

Econometrics II TA Session #3

Hiroki Kato

1 Empirical Application of Binary Model: Titanic Survivors

Brief Background. “Women and children first” is a behavioral norm, which women and children are saved first in a life-threatening situation. This code was made famous by the sinking of the Titanic in 1912. An empirical application investigates characteristics of survivors of Titanic to answer whether crews obeyed the code or not.

Data. We use an open data about Titanic survivors ¹. Although this dataset contains many variables, we use only four variables: `survived`, `age`, `fare`, and `sex`. We summarize descriptions of variables as follows:

- `survived`: a binary variable taking 1 if a passenger survived.
- `age`: a continuous variable representing passenger's age.
- `fare`: a continuous variable representing how much passenger paid.
- `sex`: a string variable representing passenger's sex.

Using `sex`, we will make a binary variable, called `female`, taking 1 if passenger is female. Instead of `sex`, we use `female` variable in regression.

```
dt <- read.csv(
  file = "../data/titanic.csv",
  header = TRUE, sep = ",", row.names = NULL, stringsAsFactors = FALSE)

dt$female <- ifelse(dt$sex == "female", 1, 0)
dt <- subset(dt, !is.na(survived)&!is.na(age)&!is.na(fare)&!is.na(female))
dt <- dt[,c("survived", "age", "fare", "female")]

set.seed(120511)
train_id <- sample(1:nrow(dt), size = 2*nrow(dt)/3, replace = FALSE)
train_dt <- dt[train_id,]
test_dt <- dt[-train_id,]

head(dt)
```

¹data source: <http://biostat.mc.vanderbilt.edu/DataSets>.

##	survived	age	fare	female
## 1	1	29.00	211.3375	1
## 2	1	0.92	151.5500	0
## 3	0	2.00	151.5500	1
## 4	0	30.00	151.5500	0
## 5	0	25.00	151.5500	1
## 6	1	48.00	26.5500	0

Model. In a binary model, a dependent (outcome) variable y_i takes only two values, i.e., $y_i \in \{0, 1\}$. A binary variable is sometimes called a *dummy* variable. In this application, the outcome variable is **survived**. Explanatory variables are **female**, **age**, and **fare**. The regression function is

$$E[\text{survived} | \text{female}, \text{age}, \text{fare}] \\ = \mathbb{P}[\text{survived} = 1 | \text{female}, \text{age}, \text{fare}] = G(\beta_0 + \beta_1 \text{female} + \beta_2 \text{age} + \beta_3 \text{fare}).$$

The function $G(\cdot)$ is arbitrary function. In practice, we often use following three specifications:

- Linear probability model (LPM): $G(\mathbf{x}_i\beta) = \mathbf{x}_i\beta$.
- Probit model: $G(\mathbf{x}_i\beta) = \Phi(\mathbf{x}_i\beta)$ where $\Phi(\cdot)$ is the standard Gaussian cumulative function.
- Logit model: $G(\mathbf{x}_i\beta) = 1/(1 + \exp(-\mathbf{x}_i\beta))$.

1.1 Linear Probability Model

The linear probability model specifies that $G(a)$ is linear in a , that is,

$$\mathbb{P}[\text{survived} = 1 | \text{female}, \text{age}, \text{fare}] = \beta_0 + \beta_1 \text{female} + \beta_2 \text{age} + \beta_3 \text{fare}.$$

This model can be estimated using the OLS method. In R, we can use the OLS method, running `lm()` function.

```
model <- survived ~ factor(female) + age + fare
LPM <- lm(model, data = train_dt)
```

However, `lm()` function does not deal with heteroskedasticity problem. To resolve it, we need to calculate heteroskedasticity-robust standard errors using the White method.

$$\hat{V}(\hat{\beta}) = \left(\frac{1}{n} \sum_i \mathbf{x}_i' \mathbf{x}_i \right)^{-1} \left(\frac{1}{n} \sum_i \hat{u}_i^2 \mathbf{x}_i' \mathbf{x}_i \right) \left(\frac{1}{n} \sum_i \mathbf{x}_i' \mathbf{x}_i \right)^{-1}$$

```
# heteroskedasticity-robust standard errors
train_dt$"(Intercept)" <- 1
X <- as.matrix(train_dt[,c("(Intercept)", "female", "age", "fare")])
u <- diag(LPM$residuals^2)
```

```

XX <- t(X) %*% X
avgXX <- XX * nrow(X)^{-1}
inv_avgXX <- solve(avgXX)

uXX <- t(X) %*% u %*% X
avguXX <- uXX * nrow(X)^{-1}

vcov_b <- (inv_avgXX %*% avguXX %*% inv_avgXX) * nrow(X)^{-1}
rse_b <- sqrt(diag(vcov_b))

label <- c("(Intercept)", "factor(female)1", "age", "fare")
names(rse_b) <- label

# homoskedasticity-based standard errors
se_b <- sqrt(diag(vcov(LPM)))

print("The Variance of OLS"); vcov(LPM)

## [1] "The Variance of OLS"

##              (Intercept) factor(female)1          age          fare
## (Intercept)    1.505787e-03  -3.905773e-04 -3.676396e-05 -5.951346e-07
## factor(female)1 -3.905773e-04   1.089299e-03  2.569835e-06 -2.154400e-06
## age            -3.676396e-05   2.569835e-06  1.264948e-06 -6.274261e-08
## fare           -5.951346e-07  -2.154400e-06 -6.274261e-08  9.167801e-08

print("The Robust variance of OLS"); vcov_b

## [1] "The Robust variance of OLS"

##              (Intercept)          female          age          fare
## (Intercept)  1.810499e-03 -3.968956e-04 -4.601203e-05  8.979498e-07
## female      -3.968956e-04  1.239665e-03  4.975911e-06 -4.566026e-06
## age         -4.601203e-05  4.975911e-06  1.476806e-06 -7.956793e-08
## fare        8.979498e-07 -4.566026e-06 -7.956793e-08  7.846876e-08

print("The Robust se using White method"); rse_b

## [1] "The Robust se using White method"

##              (Intercept) factor(female)1          age          fare
## 0.0425499596    0.0352088828    0.0012152389    0.0002801228

```

Using the package `lmtest` and `sandwich` is the easiest way to calculate heteroskedasticity-robust standard errors and *t*-statistics.

```

library(lmtest) #use function `coeftest`
library(sandwich) #use function `vcovHC`

```

```
coeftest(LPM, vcov = vcovHC(LPM, type = "HC0"))[, "Std. Error"]
```

```
##      (Intercept) factor(female)1          age          fare
##    0.0425499596    0.0352088828    0.0012152389    0.0002801228
```

Finally, we summarize results of linear probability model in table 1. We will discuss interpretation of results and goodness-of-fit of LPM later.

```
library(stargazer)
stargazer(
  LPM, LPM,
  se = list(se_b, rse_b),
  t.auto = FALSE, p.auto = FALSE,
  report = "vcs", keep.stat = c("n"),
  covariate.labels = c("Female = 1"),
  add.lines = list(
    c("Standard errors", "Homoskedasticity-based", "Heteroskedasticity-robust")),
  title = "Results of Linear Probability Model", label = "LPM",
  type = "latex", header = FALSE, font.size = "small",
  omit.table.layout = "n", table.placement = "h"
)
```

Table 1: Results of Linear Probability Model

	<i>Dependent variable:</i>	
	survived	
	(1)	(2)
Female = 1	0.512 (0.033)	0.512 (0.035)
age	−0.003 (0.001)	−0.003 (0.001)
fare	0.001 (0.0003)	0.001 (0.0003)
Constant	0.245 (0.039)	0.245 (0.043)
Standard errors	Homoskedasticity-based	Heteroskedasticity-robust
Observations	696	696

1.2 Probit and Logit Model

Unlike LPM, the probit and logit model must be estimated using the ML method. The probability of observing y_i is

$$p_\beta(y_i|\mathbf{x}_i) = \mathbb{P}(y_i = 1|x_i)^{y_i} [1 - \mathbb{P}(y_i = 1|x_i)]^{1-y_i} = G(\mathbf{x}_i\beta)^{y_i} (1 - G(\mathbf{x}_i\beta))^{1-y_i}.$$

Taking logalithm yields

$$\log p_\beta(y_i|\mathbf{x}_i) = y_i \log(G(\mathbf{x}_i\beta)) + (1 - y_i) \log(1 - G(\mathbf{x}_i\beta)).$$

The log-likelihood is

$$M_n(\beta) = \sum_{i=1}^n \log p_\beta(y_i|\mathbf{x}_i).$$

The MLE $\hat{\beta}$ holds that the score, which is the first-order derivatives with respect to β , is equal to 0. That is $\nabla_\beta M_n(\hat{\beta}) = 0$. For both logit and probit model, the Hessian matrix, $\nabla_{\beta\beta'}^2 M_n(\beta)$, is always negative definite. This implies that log-likelihood function based on both models is globally concave, and ensures that the MLE maximizes the log-likelihood function. The first-order condition of the probit model is

$$\nabla_\beta M_n(\hat{\beta}) = \sum_{i=1}^n (y_i - \Phi(\mathbf{x}_i\hat{\beta})) \frac{\phi(\mathbf{x}_i\hat{\beta})}{\Phi(\mathbf{x}_i\hat{\beta})(1 - \phi(\mathbf{x}_i\hat{\beta}))} = 0.$$

The first-order condition of the logit model is

$$\nabla_\beta M_n(\hat{\beta}) = \sum_{i=1}^n (y_i - G(\mathbf{x}_i\hat{\beta})) \mathbf{x}_i' = 0.$$

Since it is hard for us to solve this condition analytically, we obtain estimators using numerical procedure.

The asymptotic distribution of $\hat{\beta}$ is $\hat{\beta} \xrightarrow{d} N(\beta, \Sigma_\beta)$ where

$$\Sigma_\beta = - \left(\sum_i E[E[\nabla_{\beta\beta'}^2 \log p_\beta(y_i|\mathbf{x}_i)|\mathbf{x}_i]] \right)^{-1}.$$

In practice, we replace $E[E[\nabla_{\beta\beta'}^2 \log p_\beta(y_i|\mathbf{x}_i)|\mathbf{x}_i]]$ by

$$\frac{1}{n} \sum_i \nabla_{\beta\beta'}^2 \log p_{\hat{\beta}}(y_i|\mathbf{x}_i),$$

that is,

$$\hat{\Sigma}_\beta = \left(\sum_i \nabla_{\beta\beta'}^2 (-\log p_{\hat{\beta}}(y_i|\mathbf{x}_i)) \right)^{-1}.$$

In R, the function `nlm()` provides the Newton-Raphson algorithm to minimize the function². To run this function, we need to define the log-likelihood function (`LnLik`) beforehand. Moreover, since we need to give initial values in augments, we use coefficients estimated by OLS. Alternatively, we often use `glm()` function. Using this function, we do not need to define the log-likelihood function and initial values. Since estimates of `glm()` are approximate to estimates of `nlm()`, we can use this command safely. In this application, we use `nlm` function to minimize the log-likelihood function.

```
Y <- train_dt$survived
female <- train_dt$female
age <- train_dt$age
fare <- train_dt$fare

# log-likelihood
LnLik <- function(b, model = c("probit", "logit")) {

  xb <- b[1] + b[2]*female + b[3]*age + b[4]*fare

  if (model == "probit") {
    L <- pnorm(xb)
  } else {
    L <- 1/(1 + exp(-xb))
  }

  LL_i <- Y * log(L) + (1 - Y) * log(1 - L)
  LL <- -sum(LL_i)

  return(LL)
}

#Newton-Raphson
init <- c(0.169, 0.520, -0.0002, 0.001)
probit <- nlm(LnLik, init, model = "probit", hessian = TRUE)

label <- c("(Intercept)", "factor(female)1", "age", "fare")
names(probit$estimate) <- label
colnames(probit$hessian) <- label; rownames(probit$hessian) <- label

b_probit <- probit$estimate
vcov_probit <- solve(probit$hessian); se_probit <- sqrt(diag(vcov_probit))
LL_probit <- -probit$minimum

#glm function
```

²`optim()` function is another way to minimize the function. Especially, the function `optim(method = "BFGS")` provides the Quasi-Newton algorithm which carries on the spirit of Newton method.

```

model <- survived ~ factor(female) + age + fare
probit_glm <- glm(model, data = train_dt, family = binomial("probit"))

#result
print("The MLE of probit model using nlm"); b_probit

## [1] "The MLE of probit model using nlm"

##      (Intercept) factor(female)1      age      fare
## -0.740010404      1.440663450      -0.009316882      0.006302940

print("The Variance of probit model using nlm"); vcov_probit

## [1] "The Variance of probit model using nlm"

##      (Intercept) factor(female)1      age      fare
## (Intercept)      1.764185e-02      -4.735516e-03 -4.149486e-04 -2.453847e-05
## factor(female)1 -4.735516e-03      1.255295e-02  8.495496e-06 -5.592007e-06
## age              -4.149486e-04      8.495496e-06  1.512962e-05 -9.929199e-07
## fare              -2.453847e-05      -5.592007e-06 -9.929199e-07  1.737151e-06

print("The se of probit model using nlm"); se_probit

## [1] "The se of probit model using nlm"

##      (Intercept) factor(female)1      age      fare
##      0.132822608      0.112039969      0.003889681      0.001318010

print("The coefficients of probit using glm"); coef(probit_glm)

## [1] "The coefficients of probit using glm"

##      (Intercept) factor(female)1      age      fare
## -0.740094134      1.440662013      -0.009314690      0.006303577

print("The se of probit using glm"); sqrt(diag(vcov(probit_glm)))

## [1] "The se of probit using glm"

##      (Intercept) factor(female)1      age      fare
##      0.134738833      0.112061942      0.003966673      0.001326048

```

Using LogLik, we can also estimate logit model by Newton-Raphson algorithm. To compare result, we also use glm() function.

```

#Newton-Raphson
logit <- nlm(LnLik, init, model = "logit", hessian = TRUE)

label <- c("(Intercept)", "factor(female)1", "age", "fare")
names(logit$estimate) <- label
colnames(logit$hessian) <- label; rownames(logit$hessian) <- label

```

```

b_logit <- logit$estimate
vcov_logit <- solve(logit$hessian); se_logit <- sqrt(diag(vcov_logit))
LL_logit <- -logit$minimum

#glm function
logit_glm <- glm(model, data = train_dt, family = binomial("logit"))

#result
print("The MLE of logit model"); b_logit

## [1] "The MLE of logit model"
##      (Intercept) factor(female)1      age      fare
##      -1.19071868      2.36579523     -0.01665811     0.01049121

print("The Variance of logit model"); vcov_logit

## [1] "The Variance of logit model"
##      (Intercept) factor(female)1      age      fare
## (Intercept)      5.347251e-02     -1.306856e-02 -1.260674e-03 -7.166131e-05
## factor(female)1 -1.306856e-02      3.678907e-02 -4.389835e-05 -2.773805e-06
## age             -1.260674e-03     -4.389835e-05  4.703086e-05 -3.343743e-06
## fare            -7.166131e-05     -2.773805e-06 -3.343743e-06  5.199195e-06

print("The se of logit model"); se_logit

## [1] "The se of logit model"
##      (Intercept) factor(female)1      age      fare
##      0.231241234      0.191804780      0.006857905      0.002280174

print("The coefficients of logit using glm"); coef(logit_glm)

## [1] "The coefficients of logit using glm"
##      (Intercept) factor(female)1      age      fare
##      -1.19080405      2.36579304     -0.01665588     0.01049185

print("The se of logit using glm"); sqrt(diag(vcov(logit_glm)))

## [1] "The se of logit using glm"
##      (Intercept) factor(female)1      age      fare
##      0.231133819      0.191810415      0.006862245      0.002272391

```

As a result, table 2 summarizes results of probit model and logit model. t -statistics represents z -value which follows the standard normal distribution. Standard errors are in parentheses. We will discuss interpretation of results and goodness-of-fit later.


```

stargazer(
  probit_glm, logit_glm,
  coef = list(b_probit, b_logit), se = list(se_probit, se_logit),
  t.auto = FALSE, p.auto = FALSE,
  report = "vcs", keep.stat = c("n"),
  covariate.labels = c("Female = 1"),
  add.lines = list(
    c("Log-Likelihood", round(LL_probit, 3), round(LL_logit, 3))),
  title = "Results of Probit and Logit model",
  label = "probit_logit",
  type = "latex", header = FALSE, font.size = "small",
  table.placement = "h", omit.table.layout = "n"
)

```

Table 2: Results of Probit and Logit model

	<i>Dependent variable:</i>	
	survived	
	<i>probit</i>	<i>logistic</i>
	(1)	(2)
Female = 1	1.441 (0.112)	2.366 (0.192)
age	-0.009 (0.004)	-0.017 (0.007)
fare	0.006 (0.001)	0.010 (0.002)
Constant	-0.740 (0.133)	-1.191 (0.231)
Log-Likelihood	-351.507	-351.873
Observations	696	696

1.3 Interpretations

In the linear probability model, interpretations of coefficients are straight-forward. The coefficient β_1 is the change in survival probability given a one-unit increase in continuous variable x . In the case of discrete variable, the coefficient β_1 is the difference in survival probability between two groups. However, when we use the probit or logit model, it is hard for us to interpret results because the partial effect is not constant across other covariates.

As an illustration, the partial effect of continuous variable **age** is

$$\partial_{age} \mathbb{P}[\text{survived} = 1 | \text{female}, \text{age}, \text{fare}] = \begin{cases} \beta_2 & \text{if LPM} \\ \phi(\mathbf{x}_i \beta) \beta_2 & \text{if Probit} \\ \frac{\exp(-\mathbf{x}_i \beta)}{(1 + \exp(-\mathbf{x}_i \beta))^2} \beta_2 & \text{if Logit} \end{cases}.$$

The partial effect of dummy variable **female** is

$$\begin{aligned} & \mathbb{P}[\text{survived} = 1 | \text{female} = 1, \text{age}, \text{fare}] - \mathbb{P}[\text{survived} = 1 | \text{female} = 0, \text{age}, \text{fare}] \\ &= \begin{cases} \beta_1 & \text{if LPM} \\ \Phi(\beta_0 + \beta_1 + \beta_2 \text{age} + \beta_3 \text{fare}) - \Phi(\beta_0 + \beta_1 + \beta_2 \text{age} + \beta_3 \text{fare}) & \text{if Probit} \\ \Lambda(\beta_0 + \beta_1 + \beta_2 \text{age} + \beta_3 \text{fare}) - \Lambda(\beta_0 + \beta_1 + \beta_2 \text{age} + \beta_3 \text{fare}) & \text{if Logit} \end{cases}, \end{aligned}$$

where $\Lambda(a) = 1/(1 + \exp(-a))$.

The first solution is to compute the partial effect at interesting values of \mathbf{x}_i . We often use the sample average of covariates (“average” person) to plugin in the partial effect formula. This is sometimes called *marginal effect at means*. However, since it is unclear what the sample average of dummy variable represents, the marginal effect at means may be hard to explain.

The second solution is to compute the average value of partial effect across the population, that is,

$$\partial_{x_{ij}} \mathbb{P}[y_i = 1 | \mathbf{x}_i] = \beta_j E[g(\mathbf{x}_i \beta)],$$

or, in the case of discrete variable,

$$E[\mathbb{P}[y_i = 1 | x_{ij} = 1, \mathbf{x}_{i,-k}] - \mathbb{P}[y_i = 1 | x_{ij} = 0, \mathbf{x}_{i,-k}]].$$

This is called *average marginal effect* (AME). When we use dummy variables as explanatory variables, we should use this solution.

Standard errors of average marginal effect can be obtained by the Delta method. Let $h_{ij}(\hat{\beta})$ be marginal (partial) effect of the variable x_j for unit i . Then, AME is $h_j(\hat{\beta}) = E[h_{ij}(\hat{\beta})]$. The Delta method implies that $h_j(\hat{\beta}) \xrightarrow{d} N(h_j(\beta), \nabla_{\beta} h_j(\hat{\beta}) V(\beta) (\nabla_{\beta} h_j(\hat{\beta}))')$, where V is variance of β , and

$$\nabla_{\beta} h_j(\hat{\beta}) = \left(\frac{\partial h_j(\hat{\beta})}{\partial \beta_1} \quad \dots \quad \frac{\partial h_j(\hat{\beta})}{\partial \beta_k} \right)$$

When you use the `nlm` function to obtain MLE, we need to calculate standard errors manually. The `DeltaAME` function is a function returning average marginal effect and its standard errors.

When we use the `glm` function, we can use the function `margins` in the library `margins` to obtain the average marginal effect.

Table 3 shows results of linear probability model, probit model, and logit model. In the probit and logit model, coefficients report average marginal effects, t -statistics report z -statistics which follows the standard normal distribution.

All specifications shows that the survival probability of female is about 50% point higher than of male, which is statistically significant. Moreover, the survival probability is decreasing in age, which implies children are more likely to survive. However, the size of coefficient is small. Overall, crews obeyed the code of “women and children first”, but the survival probability of children is not largely different from of adult.

1.4 Model Fitness

There are two measurements of goodness-of-fit. First, the *percent correctly predicted* reports the percentage of unit whose predicted y_i matches the actual y_i . The predicted y_i takes one if $G(\mathbf{x}_i\hat{\beta}) > 0.5$, and takes zero if $G(\mathbf{x}_i\hat{\beta}) \leq 0.5$.

```
# In-sample
in_Y <- train_dt$survived
in_X <- as.matrix(train_dt[,c("(Intercept)", "female", "age", "fare")])

in_Xb_lpm <- in_X %*% matrix(coef(LPM), ncol = 1)
in_Xb_probit <- in_X %*% matrix(b_probit, ncol = 1)
in_Xb_logit <- in_X %*% matrix(b_logit, ncol = 1)

in_hatY_lpm <- ifelse(in_Xb_lpm > 0.5, 1, 0)
in_hatY_probit <- ifelse(pnorm(in_Xb_probit) > 0.5, 1, 0)
in_hatY_logit <- ifelse(1/(1 + exp(-in_Xb_logit)) > 0.5, 1, 0)

in_pcp_lpm <- round(sum(in_Y == in_hatY_lpm)/nrow(in_X), 4)
in_pcp_probit <- round(sum(in_Y == in_hatY_probit)/nrow(in_X), 4)
in_pcp_logit <- round(sum(in_Y == in_hatY_logit)/nrow(in_X), 4)

# Out-of-sample
out_Y <- test_dt$survived
test_dt$"(Intercept)" <- 1
out_X <- as.matrix(test_dt[,c("(Intercept)", "female", "age", "fare")])

out_Xb_lpm <- out_X %*% matrix(coef(LPM), ncol = 1)
out_Xb_probit <- out_X %*% matrix(b_probit, ncol = 1)
out_Xb_logit <- out_X %*% matrix(b_logit, ncol = 1)

out_hatY_lpm <- ifelse(out_Xb_lpm > 0.5, 1, 0)
out_hatY_probit <- ifelse(pnorm(out_Xb_probit) > 0.5, 1, 0)
out_hatY_logit <- ifelse(1/(1 + exp(-out_Xb_logit)) > 0.5, 1, 0)

out_pcp_lpm <- round(sum(out_Y == out_hatY_lpm)/nrow(out_X), 4)
out_pcp_probit <- round(sum(out_Y == out_hatY_probit)/nrow(out_X), 4)
out_pcp_logit <- round(sum(out_Y == out_hatY_logit)/nrow(out_X), 4)
```

Second measurement is the *pseudo R-squared*. The pseudo R-squared is obtained by $1 - \sum_i \hat{u}_i^2 / \sum_i y_i^2$, where $\hat{u}_i = y_i - G(\mathbf{x}_i \hat{\beta})$.

```
Y2 <- in_Y^2

hatu_lpm <- (in_Y - in_Xb_lpm)^2
hatu_probit <- (in_Y - pnorm(in_Xb_probit))^2
hatu_logit <- (in_Y - 1/(1 + exp(-in_Xb_logit)))^2

pr2_lpm <- round(1 - sum(hatu_lpm)/sum(Y2), 4)
pr2_probit <- round(1 - sum(hatu_probit)/sum(Y2), 4)
pr2_logit <- round(1 - sum(hatu_logit)/sum(Y2), 4)
```

Table 3 summarizes two measurements of model fitness. There is little difference among LPM, probit model, and logit model.

```
stargazer(
  LPM, probit_glm, logit_glm,
  coef = list(coef(LPM), b_probit, b_logit),
  se = list(rse_b, se_probit, se_logit),
  t.auto = FALSE, p.auto = FALSE,
  omit = c("Constant"), covariate.labels = c("Female = 1"),
  report = "vcs", keep.stat = c("n"),
  add.lines = list(
    c("Percent correctly predicted (in-sample)",
      in_pcp_lpm, in_pcp_probit, in_pcp_logit),
    c("Percent correctly predicted (out-of-sample)",
      out_pcp_lpm, out_pcp_probit, out_pcp_logit),
    c("Pseudo R-squared", pr2_lpm, pr2_probit, pr2_logit)
  ),
  omit.table.layout = "n", table.placement = "t",
  title = "Titanic Survivors: LPM, Probit (AME), and Logit (AME)",
  label = "titanic",
  type = "latex", header = FALSE
)
```

2 Empirical Application of Ordered Probit and Logit Model: Housing as Status Goods

Breif Background. Social image may affect consumption behavior. Specifically, a desire to signal high income or wealth may cause consumers to purchase status goods. In this application, we explore whether living in an upper floor serves as a status goods.

Data. We use the housing data originally coming from the American Housing Survey con-

Table 3: Titanic Survivors: LPM, Probit (AME), and Logit (AME)

	<i>Dependent variable:</i>		
	survived		
	<i>OLS</i>	<i>probit</i>	<i>logistic</i>
	(1)	(2)	(3)
Female = 1	0.512 (0.035)	1.441 (0.112)	2.366 (0.192)
age	-0.003 (0.001)	-0.009 (0.004)	-0.017 (0.007)
fare	0.001 (0.0003)	0.006 (0.001)	0.010 (0.002)
Percent correctly predicted (in-sample)	0.7802	0.7744	0.7744
Percent correctly predicted (out-of-sample)	0.7794	0.7765	0.7765
Pseudo R-squared	0.5869	0.5873	0.5869
Observations	696	696	696

ducted in 2013 ³. We use the following variable

- **Level**: ordered value of a story of respondent's living (1:Low - 4:High)
- **Levelnum**: variable we recode the response **Level** as 25, 50, 75, 100. This represents the extent of floor height.
- **lnPrice**: logged price of housing (proxy for quality of house)
- **Top25**: a dummy variable taking one if household income is in the top 25 percentile in sample.

```
house <- read.csv(file = "../data/housing.csv", header = TRUE, sep = ",")
house <- house[,c("Level", "lnPrice", "Top25")]
house$Levelnum <- ifelse(
  house$Level == 1, 25,
  ifelse(house$Level == 2, 50,
  ifelse(house$Level == 3, 75, 100)))
head(house)
```

```
##   Level  lnPrice Top25 Levelnum
## 1     3 11.51294     0         75
## 2     4 11.51294     1        100
```

³<https://www.census.gov/programs-surveys/ahs.html>. This is a repeated cross-section survey. We use the data at one time.

## 3	3	11.60824	0	75
## 4	3	11.69526	0	75
## 5	3	12.57764	0	75
## 6	3	12.64433	0	75

Model. The outcome variable is `Level` taking $\{1, 2, 3, 4\}$. Consider the following regression equation of a latent variable:

$$y_i^* = \mathbf{x}_i\beta + u_i,$$

where $\mathbf{x}_i = (\ln Price, Top25)$ and u_i is an error term. The relationship between the latent variable y_i^* and the observed outcome variable is

$$Level = \begin{cases} 1 & \text{if } -\infty < y_i^* \leq a_1 \\ 2 & \text{if } a_1 < y_i^* \leq a_2 \\ 3 & \text{if } a_2 < y_i^* \leq a_3 \\ 4 & \text{if } a_3 < y_i^* < +\infty \end{cases}.$$

Consider the probability of realization of y_i , that is,

$$\begin{aligned} \mathbb{P}(y_i = k | \mathbf{x}_i) &= \mathbb{P}(a_{k-1} - \mathbf{x}_i\beta < u_i \leq a_k - \mathbf{x}_i\beta | \mathbf{x}_i) \\ &= G(a_k - \mathbf{x}_i\beta) - G(a_{k-1} - \mathbf{x}_i\beta), \end{aligned}$$

where $a_4 = +\infty$ and $a_0 = -\infty$. Then, the likelihood function is defined by

$$p((y_i | \mathbf{x}_i), i = 1, \dots, n; \beta, a_1, \dots, a_3) = \prod_{i=1}^n \prod_{k=1}^4 (G(a_k - \mathbf{x}_i\beta) - G(a_{k-1} - \mathbf{x}_i\beta))^{I_{ik}}.$$

where I_{ik} is a indicator variable taking 1 if $y_i = k$. Finally, the log-likelihood function is

$$M(\beta, a_1, a_2, a_3) = \sum_{i=1}^n \sum_{k=1}^4 I_{ik} \log(G(a_k - \mathbf{x}_i\beta) - G(a_{k-1} - \mathbf{x}_i\beta)).$$

Usually, $G(a)$ assumes the standard normal distribution, $\Phi(a)$, or the logistic distribution, $1/(1 + \exp(-a))$.

In R, the library (package) `MASS` provides the `polr` function which estimates the ordered probit and logit model. Although we can use the `nlm` function when we define the log-likelihood function, we do not report this method. To compare results, we use the variable `Levelnum` as outcome variable, and apply the linear regression model.

```
library(MASS)
library(tidyverse) #use case_when()

ols <- lm(Levelnum ~ lnPrice + Top25, data = house)

model <- factor(Level) ~ lnPrice + Top25
```

```

oprobit <- polr(model, data = house, method = "probit")
ologit <- polr(model, data = house, method = "logistic")

a_oprobit <- round(oprobit$zeta, 3)
a_ologit <- round(ologit$zeta, 3)

xb_oprobit <- oprobit$lp
xb_ologit <- ologit$lp

hatY_oprobit <- case_when(
  xb_oprobit <= oprobit$zeta[1] ~ 1,
  xb_oprobit <= oprobit$zeta[2] ~ 2,
  xb_oprobit <= oprobit$zeta[3] ~ 3,
  TRUE ~ 4
)
hatY_ologit <- case_when(
  xb_ologit <= ologit$zeta[1] ~ 1,
  xb_ologit <= ologit$zeta[2] ~ 2,
  xb_ologit <= ologit$zeta[3] ~ 3,
  TRUE ~ 4
)

pred_oprobit <- round(sum(house$Level == hatY_oprobit)/nrow(house), 3)
pred_ologit <- round(sum(house$Level == hatY_ologit)/nrow(house), 3)

```

2.1 Interepretations

Table 4 shows results. OLS model shows that respondents whose household income is in the top 25 percentile live in 3.7% higher floor than other respondents. This implies that high earners want to live in higher floor, which may serve as a status goods. The ordered probit and logit model are in line with this result. To evaluate two models quantitatively, consider the following equation.

$$E[Levelnum|\mathbf{x}_i] = 25P[level = 1|\mathbf{x}_i] + 50P[level = 2|\mathbf{x}_i] + 75P[level = 3|\mathbf{x}_i] + 100P[level = 4|\mathbf{x}_i].$$

We compute this equation with $Top25 = 1$ and $Top25 = 0$ at mean value of $lnPrice$ and take difference.

```

quantef <- function(model) {
  b <- coef(model)
  val1 <- mean(house$lnPrice)*b[1] + b[2]
  val0 <- mean(house$lnPrice)*b[1]

  prob <- matrix(c(rep(val1, 3), rep(val0, 3)), ncol = 2, nrow = 3)
  for (i in 1:3) {

```

```

    for (j in 1:2) {
      prob[i,j] <- pnorm(model$zeta[i] - prob[i,j])
    }
  }
  Ey1 <- 25*prob[1,1] + 50*(prob[2,1]-prob[1,1]) +
    75*(prob[3,1]-prob[2,1]) + 100*(1-prob[3,1])
  Ey0 <- 25*prob[1,2] + 50*(prob[2,2]-prob[1,2]) +
    75*(prob[3,2]-prob[2,2]) + 100*(1-prob[3,2])

  return(Ey1 - Ey0)
}

ef_oprobit <- round(quantef(oprobit), 3)
ef_ologit <- round(quantef(ologit), 3)

```

As a result, we obtain similar values to OLSE. In the ordered probit model, earners in the top 25 percentile live in 4.2% higher floor than others. In the ordered logit model, earners in the top 25 percentile live in 5.9% higher floor than others. Note that, in this application, model fitness seems to be bad because the percent correctly predicted is low (16.7%).

```

stargazer(
  ols, oprobit, ologit,
  report = "vcstp", keep.stat = c("n"),
  omit = c("Constant"),
  add.lines = list(
    c("Cutoff value at 1|2", "", a_oprobit[1], a_ologit[1]),
    c("Cutoff value at 2|3", "", a_oprobit[2], a_ologit[2]),
    c("Cutoff value at 3|4", "", a_oprobit[3], a_ologit[3]),
    c("Quantitative Effect of Top25", "", ef_oprobit, ef_ologit),
    c("Percent correctly predicted", "", pred_oprobit, pred_ologit)
  ),
  omit.table.layout = "n", table.placement = "t",
  title = "Floor Level of House: Ordered Probit and Logit Model",
  label = "housing",
  type = "latex", header = FALSE
)

```

3 Empirical Application of Multinomial Model: Gender Discrimination in Job Position

Brief Background. Recently, many developed countries move toward women's social advancement, for example, an increase of number of board member. In this application, we explore whether the U.S. bank hindered the entrance of female into the workforce.

Table 4: Floor Level of House: Ordered Probit and Logit Model

	<i>Dependent variable:</i>		
	Levelnum	Level	
	<i>OLS</i>	<i>ordered probit</i>	<i>ordered logistic</i>
	(1)	(2)	(3)
lnPrice	0.348 (0.430) t = 0.810 p = 0.418	0.012 (0.016) t = 0.777 p = 0.438	0.019 (0.026) t = 0.745 p = 0.457
Top25	3.714 (1.723) t = 2.156 p = 0.032	0.156 (0.064) t = 2.426 p = 0.016	0.239 (0.106) t = 2.259 p = 0.024
Cutoff value at 1 2		-0.149	-0.25
Cutoff value at 2 3		0.246	0.384
Cutoff value at 3 4		0.97	1.574
Quantitative Effect of Top25		4.17	5.488
Percent correctly predicted		0.167	0.167
Observations	1,612	1,612	1,612

Data. We use a built-in dataset called `BankWages` in the library `AER`. This dataset contains choice of three job position: `custodial`, `admin` and `manage`. The rank of position is `custodial < admin < manage`. Other variables are `education`, `gender`, and `minority`. We use former two variables as explanatory variables.

```
library(AER)
data(BankWages)
dt <- BankWages
dt$job <- as.character(dt$job)
dt$job <- factor(dt$job, levels = c("admin", "custodial", "manage"))
head(BankWages, 5)
```

```
##      job education gender minority
## 1 manage         15   male       no
## 2  admin         16   male       no
## 3  admin         12 female       no
## 4  admin          8 female       no
## 5  admin         15   male       no
```

Model. The outcome variable y_i takes three values $\{0, 1, 2\}$. Then, the multinomial logit model has the following response probabilities

$$P_{ij} = \mathbb{P}(y_i = j | \mathbf{x}_i) = \begin{cases} \frac{\exp(\mathbf{x}_i \beta_j)}{1 + \sum_{k=1}^2 \exp(\mathbf{x}_i \beta_k)} & \text{if } j = 1, 2 \\ \frac{1}{1 + \sum_{k=1}^2 \exp(\mathbf{x}_i \beta_k)} & \text{if } j = 0 \end{cases}.$$

The log-likelihood function is

$$M_n(\beta_1, \beta_2) = \sum_{i=1}^n \sum_{j=0}^2 d_{ij} \log(P_{ij}),$$

where d_{ij} is a dummy variable taking 1 if $y_i = j$.

In R, some packages provide the multinomial logit model. In this application, we use the `multinom` function in the library `nnet`.

```
library(nnet)
est_mlogit <- multinom(job ~ education + gender, data = dt)

# observations and percent correctly predicted
pred <- est_mlogit$fitted.value
pred <- colnames(pred)[apply(pred, 1, which.max)]
n <- length(pred)
pcp <- round(sum(pred == dt$job)/n, 3)

# Log-likelihood and pseudo R-sq
loglik1 <- as.numeric(nnet::logLik.multinom(est_mlogit))
est_mlogit0 <- multinom(job ~ 1, data = dt)
loglik0 <- as.numeric(nnet::logLik.multinom(est_mlogit0))
pr2 <- round(1 - loglik1/loglik0, 3)
```

3.1 Interpretations

Table 5 summarizes the result of multinomial logit model. The coefficient represents the change of $\log(P_{ij}/P_{i0})$ in corresponding covariate. For example, education decreases the log-odds between `custodial` and `admin`, $\log(P_{i,custodial}/P_{i,admin})$ by -0.562. This implies that those who received higher education are more likely to obtain the position `admin`. Highly-educated workers are also more likely to obtain the position `manage`. Moreover, a female dummy decrease the log-odds between `manage` and `admin` by -0.748, which implies that females are less likely to obtain higher position `manage`. From this result, we conclude that the U.S. bank discouraged females to assign higher job position.

Finally, we should check the model fitness. The predicted position is the outcome with the highest estimated probability. The multinomial logit model correctly predicts many cases (correction rate: 85.2%).

Table 5: Multinomial Logit Model of Job Position

	<i>Dependent variable:</i>	
	custodial	manage
	(1)	(2)
Education	−0.562 (0.098) t = −5.721 p = 0.000	1.661 (0.247) t = 6.715 p = 0.000
Female = 1	−10.976 (27.808) t = −0.395 p = 0.694	−0.748 (0.429) t = −1.743 p = 0.082
Constant	5.030 (1.130) t = 4.450 p = 0.00001	−26.730 (3.874) t = −6.899 p = 0.000
Observations	474	
Percent correctly predicted	0.852	
Log-likelihood	−144.928	
Pseudo R-sq	0.546	

```

stargazer(
  est_mlogit,
  covariate.labels = c("Education", "Female = 1"),
  report = "vcstp", omit.stat = c("aic"),
  add.lines = list(
    c("Observations", n, ""),
    c("Percent correctly predicted", pcpr, ""),
    c("Log-likelihood", round(loglik1, 3), ""),
    c("Pseudo R-sq", pr2, "")
  ),
  omit.table.layout = "n", table.placement = "t",
  title = "Multinomial Logit Model of Job Position",
  label = "job",
  type = "latex", header = FALSE
)

```