

# Discrete Choice Models

## Binary Models

Mohammad Arshad Rahman  
(webpage: <https://www.arshadrahman.com>)

Course: Econometrics  
Chennai Mathematical Institute





## Logistic Distribution

Suppose  $X \sim \mathcal{L}(a, b)$ , then the pdf of  $X$  is given by,

$$f_X(x|a, b) = \frac{1}{b} \frac{\exp \left\{ -\frac{(x-a)}{b} \right\}}{\left[ 1 + \exp \left\{ -\frac{(x-a)}{b} \right\} \right]^2}, \quad -\infty < x, a, < \infty; b > 0,$$

and the *cdf* is given by,

$$F_X(x|a, b) = \frac{1}{\left[ 1 + \exp \left\{ -\frac{(x-a)}{b} \right\} \right]}.$$



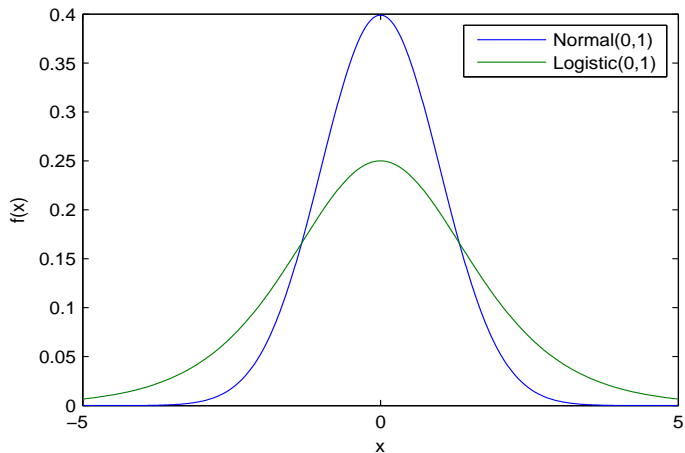


Figure 1: Pdf's for the Normal and Logistic distributions.

## Binary Data Models

Binary models are designed to deal with situations where the outcome (dependent) variable is dichotomous i.e., takes only two values, typically coded as 1 for 'success' and 0 for 'failure'.

For example: when deciding to raise funds for firm expansion, we may model the outcome as 'success' if a firm chooses to raise funds by issuing new stocks and a 'failure' if it raises fund by issuing bonds.

Other examples: Choice of one of the parties in a two-party election, affiliation to union membership, and predicting bank failure.













































## Mortgage Data

The Boston HMDA data set was collected by researchers at the Federal Reserve Bank of Boston. The data set combines information from mortgage applications and a follow-up survey of the banks and other lending institutions that received these mortgage applications.

The data pertain to mortgage applications made in 1990 in the greater Boston metropolitan area. The full data set has 2925 observations, consisting of all mortgage applications by blacks and Hispanics plus a random sample of mortgage applications by whites.



## Data Summary

**Table 1:** Variable definitions and data summary.

VARIABLES	DESCRIPTION	MEAN/COUNT	STD/PROP
deny	mortgage application, = 1 if denied, 0 otherwise	285	0.12
pirat	payments-to-income ratio	0.33	0.12
hirat	inhouse expense-to-total-income ratio	0.25	0.01
lvrat	loan-to-value ratio	0.74	0.18
chist	consumer credit score, values = 1 to 6		
mhist	mortgage credit score, values = 1, 2		
phist	public bad credit record, = 1 if yes, 0 otherwise	175	0.07
unemp	unemployment rate	3.77	2.03
selfemp	self employed, = 1 if yes, 0 otherwise	277	0.12
insurance	denied mortgage insurance, = 1 if yes, 0 otherwise	48	0.02
condomin	condominium/single family residence, = 1 if yes, 0 otherwise	686	0.29
afam (black)	family is black, = 1 if yes, 0 otherwise	339	0.14
single	marital status is single, = 1 if yes, 0 otherwise	936	0.39
hschool	high school and above, = 1 if yes, 0 otherwise	2341	0.98









## LPM: Model 1

Table 2: Estimates from linear probability model.

	Coef	Std Err	t-stat
Intercept	-0.080	0.021	-3.777
pirat	0.603	0.061	9.920

Residual standard error  $\hat{\sigma}$ : 0.3183 on 2378 degrees of freedom. R-squared: 0.03974, Adj. R-squared: 0.03933. F-statistic: 98.41 on 1 and 2378 DF, p-value:  $< 2.2e-16$ .

As per the estimated model, there is a positive relation between (P/I) ratio and the probability of a denied mortgage application. A 1 percentage point increase in (P/I) ratio leads to an increase in the probability of a loan denial by  $0.604 \times 0.01 = 0.00604 \simeq 0.6\%$ .



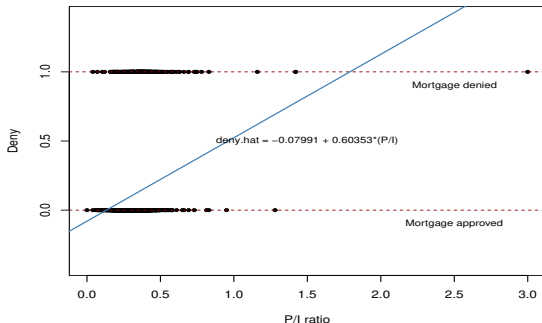


Figure 3: Scatter plot and estimated regression line.

According to the estimated model, a payment-to-income ratio of 1 is associated with an expected probability of mortgage application denial of roughly 50%.



## LPM: Model 2

```
# rename the variable 'afam' for consistency
colnames(HMDA)[colnames(HMDA) == "afam"] <- "black"
```

```
# estimate the model
denymod2 <- lm(deny ~ pirat + black, data = HMDA)
coeftest(denymod2, vcov. = vcovHC)
```

t test of coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-0.090514	0.033430	-2.7076	0.006826**
pirat	0.559195	0.103671	5.3939	7.575e-08***
blackyes	0.177428	0.025055	7.0815	1.871e-12

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









# Probit Model 1

Table 4: Estimates from Probit regression model.

	Coef	Std Err	z-stat
Intercept	-2.194	0.138	-15.927
pirat	2.968	0.386	7.694

Null deviance: 1744.2 on 2379 degrees of freedom. Residual deviance: 1663.6 on 2378 degrees of freedom. AIC: 1667.6

The coefficient for P/I ratio is significant, which implies that applicants with a high P/I ratio face a higher risk of being rejected.

## Probit Model 1: Scatter Plot and Regression Line

```
# plot dataplot(x = HMDA$pirat, y = HMDA$deny,
main = "Probit Model of the Probability of Denial, given P/I Ratio",
xlab = "P/I ratio", ylab = "Deny", pch = 20,
ylim = c(-0.4, 1.4), cex.main = 0.85)

# add horizontal dashed lines and text
abline(h = 1, lty = 2, col = "darkred")
abline(h = 0, lty = 2, col = "darkred")
text(2.5, 0.9, cex = 0.8, "Mortgage denied")
text(2.5, -0.1, cex = 0.8, "Mortgage approved")
# add estimated regression line
x <- seq(0, 3, 0.01)
y <- predict(denyprobbit, list(pirat = x), type = "response")
lines(x, y, lwd = 1.5, col = "steelblue")
```

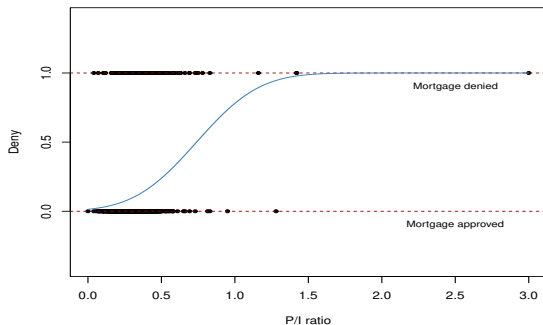


Figure 4: Probit Model of the Probability of Denial, given P/I Ratio

The estimated regression function has a stretched “S”-shape which is typical for the CDF of a continuous random variable with symmetric PDF like that of a normal random variable.

The function is nonlinear and flattens out for large and small values of  $P/I$  ratio. The functional form thus also ensures that the predicted conditional probabilities of a denial lie between 0 and 1.

## Covariate Effect

Probit (and Logit) are non-linear models. So, the coefficients do not give the covariate effect and we need to compute it.

Let's suppose P/I ratio is increased from 0.3 to 0.4. What is the predicted change in the denial probability? This change in probability is computed as follows:

$$\Phi(-2.194 + 2.968 * 0.4) - \Phi(-2.194 + 2.968 * 0.3) = 0.061.$$

We find that an increase in the payment-to-income ratio from 0.3 to 0.4 is predicted to increase the probability of denial by approximately 6.08%.

# Probit Model 1: Covariate Effect

```
# 1. compute predictions for P/I ratio = 0.3, 0.4
predictions1 <- predict(denyprobbit,
newdata = data.frame("pirat" = c(0.3, 0.4)), type = "response")

# 2. Compute difference in probabilities
diff(predictions1)

# Output
0.06081433
```

## Probit Model 2: Is there racial bias?

We continue by using an augmented Probit model to estimate the effect of race on the probability of a mortgage application denial.

```
denyprobit2 <- glm(deny ~ pirat + black,  
family = binomial(link = "probit"), data = HMDA)  
coeftest(denyprobit2, vcov. = vcovHC, type = "HC1")
```

# Output

z test of coefficients:

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-2.258787	0.176608	-12.7898	< 2.2e-16***
pirat	2.741779	0.497673	5.5092	3.605e-08***
blackyes	0.708155	0.083091	8.5227	< 2.2e-16***

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

## Probit Model 2: Is there racial bias?

Table 5: Estimates from Probit regression model.

	Coef	Std Err	z-stat
Intercept	-2.259	0.137	-16.525
pirat	2.742	0.380	7.206
black	0.709	0.083	8.496

Null deviance: 1744.2 on 2379 degrees of freedom. Residual deviance: 1594.3 on 2377 degrees of freedom. AIC: 1600.3

# Estimated model:

$P(\text{deny}=1|x's) = \Phi(-2.259 + 2.742*(P/I) \text{ ratio} + 0.709*\text{black})$



## Probit Model 2

All coefficients are statistically significant. The coefficient for African Americans is positive and indicates that African Americans have a higher probability of denial than White Americans, *ceteris paribus*. Also, applicants with a high P/I ratio face a higher risk of being rejected.

Question: How big is the estimated difference in denial probabilities between two hypothetical applicants with the same payments-to-income ratio?

## Probit Model 2: Covariate Effect

We continue by using an augmented Probit model to estimate the effect of race on the probability of a mortgage application denial.

```
predictions2 <- predict(denyprob2, newdata  
= data.frame("black" = c("no", "yes"), "pirat" = c(0.3,0.3)),  
type = "response")
```

```
# 2. compute difference in probabilities  
diff(predictions2)
```

```
# Output  
0.1578117
```

We find that the white applicant faces a denial probability of only 7.546%, while the African American is rejected with a probability of 23.327%, a difference of 15.781 percentage points.

# Logit Model

As for Probit regression, there is no simple interpretation of the model coefficients and it is best to consider predicted probabilities or differences in predicted probabilities.

However, for the logit model the coefficients represent the log-odds (or logarithm of odds ratio) of success. As such, it is preferred in many applications.

We now analyze the mortgage lending problem using the Logit model.

## Logit Model 1

We continue by using an augmented Probit model to estimate the effect of race on the probability of a mortgage application denial.

```
# Estimating a logit model
denylogit <- glm(deny ~ pirat, family = binomial(link =
"logit"), data = HMDA)
coeftest(denylogit, vcov. = vcovHC, type = "HC1")

# Output
z test of coefficients:

              Estimate   Std. Error   z value   Pr(>|z|)
(Intercept) -4.02843     0.35898    -11.2218   < 2.2e-16***
pirat         5.88450     1.00015     5.8836    4.014e-09 ***

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

# Logit Model 1

Table 6: Estimates from Logit regression model.

	Coef	Std Err	z-stat
Intercept	-4.028	0.269	-15.000
pirat	5.884	0.734	8.021

Null deviance: 1744.2 on 2379 degrees of freedom. Residual deviance: 1660.2 on 2378 degrees of freedom. AIC: 1664.2

Both models produce very similar estimates of the probability that a mortgage application will be denied depending on the applicants payment-to-income ratio. This is shown in Figure 5 in the next slide.

## Probit & Logit: Probability of Denial

```
plot(x = HMDA$pirat, y = HMDA$deny,
#main = "Probit and Logit Models of the Probability of Denial,
Given P/I Ratio", xlab = "P/I ratio", ylab = "Deny",
pch = 20,      ylim = c(-0.4, 1.4), cex.main = 0.9)

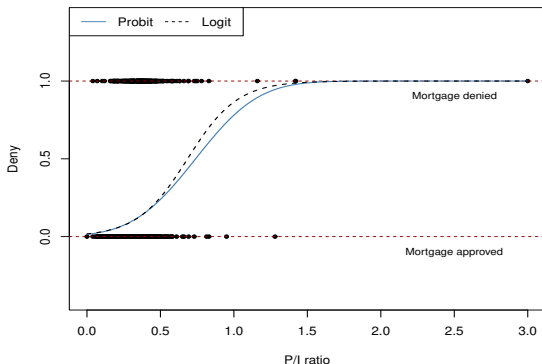
# add horizontal dashed lines and text
abline(h = 1, lty = 2, col = "darkred")
abline(h = 0, lty = 2, col = "darkred")
text(2.5, 0.9, cex = 0.8, "Mortgage denied")
text(2.5, -0.1, cex = 0.8, "Mortgage approved")
```

## Probit & Logit: Probability of Denial

```
# add estimated regression line of Probit and Logit models
x <- seq(0, 3, 0.01)
y_probit <- predict(denyprobit, list(pirat = x),
type = "response")
y_logit <- predict(denylogit, list(pirat = x), type = "response")
lines(x, y_probit, lwd = 1.5, col = "steelblue")
lines(x, y_logit, lwd = 1.5, col = "black", lty = 2)

# add a legend
legend("topleft", horiz = TRUE, legend = c("Probit",
"Logit"), col = c("steelblue", "black"), lty = c(1, 2))
```

# Probit & Logit: Probability of Denial



**Figure 5:** Probit and Logit models of the probability of denial, given the P/I ratio.



## Logit Model 2: Is there racial bias? Revisited

We now extend the Logit model of mortgage denial with the additional regressor black.

```
# estimate a Logit regression with multiple regressors
denylogit2 <- glm(deny ~ pirat + black, family
= binomial(link = "logit"), data = HMDA)
coeftest(denylogit2, vcov. = vcovHC, type = "HC1")
```

# Output

z test of coefficients:

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-4.12556	0.34597	-11.9245	< 2.2e-16***
pirat	5.37036	0.96376	5.5723	2.514e-08***
blackyes	1.27278	0.14616	8.7081	< 2.2e-16***

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

## Logit Model 2

Table 7: Estimates from Logit regression model.

	Coef	Std Err	z-stat
Intercept	-4.126	0.268	-15.370
pirat	5.370	0.728	7.374
black	1.279	0.146	8.706

Null deviance: 1744.2 on 2379 degrees of freedom. Residual deviance: 1591.4 on 2377 degrees of freedom. AIC: 1597.4

# Estimated model:

$P(\text{deny}=1|x's) = F(-4.126 + 5.370*(P/I) \text{ ratio} + 1.279*\text{black})$

## Logit Model 2

Similar to the Probit model, all model coefficients are highly significant and we obtain positive estimates for the coefficients on P/I ratio and black.

Question: How big is the estimated difference in denial probabilities between two hypothetical applicants with the same payments-to-income ratio?

## Logit Model 2

```
predictions3 <- predict(denylogit2, newdata =  
data.frame("black" = c("no", "yes"), "pirat"=c(0.3,0.3)),  
type = "response")
```

```
predictions3
```

```
# 2. Compute difference in probabilities  
diff(predictions3)
```

```
# Output  
0.1492945
```

We find that the white applicant faces a denial probability of only 7.485%, while the African American is rejected with a probability of 22.414%, a difference of 14.929 percentage points.

## Model Comparison

We compute the McFadden's R-squared or pseudo R-squared for the Probit and Logit models augmented with the additional regressor `black`.

```
# compute the null Probit model
denyprobit_null <- glm(formula = deny ~ 1, family
= binomial(link = "probit"), data = HMDDA)

# compute the pseudo-R2 using 'logLik'
1 - logLik(denyprobit2)[1]/logLik(denyprobit_null)[1] #> [1]

# Output
0.08594259
```



Thank you!