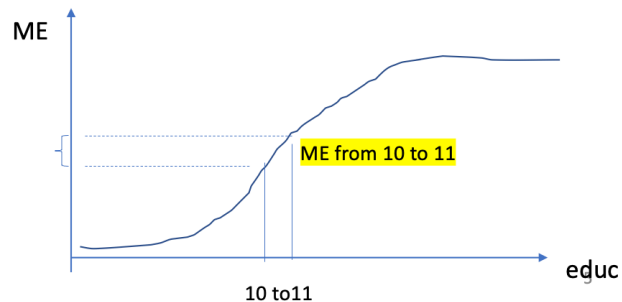For a continuous variable x1 education for example: How does one report the marginal effects (ME) then given that it depends on xs and starting point?

Report it for a fictitious person that would have all x's at the average, that is, for $(\overline{educ})$=12.2, $(\overline{kids})$=0.238 etc etc, all average of all x's, in this case, ME education is 0.0537, or 5.37 percentage points ---see next cell on how to get estimated ME

```
In [5]:  #Marginal Effects (ME) at Mean X
         logitmfx(inlf ~ nwifeinc + educ + age + kidslt6 + kidsge6, lec24df, atmean =
```

```
Error in eval(expr, envir, enclos): object 'lec24df' not found
Traceback:

1. logitmfx(inlf ~ nwifeinc + educ + age + kidslt6 + kidsge6, lec24df,
   .      atmean = TRUE)
2. logitmfxest(formula, data, atmean, robust, clustervar1, clustervar2,
   .      start, control)
3. is.data.frame(data)
```

```
In [9]:  #compuyte the average ME, not at the mean of all X's as above

         #Average Marginal Effects
         logitmfx(inlf ~ nwifeinc + educ + age + kidslt6 + kidsge6, mydata, atmean =
```

```
Call:
logitmfx(formula = inlf ~ nwifeinc + educ + age + kidslt6 + kidsge6,
    data = mydata, atmean = FALSE)

Marginal Effects:
              dF/dx  Std. Err.        z     P>|z|
nwifeinc -0.0072699  0.0017413 -4.1749 2.981e-05 ***
educ      0.0537496  0.0095011  5.6572 1.539e-08 ***
age      -0.0119918  0.0028182 -4.2551 2.089e-05 ***
kidslt6  -0.3092968  0.0480215 -6.4408 1.188e-10 ***
kidsge6  -0.0138274  0.0141772 -0.9753    0.3294
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```
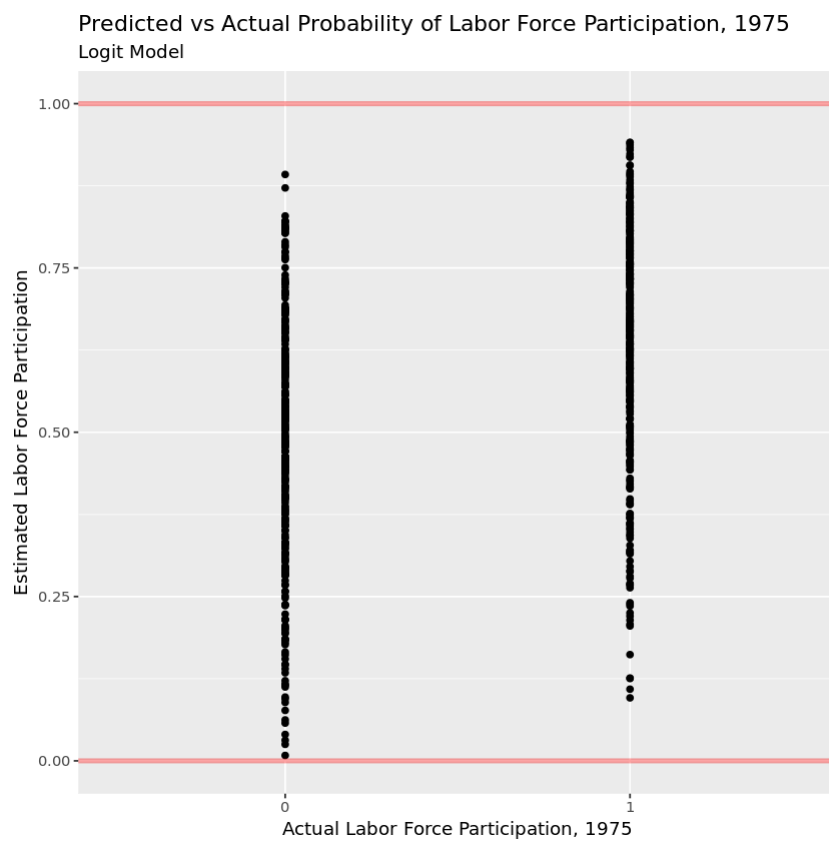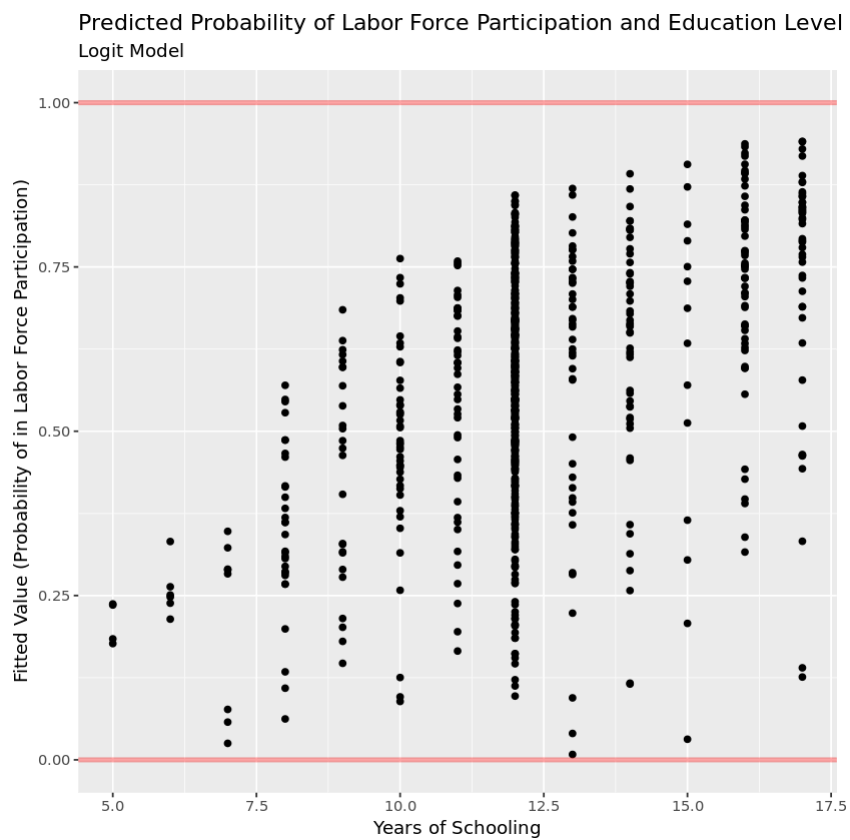
In [7]:
```r
#generate predictions

mydata <- mutate(mydata, log_fit = logit$fitted.values) # add in the logit 1


#Reproduce figures for logit
# no need to use the second approach as we're always within [0,1] with logit
# set data and aesthetics (x and y vars here since the same for all elements
ggplot(mydata, aes(x = educ, y = log_fit)) +
  # First add points, color determined by whether in or out of [0,1]
  geom_point() + # add points
  # add horizontal lines, width slightly wider, making partially transparent
  geom_hline(yintercept=0, size = 1.4, alpha = 0.35, color = "red") + # add
  geom_hline(yintercept=1, size = 1.4, alpha = 0.35, color = "red") + # add
  # generate labels
  labs(title = "Predicted Probability of Labor Force Participation and Educa
       subtitle = "Logit Model",
       x = "Years of Schooling",
       y = "Fitted Value (Probability of in Labor Force Participation)")

# actual vs predicted
ggplot(mydata, aes(x = factor(inlf), y = log_fit)) +
  # First add points, color determined by whether in or out of [0,1]
  geom_point() +
  # add horizontal lines, width slightly wider, making partially transparent
  geom_hline(yintercept=0, size = 1.4, alpha = 0.35, color = "red") + # add
  geom_hline(yintercept=1, size = 1.4, alpha = 0.35, color = "red") + # add
  # generate labels
  labs(title = "Predicted vs Actual Probability of Labor Force Participation
       subtitle = "Logit Model",
       x = "Actual Labor Force Participation, 1975",
       y = "Estimated Labor Force Participation")
```

```
Warning message:
"Using `size` aesthetic for lines was deprecated in ggplot2 3.4.0.
i Please use `linewidth` instead."
```

## Predicted Probability of Labor Force Participation and Education Level
Logit Model



## Predicted vs Actual Probability of Labor Force Participation, 1975
Logit Model



# Marginal Effect for a Discrete variable X

For a discrete variable x1 city for example:

We need to compute the difference in probability, that is ME city= Prob(y=1| x, city=1) – Prob(y=1| x, city=0)

And once again we evaluate all at the average of all other x's

(*) dy/dx is for discrete change of dummy variable from 0 to 1

```
In [10]:  #run a logit with a city dummy variable
          logit2 <- glm(inlf ~ nwifeinc + educ + age + kidslt6 + kidsge6+city, mydata,
          summary(logit2)

          #Average Marginal Effects
          logitmfx(inlf ~ nwifeinc + educ + age + kidslt6 + kidsge6+city, mydata, atme

          #Marginal Effects at Mean X
          logitmfx(inlf ~ nwifeinc + educ + age + kidslt6 + kidsge6+city, mydata, atme
```

```
Call:
glm(formula = inlf ~ nwifeinc + educ + age + kidslt6 + kidsge6 +
    city, family = binomial(link = "logit"), data = mydata)

Coefficients:
            Estimate Std. Error z value Pr(>|z|)
(Intercept)  0.725440   0.789091   0.919    0.358
nwifeinc    -0.035075   0.008067  -4.348 1.37e-05 ***
educ         0.257560   0.040910   6.296 3.06e-10 ***
age         -0.057689   0.012800  -4.507 6.58e-06 ***
kidslt6     -1.484777   0.198075  -7.496 6.58e-14 ***
kidsge6     -0.066625   0.067901  -0.981    0.326
city         0.019103   0.174730   0.109    0.913
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 1029.75  on 752  degrees of freedom
Residual deviance:  908.36  on 746  degrees of freedom
AIC: 922.36

Number of Fisher Scoring iterations: 4
```

```
Call:
logitmfx(formula = inlf ~ nwifeinc + educ + age + kidslt6 + kidsge6 +
    city, data = mydata, atmean = FALSE)

Marginal Effects:
              dF/dx  Std. Err.        z     P>|z|
nwifeinc -0.0073083  0.0017781  -4.1102 3.953e-05 ***
educ      0.0536651  0.0095296   5.6314 1.788e-08 ***
age      -0.0120200  0.0028312  -4.2456 2.181e-05 ***
kidslt6  -0.3093674  0.0480343  -6.4406 1.190e-10 ***
kidsge6  -0.0138819  0.0141867  -0.9785    0.3278
city      0.0039811  0.0364217   0.1093    0.9130
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

dF/dx is for discrete change for the following variables:

[1] "city"
Call:
logitmfx(formula = inlf ~ nwifeinc + educ + age + kidslt6 + kidsge6 +
    city, data = mydata, atmean = TRUE)

Marginal Effects:
              dF/dx  Std. Err.        z     P>|z|
nwifeinc -0.0085755  0.0019749  -4.3423 1.410e-05 ***
educ      0.0629700  0.0099917   6.3023 2.933e-10 ***
age      -0.0141041  0.0031252  -4.5130 6.393e-06 ***
kidslt6  -0.3630078  0.0486205  -7.4662 8.258e-14 ***
kidsge6  -0.0162889  0.0166006  -0.9812    0.3265
city      0.0046722  0.0427535   0.1093    0.9130
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

dF/dx is for discrete change for the following variables:

[1] "city"
```

For a city relative to not a city the probability of a woman being in the labor force increases by 0.004,

but not significantly because the p value of the marginal effect is 0.913

and confidence interval for city Marginal effect covers zero : lower= -0.06737453 upper=0.075335006

# Estimation of Logit - by Maximum Likelihood

Maximum Likelihood

Derivation: for each observation of a woman

Suppose woman i working Yi=1, then, the prob is Pr(Yi=1|xi]= $\Lambda(\beta\_0+\beta\_1\ x1i\ )$

Suppose woman j is not working, Yi=0, then the prob of that is Pr(Yj=0|xj]= 1- $\Lambda(\beta\_0+\beta\_1\ x1j\ )$

Maximum Likelihood

The Probability of observing i working and j not is equal to the product below which is the

Likelihood

= $(\Lambda\ \beta\_0+\beta\_1\ x1i\ ) * [1- \Lambda(\beta\_0+\beta\_1\ x1j\ )\ ]$

= Pr(Yi=1|xi] times Pr(Yj=0|xj]

- Put all the working in data together and all the non working
- The prob to see what we see in the sample is the product of the prob of all the working i's

Likelihood = $\Pi_i\ (\Lambda(\ \beta_0 + \beta_1\ x_{1i}\ )\ )\ \Pi_j\ [1- \Lambda(\ \beta_0+\beta_1\ x_{1j}\ )\ ]$

*all* $Y_i$ = inlf$_i$ =1  if women in labor market

$Y_i$ = inlf$_i$ =0  if  not in labor market

and the product of the prob of all the non working j's.

$$ L = \prod_j \left[\frac{e^{X_i\beta}}{1+e^{X_i\beta}}\right]^{y_i} \prod_i [1 - \frac{e^{X_i\beta}}{1+e^{X_i\beta}}]^{1-y_i} $$

21

- Logging all that

log Likelihood =

LL = $\sum_i ln[\Lambda(\beta_o + X_i\ \beta)]$ + $\sum_j ln\ [1 - \Lambda(\beta_o + X_j\ \beta)]$

*all* $Y_i$ = inlf$_i$ =1
if women in labor market

$Y_i$ = inlf$_j$ =0  if
not in labor market

$$ logL = \sum_j y_i * log[\frac{e^{X_i\beta}}{1+e^{X_i\beta}}] + \sum_i (1 - y_i)log[1 - \frac{e^{X_i\beta}}{1+e^{X_i\beta}}] $$

Loading [MathJax]/extensions/Safe.js

# Estimation Logit, Max Likelihood

- Put all the working in data together and all the non working
- The prob to see what we see in the sample is the product of the prob of all the working i's and the product of the prob of all the non working j's.
- If we log all of that we get
- **log Likelihood =**

$$LL = \sum_i ln[\Lambda(\beta_o + X_i \, \beta)] \quad + \quad \sum_j ln\,[1 - \Lambda(\beta_o + X_j \, \beta)]$$

*for all* $Y_i$ = inlf$_i$ =1   if women
in labor market

*for all* $Y_i$ = inlf$_i$ =0   if  not in labor market

$$logL = \sum_j y_i * log[\frac{e^{X_i\beta}}{1+e^{X_i\beta}}] + \sum_i (1 - y_i)log[1 - \frac{e^{X_i\beta}}{1+e^{X_i\beta}}]$$

In [8]:
```
#estimate a model with lots of X's

logit_u <- glm(inlf ~ nwifeinc + educ + age + kidslt6 + kidsge6+city+hushrs+
summary(logit_u)
```

```
Call:
glm(formula = inlf ~ nwifeinc + educ + age + kidslt6 + kidsge6 +
    city + hushrs + husage + huseduc + huswage, family = binomial(link = "lo
git"),
    data = mydata)

Coefficients:
             Estimate Std. Error z value Pr(>|z|)
(Intercept)  2.1072318  0.9407917   2.240   0.0251 *
nwifeinc    -0.0182788  0.0128726  -1.420   0.1556
educ         0.2893468  0.0478669   6.045 1.50e-09 ***
age         -0.0383568  0.0224972  -1.705   0.0882 .
kidslt6     -1.5370349  0.2009480  -7.649 2.03e-14 ***
kidsge6     -0.0648634  0.0684488  -0.948   0.3433
city         0.0147352  0.1809473   0.081   0.9351
hushrs      -0.0003818  0.0001706  -2.238   0.0252 *
husage      -0.0283468  0.0224390  -1.263   0.2065
huseduc     -0.0354425  0.0365281  -0.970   0.3319
huswage     -0.0434876  0.0372837  -1.166   0.2435
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 1029.75  on 752  degrees of freedom
Residual deviance:  900.47  on 742  degrees of freedom
AIC: 922.47

Number of Fisher Scoring iterations: 4
```

What do you see in the output above?

Loading [MathJax]/extensions/Safe.js

AIC reported, a good measure of fit that is also used for model comparison

Akaike information Criterion (AIC) , not R squared any more, no more minimizing SSR

now we are maximizing log Likelihood as the estimation criterion, what are the parameters that make the sample we see the most likely?

AIC: 922.36

obtained by

Akaike Information Criterion

AIC=ln(ei2/n)+(2k/n)=ln(SSR/n)+(2k/n)

## Hypothesis testing for one coefficient?

In [9]:
```
#Hypothesis testing for one coefficient

#Single parameter test-  use normal  z below


#Coefficients:
#                       Estimate  Std. Error    z value  Pr(>|z|)

#(Intercept)            0.725440   0.789091    0.919    0.358

#nwifeinc              -0.035075   0.008067   -4.348    1.37e-05 ***

#educ                   0.257560   0.040910    6.296    3.06e-10 ***

#age                   -0.057689   0.012800   -4.507    6.58e-06 ***

#kidslt6               -1.484777   0.198075   -7.496    6.58e-14 ***

#kidsge6               -0.066625   0.067901   -0.981    0.326

#city                   0.019103   0.174730    0.109    0.913
```

For example, reject that educaition coefficient is zero. z stat is 6.29 p value 3.06e-10 ***

## Hypothesis Testing for multiple coefficients?

likelihood ratio test in step 2

and critical values of a chi squared distribution in step 3

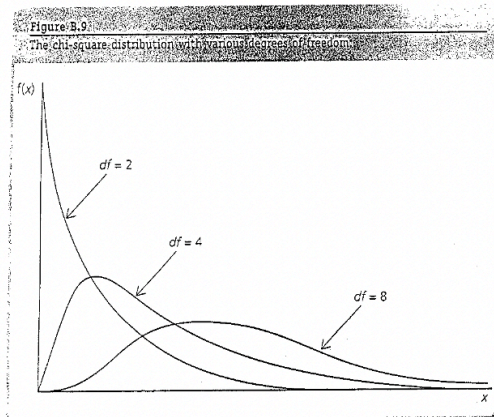# Hypothesis Testing for multiple betas

**LIKELIHOOD RATIO TEST**

Example: $H_o$ $beta_1=beta_2=beta_3=beta_4=... = 0$    (q = number of restrictions)

Under the null hypothesis:

**LR=2 [Log likelihood unrestricted – Log Likelihood restricted]**

is distributed $\chi^2$(q)    **Chi square with q= degrees of freedom**

## The chi -square distribution and table



Figure B.9
The chi-square distribution with various degrees of freedom

$\chi^2$(q)

| | TABLE G.4 Critical Values of the Chi-Square Distribution | | |
|---|---|---|---|
| | Significance Level | | |
| | .10 | .05 | .01 |
| 1 | 2.71 | 3.84 | 6.63 |
| 2 | 4.61 | 5.99 | 9.21 |
| 3 | 6.25 | 7.81 | 11.34 |
| 4 | 7.78 | 9.49 | 13.28 |
| 5 | 9.24 | 11.07 | 15.09 |
| 6 | 10.64 | 12.59 | 16.81 |
| 7 | 12.02 | 14.07 | 18.48 |
| 8 | 13.36 | 15.51 | 20.09 |
| 9 | 14.68 | 16.92 | 21.67 |
| 10 | 15.99 | 18.31 | 23.21 |
| 11 | 17.28 | 19.68 | 24.72 |
| 12 | 18.55 | 21.03 | 26.22 |
| 13 | 19.81 | 22.36 | 27.69 |
| 14 | 21.06 | 23.68 | 29.14 |
| 15 | 22.31 | 25.00 | 30.58 |
| 16 | 23.54 | 26.30 | 32.00 |
| 17 | 24.77 | 27.59 | 33.41 |
| 18 | 25.99 | 28.87 | 34.81 |
| 19 | 27.20 | 30.14 | 36.19 |
| 20 | 28.41 | 31.41 | 37.57 |
| 21 | 29.62 | 32.67 | 38.93 |
| 22 | 30.81 | 33.92 | 40.29 |
| 23 | 32.01 | 35.17 | 41.64 |
| 24 | 33.20 | 36.42 | 42.98 |
| 25 | 34.38 | 37.65 | 44.31 |
| 26 | 35.56 | 38.89 | 45.64 |
| 27 | 36.74 | 40.11 | 46.96 |
| 28 | 37.92 | 41.34 | 48.28 |
| 29 | 39.09 | 42.56 | 49.59 |
| 30 | 40.26 | 43.77 | 50.89 |

Degrees of Freedom

*Example:* The 5% critical value with $df$ = 8 is 15.51.
*Source:* This table was generated using the Stata® function invchi.

5 Steps as usual in hypothesis Testing

# STEPS in Hypothesis testing

- Specify the null and the alternative hypothesis
- Run logit with all xs on the right = unrestricted model
  - Get the Log Likelihood value for the unrestricted $L_{UR}$
- Then run logit omitting 4 x's, we are testing whether those betas for those x's are zero – this is the restricted model
  - Get the Log Likelihood value for the restricted $L_R$

- Compute Likelihood Ratio Test Statistic= **LR**=2 ($L_{UR}$-$L_R$)
- Compare with critical value of $\chi^2$ with 4 degrees of freedom for significance level chosen
- If critical value less than **LR** then we reject the null. Otherwise cannot reject the null

## now coding and computing and doing the actual test

```
In [10]:   #step 1 Null that coefficients on the four husbands charactetistics, all fou

           #step 2


           #likelihood testing

           #run unrestricted model
           logit_u <- glm(inlf ~ nwifeinc + educ + age + kidslt6 + kidsge6+city+hushrs+
           summary(logit_u)

           #get the log likelihood of the unrestricted model
```

```
Call:
glm(formula = inlf ~ nwifeinc + educ + age + kidslt6 + kidsge6 +
    city + hushrs + husage + huseduc + huswage, family = binomial(link = "lo
git"),
    data = mydata)

Coefficients:
             Estimate Std. Error z value Pr(>|z|)
(Intercept)  2.1072318  0.9407917   2.240   0.0251 *
nwifeinc    -0.0182788  0.0128726  -1.420   0.1556
educ         0.2893468  0.0478669   6.045 1.50e-09 ***
age         -0.0383568  0.0224972  -1.705   0.0882 .
kidslt6     -1.5370349  0.2009480  -7.649 2.03e-14 ***
kidsge6     -0.0648634  0.0684488  -0.948   0.3433
city         0.0147352  0.1809473   0.081   0.9351
hushrs      -0.0003818  0.0001706  -2.238   0.0252 *
husage      -0.0283468  0.0224390  -1.263   0.2065
huseduc     -0.0354425  0.0365281  -0.970   0.3319
huswage     -0.0434876  0.0372837  -1.166   0.2435
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 1029.75  on 752  degrees of freedom
Residual deviance:  900.47  on 742  degrees of freedom
AIC: 922.47

Number of Fisher Scoring iterations: 4
```

In [11]:
```r
#run the restricted model
#no husband charct as regressors

logit_r <- glm(inlf ~ nwifeinc + educ + age + kidslt6 + kidsge6+city, mydata
summary(logit_r)


#get the log likelihood of restricted model
```

```
Call:
glm(formula = inlf ~ nwifeinc + educ + age + kidslt6 + kidsge6 +
    city, family = binomial(link = "logit"), data = mydata)

Coefficients:
             Estimate Std. Error z value Pr(>|z|)
(Intercept)  0.725440   0.789091   0.919    0.358
nwifeinc    -0.035075   0.008067  -4.348 1.37e-05 ***
educ         0.257560   0.040910   6.296 3.06e-10 ***
age         -0.057689   0.012800  -4.507 6.58e-06 ***
kidslt6     -1.484777   0.198075  -7.496 6.58e-14 ***
kidsge6     -0.066625   0.067901  -0.981    0.326
city         0.019103   0.174730   0.109    0.913
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 1029.75  on 752  degrees of freedom
Residual deviance:  908.36  on 746  degrees of freedom
AIC: 922.36

Number of Fisher Scoring iterations: 4
```

In [12]: 
```
#get both log likelihood values for the test statistics we will compute to e

#get log likelihood value unrestricted
logLik(logit_u)
```

'log Lik.' -450.2368 (df=11)

In [13]: 
```
#get log likelihood value restricted
logLik(logit_r)
```

'log Lik.' -454.1793 (df=7)

# compute the chi square stat

# By hand, you will do this in Pset 5:

LR = 2 (loglikelihood UR – loglikelihood R) = 2 *(-450.237- + 454.179)=2 *3.94
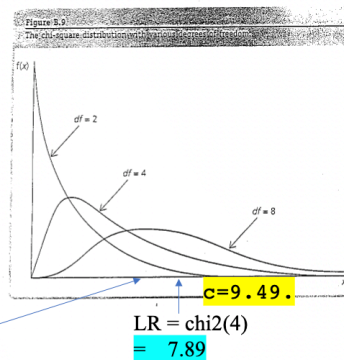
So LR = chi2(4) = 7.89

# step 3 go to the table and get the critical value for a certain significance level

# see below

Step 3: get critical value from Chi squared Table:

at 10% c=7.78

At 5% c=9.49



Figure B.9
The chi-square distribution with various degrees of freedom.

c=9.49.

LR = chi2(4)
= 7.89

**Step 4:** get critical value from Chi squared Table: at 10% c=7.78 At 5% c=9.49. reject at 10% cannot reject at 5%.

**Step 5:** conclude with a sentence. At 5%, there is no statistical evidence that husbands characteristics matter for prob woman in labor force controlling for non wife income, kids, woman educ, etc

TABLE G.4
Critical Values of the Chi-Square Distribution

| Degrees of Freedom | Significance Level | | |
|---|---|---|---|
| | .10 | .05 | .01 |
| 1 | 2.71 | 3.84 | 6.63 |
| 2 | 4.61 | 5.99 | 9.21 |
| 3 | 6.25 | 7.81 | 11.34 |
| 4 | 7.78 | 9.49 | 13.28 |
| 5 | 9.24 | 11.07 | 15.09 |
| 6 | 10.64 | 12.59 | 16.81 |
| 7 | 12.02 | 14.07 | 18.48 |
| 8 | 13.36 | 15.51 | 20.09 |
| 9 | 14.68 | 16.92 | 21.67 |
| 10 | 15.99 | 18.31 | 23.21 |
| 11 | 17.28 | 19.68 | 24.72 |
| 12 | 18.55 | 21.03 | 26.22 |
| 13 | 19.81 | 22.36 | 27.69 |
| 14 | 21.06 | 23.68 | 29.14 |
| 15 | 22.31 | 25.00 | 30.58 |
| 16 | 23.54 | 26.30 | 32.00 |
| 17 | 24.77 | 27.59 | 33.41 |
| 18 | 25.99 | 28.87 | 34.81 |
| 19 | 27.20 | 30.14 | 36.19 |
| 20 | 28.41 | 31.41 | 37.57 |
| 21 | 29.62 | 32.67 | 38.93 |
| 22 | 30.81 | 33.92 | 40.29 |
| 23 | 32.01 | 35.17 | 41.64 |
| 24 | 33.20 | 36.42 | 42.98 |
| 25 | 34.38 | 37.65 | 44.31 |
| 26 | 35.56 | 38.89 | 45.64 |
| 27 | 36.74 | 40.11 | 46.96 |
| 28 | 37.92 | 41.34 | 48.28 |
| 29 | 39.09 | 42.56 | 49.59 |
| 30 | 40.26 | 43.77 | 50.89 |

*Example:* The 5% critical value with $df = 8$ is 15.51.
*Source:* This table was generated using the Stata® function invchi.

34

# step 4

# at 10% c=7.78 < LR=7.89 so we reject the null at 10%

# at 5% c=9.49 > LR = 7.89, so we cannot reject the null at 5%

# Step5: conclude with a sentence. At 5%, there is no statistical evidence that husbands characteristics matter for prob woman in labor force controlling for non wife income, kids, woman educ, etc

## all together

Step 1: H0 Beta_hushrs=Beta_husage=Beta_huseduc=Beta_huswage=0

```
    H1    not H0
```

Step 1: under the null 2 (loglikelihood UR – loglikelihood R) follows a Chi Square with q degrees of freedom

Step 2:

By hand, you will do this in Pset 5:

LR = 2 (loglikelihood UR – loglikelihood R) = 2 *(-450.237- + 454.179)=2 *3.94

So LR = chi2(4) = 7.89

Step 3: get critical value from Chi squared Table: at 10% c=7.78 At 5% c=9.49. reject at 10% cannot reject at 5%.

Step 4/5: conclude with a sentence. At 5%, there is no statistical evidence that husbands charct matter for prob woman in labor force controlling for non wife income, kids, woman educ, etc

In [14]:
```
#in your career you can use a canned command, not in this class though...

##in R: various equivalent specifications of the LR test
lrtest(logit_u, logit_r)
```

A anova: 2 × 5

| | #Df | LogLik | Df | Chisq | Pr(>Chisq) |
|---|---|---|---|---|---|
| | <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| **1** | 11 | -450.2368 | NA | NA | NA |
| **2** | 7 | -454.1793 | -4 | 7.885107 | 0.09587869 |

Loading [MathJax]/extensions/Safe.js

# In R- for your future work in Metrics in life ☺

##in R: various equivalent specifications of the LR test
**lrtest(logit_u, logit_r)**
You get the output in R then:
Likelihood ratio test
Model 1: inlf ~ nwifeinc + educ + age + kidslt6 + kidsge6 + city + hushrs + husage + huseduc + huswage
Model 2: inlf ~ nwifeinc + educ + age + kidslt6 + kidsge6 + city

| #Df | LogLik | Df | Chisq | Pr(>Chisq) |
|-----|--------|-----|-------|-----------|
| 1 11 | -450.24 | | | |
| 2 7 | -454.18 | -3.94 | 7.8851 | 0.09588 . |

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

**Here would not even reject at 10% because p value 0.0958**

the end

In [ ]: