

Modelos de Regressão de Resposta Qualitativa

PROF. DR. VASCONCELOS WAKIM

Modelos de Regressão de Resposta Qualitativa

- Até o momento, trabalhou-se a ideia de Y da equação ser contínuo, ou seja, uma variável quantitativa;
- Neste tipo de modelo proposto hoje, buscar-se-á trabalhar a ideia de uma variável Y binária (0 ou 1)



Modelos de Regressão de Resposta Qualitativa

- Um modelo binário (variável Y) permite identificar uma situação de decisão de um indivíduo → trabalhar ou não trabalhar → usa-se uma variável dummy
- Pode-se querer estudar por exemplo, qual a probabilidade de uma pessoa votar em um partido de esquerda ou de direita, a partir de suas características pessoais.
- Assim, ter-se-ia o Y binário (0 = Esquerda e 1 = Direita) → X's podem ser binários ou contínuas



Modelos de Regressão de Resposta Qualitativa

- Pode-se ter também um Y tricotômico, com mais de duas opções (3 ou mais).
Por exemplo, votar em presidente da república, cujo ele pertença a: 0 = esquerda; 1 = Direita; 3 = Centro
- Qual a diferença entre os modelos tradicionais que tem o Y quantitativo para os modelos em que o Y é dicotômico?



Modelos de Regressão de Resposta Qualitativa

- A diferença é: no tradicional estamos estimando um valor médio de Y dados os valores que X assumem, podendo estes serem quantitativos ou qualitativos.
- Em modelos em que o Y é qualitativo (binário), buscamos identificar a **PROBABILIDADE** de que algo possa acontecer.
- Estes modelos binários são conhecidos como: **Modelos de Probabilidade**



Modelos de Regressão de Resposta Qualitativa

- Existem alguns modelos de resposta qualitativa (binárias), são eles:
- Modelo de Probabilidade Linear (MPL)
- Modelo Logit
- Modelo Probit
- Modelo Tobit



Modelo de Probabilidade Linear (MPL)

- Considere a seguinte expressão:

$$Y_i = \beta_1 + \beta_2 X_i + u_i$$

- X_i é a renda das famílias e $Y_i = 1$ se as famílias tem imóvel $Y_i = 0$ se as famílias não tem imóvel
- A expressão acima parece um MQO tradicional, mas como Y é binário, é chamado de MPL.



Modelo de Probabilidade Linear (MPL)

- Sendo a variável Y binária, tem-se que: $E(Y_i|X_i)$ deve ser interpretado como **PROBABILIDADE CONDICIONAL DE QUE O EVENTO OCORRA** dado os valores que X_i assume, isto porque, $\Pr(Y_i = 1|X_i)$.
- No caso anterior, estamos buscando saber qual a probabilidade de uma família ter casa, dado o seu nível de renda (X_i)



Modelo de Probabilidade Linear (MPL)

- Desta forma, assumindo que $E(u_i) = 0$, tem-se que:
- $E(Y_i|X_i) = \beta_1 + \beta_2 X_i$
- Se P_i = probabilidade de que $Y_i = 1$ (de que o evento ocorra) e que $1 - P_i$ = probabilidade do evento não acontecer $Y_i = 0$ tem a seguinte distribuição de probabilidade:
- Y_i segue uma distribuição de probabilidade de Bernoulli

Y_i	Probabilidade
0	$1 - P_i$
1	P_i
Total	1



Modelo de Probabilidade Linear (MPL)

- Aplicando Esperança Matemática têm-se:
- $E(Y_i) = 0(1 - P_i) + 1(P_i) = P_i$
- Se comprarmos: $E(Y_i|X_i) = \beta_1 + \beta_2 X_i$ com $E(Y_i) = 0(1 - P_i) + 1(P_i) = P_i$ pode-se igualar obtendo:
- $E(Y_i|X_i) = \beta_1 + \beta_2 X_i = P_i$



Modelo de Probabilidade Linear (MPL)

- A distribuição de probabilidade Bernoulli prevê que a probabilidade de ocorrência do valor 1 é totalmente aleatório.
- Se houver n experimentos, cada evento tem probabilidade p de sucesso e de $(1 - p)$ de insucesso.
- P_i deve estar entre 0 e 1
- $0 \leq E(Y_i|X_i) \leq 1$



Modelo de Probabilidade Linear (MPL)

- O MPL apresenta ausência de normalidade dos resíduos, pois estes também assumiram apenas dois valores (1 ou 0) conforme o Y ;
- $u_i = Y_i - \beta_1 - \beta_2 X_i$, sendo que a distribuição de probabilidade de X_i é:

	u_i	Probabilidade
Quando $Y_i = 1$	$1 - \beta_1 - \beta_2 X_i$	P_i
Quando $Y_i = 0$	$-\beta_1 - \beta_2 X_i$	$(1 - P_i)$

- Não se pode inferir que u_i é normalmente distribuído, porque também segue distribuição de Bernoulli.



Modelo de Probabilidade Linear (MPL)

- O MPL apresenta variância heterocedástica no termo de erro (u_i)
- Mesmo se $E(u_i) = 0$ e a $Cov(u_i, u_j) = 0 \forall i \neq j$, os termos de erro serão heterocedásticos.
- Aplicando o conceito de Variância
- $Var(u_i) = P_i(1 - P_i)$, no MPL o termo de erro é heterocedástico. Para corrigir este problema, basta fazer uma ponderação no modelo MPL



Modelo de Probabilidade Linear (MPL)

- $\sqrt{E(Y_i|X_i)[1 - E(Y_i|X_i)]} = \sqrt{P_i(1 - P_i)} = \sqrt{w_i}$

$$\frac{Y_i}{\sqrt{w_i}} = \frac{\beta_1}{\sqrt{w_i}} + \beta_2 \frac{X_i}{\sqrt{w_i}} + \frac{u_i}{\sqrt{w_i}}$$

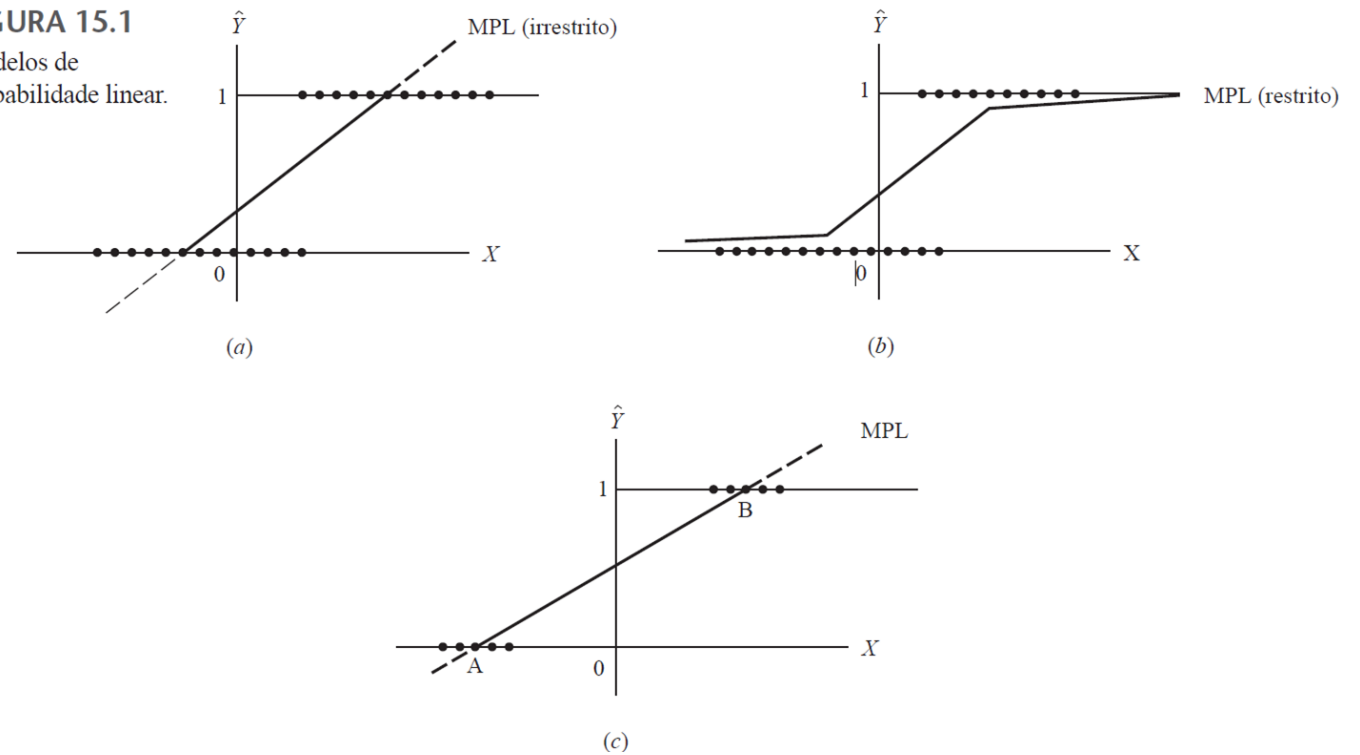
- Com esta ponderação, o termo erro torna-se homocedástico. Assim, pode-se estimar o MPL por meio do MQO tradicional que também pode ser chamado de MQP (Mínimos quadrados ponderados)



Modelo de Probabilidade Linear (MPL)

FIGURA 15.1

Modelos de probabilidade linear.

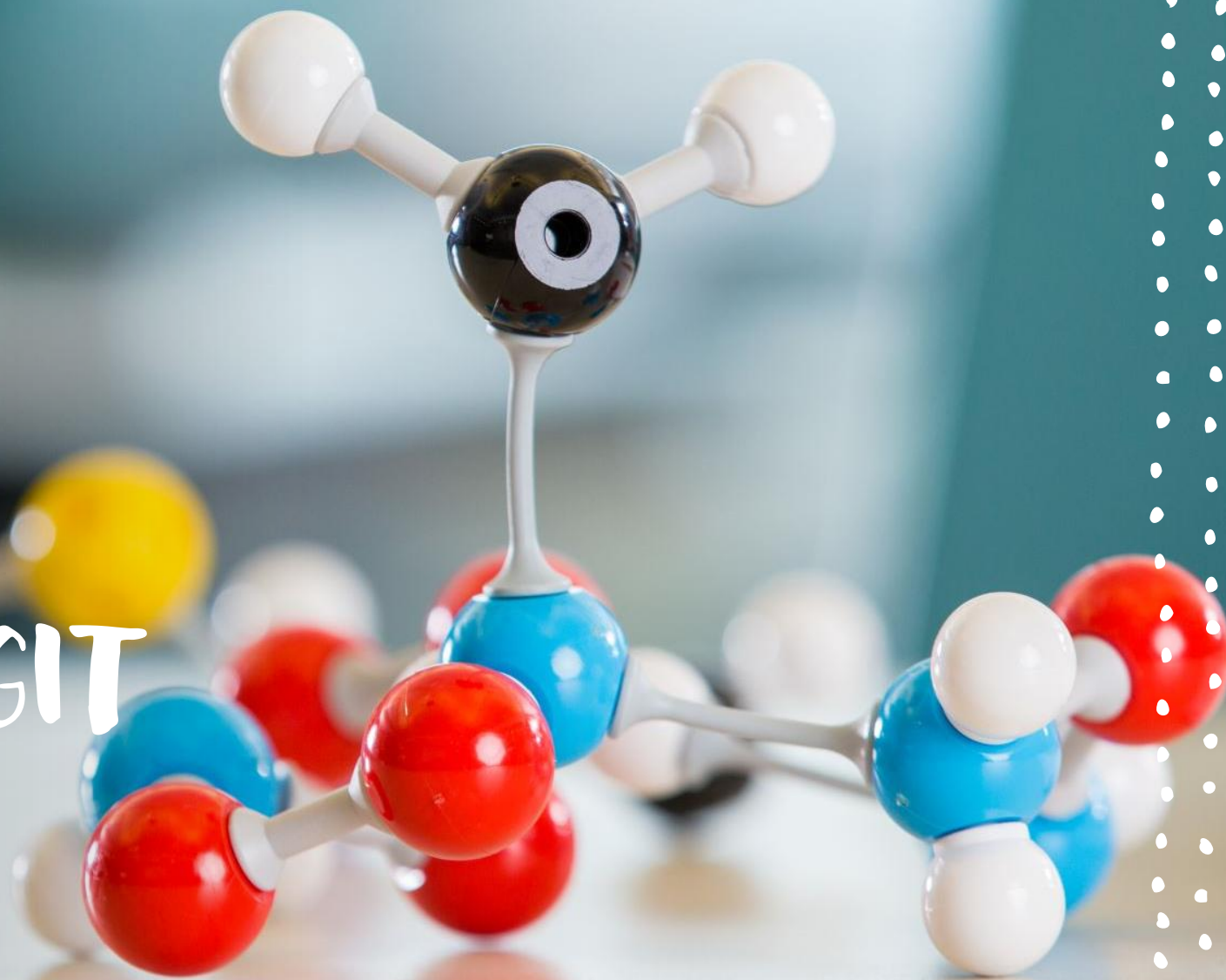


- Em modelos dicotômicos, o R^2 passa a não ser uma medida interessante para verificação do ajuste do modelo, por o R^2 vai ficar muito abaixo de 1.

Por isso, tem-se modelos alternativos de estimação → LOGIT, PROBIT e TOBIT



MODELO LOGIT

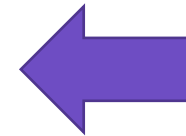


Modelo LOGIT

$$P_i = \beta_1 + \beta_2 X_i$$

- X_i é a renda da família e $P_i = E(Y_i = 1|X_i)$ mostra a probabilidade da família ter casa própria. Esta representação pode ser reescrita como:
- $P_i = \frac{1}{1+e^{-(\beta_1+\beta_2 X_i)}}$ esta expressão pode ser reescrita como:

$$P_i = \frac{1}{1 + e^{-z_i}} = \frac{e^z}{1 + e^z}$$



Esta função é conhecida como
Função de Distribuição Logística

$$z_i = \beta_1 + \beta_2 X_i$$



Modelo LOGIT

- Z_i varia de $-\infty$ a $+\infty$
- P_i varia de 0 a 1 e estão não linearmente com Z_i
- A expressão $P_i = \frac{1}{1+e^{-(\beta_1+\beta_2 X_i)}}$ pode ser linearizada para ser estimada por MQO
- Se P_i é a prob de se ter casa própria e $1 - P_i$ é a prob de não ter, então têm-se:

$$1 - P_i = \frac{1}{1+e^{Z_i}}$$

Reescrevendo



Modelo LOGIT

$$\frac{P_i}{1-P_i} = \frac{1+e^{Z_i}}{1+e^{-Z_i}}$$

- Em que $\frac{P_i}{1-P_i}$ é a razão de chances em favor de se ter casa própria.
- Se $P_i = 0,8$ significa que as chances são de 4 para 1 a favor de a família ter casa própria.
- Tomando o log na equação acima temos:



Modelo LOGIT

$$L_i = \ln\left(\frac{P_i}{1-P_i}\right) = Z_i$$

$$Z_i = \beta_1 + \beta_2 X_i$$

$$L_i = \ln\left(\frac{P_i}{1-P_i}\right) = \beta_1 + \beta_2 X_i$$

EQUAÇÃO BÁSICA DO MODELO LOGIT



ESTIMAÇÃO DO MODELO LOGIT



Estimação do Modelo LOGIT

$$L_i = \ln\left(\frac{P_i}{1-P_i}\right) = \beta_1 + \beta_2 X_i + u_i$$

- Partimos de 2 opções para estimação:
- 1) Dados em nível e/ou individual ou micro
- 2) Dados agrupados ou replicados

Dados em nível e/ou individual ou micro

- Se o pesquisador possui dados individual de cada indivíduo da amostra, fazer a estimação por meio do MQO é inviável. Se $P_i = 1$ família ter casa própria e $P_i = 0$ não ter casa própria, colocando estes valores na equação logística, têm-se:

$$L_i = \ln\left(\frac{1}{0}\right)$$

Ter casa própria

$$L_i = \ln\left(\frac{0}{1}\right)$$

Não ter casa própria

Por meio do MQO não faz sentido a estimação, assim, deve-se utilizar o **Método de Máxima Verossimilhança (MV)** para estimar os parâmetros.



Dados Agrupados ou Replicados

- Se o pesquisador possui dados agregados das famílias de forma agrupada ou replicada (observações repetidas), calcula-se a probabilidade de ter casa própria a partir da expressão:
- $\hat{P}_i = \frac{n_i}{N_i}$
- Em que N_i é o número de famílias da amostra e n_i número de famílias donas de casas próprias

Dados Agrupados ou Replicados

- Partindo desta lógica, pode-se estimar um logit estimado:

$$\hat{L}_i = \ln\left(\frac{\hat{P}_i}{1-\hat{P}_i}\right) = \hat{\beta}_1 + \hat{\beta}_2 X_i$$

- Desta forma obtêm-se uma boa estimativa para L_i
- Considerando que \hat{u}_i é binário:

$$u_i \sim N \left[0, \frac{1}{N_i P_i (1 - P_i)} \right]$$

Dados Agrupados ou Replicados

- No caso do MPL o termo de erro é heterocedástico, deve-se estimar via MQP, assim reescrevemos a variância como:

$$\hat{\sigma}^2 = \frac{1}{N_i \hat{P}_i (1 - \hat{P}_i)}$$

- Quando reescrevemos desta forma, estamos ponderando por w_i as variáveis e contornando o problema de heterocedasticidade do modelo.





Medida de Ajustamento do Modelo



Medida de Ajustamento

- Normalmente o Logit é estimado pelo método de Máxima Verossimilhança que é um método para grandes amostras e erros padrões assintóticos;
- Usamos a estatística Z (normal) ao invés da t para analisar a significância dos coeficientes;
- R^2 não é adequado para modelos binários. Deve-se usar o **pseudo R^2** , também chamado de **R^2 de MacFadden**.
- Outra medida de ajuste é o **Count R^2** . Que busca medir a relação de acertos do modelo quando o 1 acontece.



Medida de Ajustamento

- Matematicamente: $Count R^2 = \frac{n^{\circ} \text{ de previsões corretas}}{n^{\circ} \text{ total de observações}}$





*Interpretando os
resultados
encontrados*



Análise das variáveis via Efeito Marginal

- Considerando que o logit é um modelo não linear, a interpretação dos betas estimados não é de forma direta conforme ocorre no MCRL.
- Para isto deve-se calcular o chamado **Efeito Marginal (MFX)**, que mede o impacto de variação de uma unidade na variável explicativa sobre a probabilidade de ocorrência do evento $P_i(Y_i = 1)$
- Normalmente o Efeito Marginal é dado pela deriva da função em relação à variável



Análise das variáveis via Efeito Marginal

- $EM_{X_k} = \frac{\partial P_i}{\partial X_k}$

- Em que $P_i = P(Y = 1)$ e $P_i = F(X_i\beta)$

- $EM_{X_k} = \frac{\partial F(X_i\beta)}{\partial X_k} = \frac{\partial F(X_i\beta)}{\partial (X_i\beta)} \times \frac{\partial (X_i\beta)}{\partial X_k} = f(X_i\beta) \times \frac{\partial (X_i\beta)}{\partial X_k}$

VARIÁVEIS CONTÍNUAS



Análise das variáveis via Efeito Marginal

- O efeito marginal para variáveis dummies é a mudança na $P(Y = 1)$ quando a variável *dummy* D_j passa de 0 para 1.

$$EM_{X_k} = P(Y = 1|D = 1) - P(Y = 1|D = 0)$$



Análise das variáveis via Efeito Marginal

- O efeito marginal no modelo Logit e Probit, variam de acordo com os valores das variáveis explicativas, sendo assim, podem ser calculados de 3 formas:
- **1) Média do Efeito Marginal:** avalia-se o efeito para cada observação e depois calcula-se a média;
- **2) Efeito Marginal Médio:** calculado com os valores médios das variáveis explicativas;
- **3) Efeito marginal em postos específicos da amostra:** calculado com valores específicos das variáveis explicativas. Uma variante do nº 2.





*Avaliação do modelo
estimado*

Avaliação do modelo estimado

- Tecnicamente, avalia-se o modelo seguindo a mesma lógica de um MCRL;
- Avalia-se: sinais dos coeficientes. No caso do modelo logit e probit, o efeito de X sobre Y somente poderá ser calculado a partir do Efeito Marginal. (Ele não é direto);
- Faz-se o teste de significância dos coeficientes (estatística Z normal) idem p-value no MCRL;



Avaliação do modelo estimado

- Teste de significância global do modelo. Este é semelhante ao teste F, em que busca estimar um modelo irrestrito (com todas as variáveis explicativas) e outro apenas com a constante e aplica-se um teste de qui-quadrado.

$$X_c^2 = -2 \ln \lambda = 2(\ln L_{IR} - \ln L_R) \sim X_{k-1}^2$$

$$\lambda = \frac{L_R}{L_{IR}}$$

Ln = logaritmo natural e k

= nº de parâmetros estimados do modelo irrestrito

L_R = valor da função de verossimilhança d modelo restrito

L_{IR} = valor da função de verossimilhança d modelo irrestrito



Avaliação do modelo estimado

- Medidas de qualidade do ajustamento no modelo binário é o R^2 de MacFadden (1974)

$$R_{McFadden}^2 = 1 - \frac{\ln L}{\ln L_0}$$

- $\ln L$ é o log da função de verossimilhança com as variáveis explicativas; $\ln L_0$ é o log da função de verossimilhança apenas com o intercepto
- $R_{McFadden}^2$ é conhecido também como Índice de Razão de Verossimilhança (*likelihood ratio index*) - varia entre 0 e 1. interpretado como o R^2 do MCRL. Mas ele não explica a proporção da variação de Y.



Avaliação do modelo estimado

- Proporção de previsões corretas. Também semelhante ao R^2 do MCRL. É dado pela proporção de previsões corretas feitas pelo modelo Logit.
- Sendo $Y_i = 0$ ou $Y_i = 1$ define-se uma regra:
- $\hat{Y}_i = 1$ se $\hat{P}_i > 0,5$
- $\hat{Y}_i = 0$ se $\hat{P}_i \leq 0,5$
- \hat{Y}_i é binária que representa Y_i , onde busca-se contar os n° de acertos quando $\hat{Y}_i = 0$ ou $\hat{Y}_i = 1$



Avaliação do modelo estimado

- Os valores de **A** e **D** → são as previsões corretas do modelo;
- **B** e **C** → são as previsões erradas
- A proporção de acerto do modelo é dada por:

$$Previsão_{correta} = \left(\frac{A + D}{n} \right) \times 100$$

Total	Previsão		<i>n</i> Total
	E	F	
Observado	$Y_i = 0$	$Y_i = 1$	
$Y_i = 0$	A	B	n_0
$Y_i = 1$	C	D	n_1

A equação mostra o percentual de acerto das previsões do modelo, e não o % de variação de Y .





Aplicação do Modelo Logit

Aplicação

Considere uma amostra de dados relacionada com a inserção da mulher no mercado de trabalho contendo informações sobre 2000 mulheres das quais 1343 estão empregadas e 657 estão desempregadas (Baum, 2006, p.251). Dispõe-se de dados relacionados com as seguintes variáveis:

wk = participação na força de trabalho (empregada = 1)

kids = número de crianças

age = idade em anos

esc = escolaridade em anos

cas = estado civil (casada = 1)

wage = salário da mulher

lw = log salário da mulher



1. Estadísticas Descriptivas

. sum

Variable	Obs	Mean	Std. Dev.	Min	Max
age	2000	36.208	8.28656	20	59
esc	2000	13.084	3.045912	10	20
cas	2000	.6705	.4701492	0	1
kids	2000	1.6445	1.398963	0	5
wage	1343	23.69217	6.305374	5.88497	45.80979
lw	1343	3.126703	.2865111	1.772402	3.824498
wk	2000	.6715	.4697852	0	1

2. Estimação do Modelo Logit/Probit

A estimação pode ser feita usando o *menu* ou por meio dos comandos.

. Pelo *menu*: Statistics → Binary outcome → Logistic regression/Probit regression →

definir variáveis → ok

. Pelo comando *logit*: *logit wk age esc cas kids*

```
. logit wk age esc cas kids
```

```
Iteration 0:  log likelihood = -1266.2225
Iteration 1:  log likelihood = -1040.6658
Iteration 2:  log likelihood = -1027.9567
Iteration 3:  log likelihood = -1027.9145
Iteration 4:  log likelihood = -1027.9144
```

```
Logistic regression
```

```
Number of obs   =      2000
LR chi2(4)      =      476.62
Prob > chi2     =      0.0000
Pseudo R2      =      0.1882
```

```
Log likelihood = -1027.9144
```

wk	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
age	.0579303	.007221	8.02	0.000	.0437773 .0720833
esc	.0982513	.0186522	5.27	0.000	.0616936 .134809
cas	.7417775	.1264705	5.87	0.000	.4938998 .9896552
kids	.7644882	.0515289	14.84	0.000	.6634935 .865483
_cons	-4.159247	.3320401	-12.53	0.000	-4.810034 -3.508461



3. Razão de Chances e *odds ratio*

Uma interpretação de mais fácil entendimento, *no caso do logit*, é pela *odds ratio* que é obtida tomando o antilogaritmo ($e^{\hat{\beta}}$) dos coeficientes. O Stata oferece a opção do resultado em termos de *odds ratio*. Existem duas formas de obter os *odds ratios*. Uma é pela opção *or* no logit e a outra é pelo comando *logistic* ao invés de *logit*. Probit não tem *odds ratio*.

```
. logit wk age esc cas kids, or
```

```
Iteration 0:  log likelihood = -1266.2225
Iteration 1:  log likelihood = -1040.6658
Iteration 2:  log likelihood = -1027.9567
Iteration 3:  log likelihood = -1027.9145
Iteration 4:  log likelihood = -1027.9144
```

```
Logistic regression
```

```
Number of obs   =      2000
LR chi2(4)      =      476.62
Prob > chi2     =      0.0000
Pseudo R2      =      0.1882
```

```
Log likelihood = -1027.9144
```

wk	Odds Ratio	Std. Err.	z	P> z	[95% Conf. Interval]
age	1.059641	.0076517	8.02	0.000	1.04475 1.074745
esc	1.10324	.0205779	5.27	0.000	1.063636 1.144318
cas	2.099664	.2655457	5.87	0.000	1.638694 2.690307
kids	2.147895	.1106786	14.84	0.000	1.941563 2.376153
_cons	.0156193	.0051862	-12.53	0.000	.0081476 .029943

Com o comando OR após a rotina, os betas estimados proporcionam a razão de chances do evento ocorrer. Por exemplo, a variável **CAS** que mostra que as mulheres casadas tem 2,1 vezes mais chances de estar no mercado de trabalho do que a solteira, mantendo as demais constantes.



3. Razão de Chances e *odds ratio*

Logistic regression

Number of obs = 2000

LR chi2(4) = 476.62

Prob > chi2 = 0.0000

Pseudo R2 = 0.1882

Log likelihood = -1027.9144

wk	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
age	.0579303	.007221	8.02	0.000	.0437773	.0720833
esc	.0982513	.0186522	5.27	0.000	.0616936	.134809
cas	.7417775	.1264705	5.87	0.000	.4938998	.9896552
kids	.7644882	.0515289	14.84	0.000	.6634935	.865483
_cons	-4.159247	.3320401	-12.53	0.000	-4.810034	-3.508461

Outra forma é pelo antilog do beta $(e^{\beta} - 1) * 100$. Exemplo:
CAS BETA = 0,7417775. $(e^{0,7417775} - 1) * 100 = 110\%$. Este resultado mostra que as chances das mulheres casadas estarem no mercado de trabalho é 110% maior do que as chances das solteiras.



4. Efeito Marginal

O efeito marginal da variável explicativa pode ser calculado pelos comandos:

`. mfx` – calcula o efeito marginal no ponto médio da amostra ($X = \bar{X}$).

`. mfx, at()` – calcula o efeito marginal em pontos específicos

`. margeff` – calcula efeito marginal médio.

`. margeff, at()` – calcula efeito marginal médio em pontos específicos.

```
. mfx
```

```
Marginal effects after logit
```

```
  y = Pr(wk) (predict)
```

```
  = .72678588
```

variable	dy/dx	Std. Err.	z	P> z	[95% C.I.]	X
age	.0115031	.00142	8.08	0.000	.008713	.014293		36.208
esc	.0195096	.0037	5.27	0.000	.01226	.02676		13.084
cas*	.1545671	.02703	5.72	0.000	.101592	.207542		.6705
kids	.151803	.00938	16.19	0.000	.133425	.170181		1.6445

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



Marginal effects after logit

y = Pr(wk) (predict)

= .72678588

variable	dy/dx	Std. Err.	z	P> z	[95% C.I.]	X
age	.0115031	.00142	8.08	0.000	.008713	.014293		36.208
esc	.0195096	.0037	5.27	0.000	.01226	.02676		13.084
cas*	.1545671	.02703	5.72	0.000	.101592	.207542		.6705
kids	.151803	.00938	16.19	0.000	.133425	.170181		1.6445

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

No caso da variável **CAS** ela é binária (0 e 1), o comando mfx calcula o efeito marginal pela mudança de 0 para 1.

A interpretação é: considerando o ponto médio da amostra, um aumento da unidade aumenta a probabilidade da mulher estar no mercado de trabalho em 1,1 ponto percentual;

Um aumento de um na escolaridade, a probabilidade da mulher estar empregada é de 1,9 pontos percentuais.



a) Estatísticas e Tabela de Classificação

. estat classification (ou estat class)

`. estat class`

Logistic model for wk

Classified	True		Total
	D	~D	
+	1177	361	1538
-	166	296	462
Total	1343	657	2000

Classified + if predicted $\Pr(D) \geq .5$

True D defined as wk != 0

Sensitivity	Pr(+ D)	87.64%
Specificity	Pr(- ~D)	45.05%
Positive predictive value	Pr(D +)	76.53%
Negative predictive value	Pr(~D -)	64.07%

False + rate for true ~D	Pr(+ ~D)	54.95%
False - rate for true D	Pr(- D)	12.36%
False + rate for classified +	Pr(~D +)	23.47%
False - rate for classified -	Pr(D -)	35.93%

Correctly classified 73.65%

`. fitstat`

Measures of Fit for logit of wk

Log-Lik Intercept Only:	-1266.223	Log-Lik Full Model:	-1027.914
D(1995):	2055.829	LR(4):	476.616
McFadden's R2:	0.188	Prob > LR:	0.000
ML (Cox-Snell) R2:	0.212	McFadden's Adj R2:	0.184
McKelvey & Zavoina's R2:	0.349	Cragg-Uhler (Nagelkerke) R2:	0.295
Variance of y*:	5.052	Efron's R2:	0.216
Count R2:	0.737	Variance of error:	3.290
AIC:	1.033	Adj Count R2:	0.198
BIC:	-13107.972	AIC*n:	2065.829
BIC used by Stata:	2093.833	BIC':	-446.213
		AIC used by Stata:	2065.829