

Mery Natali Silva Abreu^{I,II}

Arminda Lucia Siqueira^{III}

Waleska Teixeira Caiaffa^{I,II}

Regressão logística ordinal em estudos epidemiológicos

Ordinal logistic regression in epidemiological studies

RESUMO

Os modelos de regressão logística ordinal vêm sendo aplicados com sucesso na análise de estudos epidemiológicos. Entretanto, a verificação da adequação de cada modelo tem recebido atenção limitada. O artigo apresenta uma breve análise dos principais modelos de regressão logística ordinal e as estratégias para ajuste s, as técnicas de verificação de qualidade do ajuste, bem como os comandos para execução nos softwares R e Stata. A metodologia é ilustrada com aplicação dos dados do *Second National Health and Nutrition Examination Survey* (NHANES II), o conhecido levantamento de saúde e nutrição.

DESCRITORES: Estatística como Assunto. Modelos Logísticos. Análise de Regressão. Métodos Epidemiológicos.

ABSTRACT

Ordinal logistic regression models have been developed for analysis of epidemiological studies. However, the adequacy of such models for adjustment has so far received little attention. In this article, we reviewed the most important ordinal regression models and common approaches used to verify goodness-of-fit, using R or Stata programs. We performed formal and graphical analyses to compare ordinal models using data sets on health conditions from the National Health and Nutrition Examination Survey (NHANES II).

DESCRIPTORS: Statistics as Topic. Logistic Models. Regression Analysis. Epidemiologic Methods.

^I Programa de Pós-graduação em Saúde Pública. Faculdade de Medicina (FM). Universidade Federal de Minas Gerais (UFMG). Belo Horizonte, MG, Brasil

^{II} Grupo de Pesquisa em Epidemiologia e Observatório de Saúde Urbana. FM-UFMG. Belo Horizonte, MG, Brasil

^{III} Departamento de Estatística. Instituto de Ciências Exatas. UFMG. Belo Horizonte, MG, Brasil

Correspondência | Correspondence:

Mery Natali Silva Abreu
Av. Alfredo Balena, 190, 6º andar
Sala 625 Santa Efigênia
31130-100 Belo Horizonte, MG, Brasil
E-mail: merynatali@yahoo.com.br

Recebido: 10/5/2007

Revisado: 23/5/2008

Aprovado: 11/6/2008

INTRODUÇÃO

Os modelos de regressão logística ordinal vêm sendo aplicados nos últimos anos na análise de dados cuja resposta ou desfecho é apresentado em categorias com ordenação. A informação ordenada, na forma de escore, tem sido cada vez mais utilizada em estudos epidemiológicos, tais como, qualidade de vida em escalas intervalares, indicadores de condição de saúde e mesmo de gravidade das doenças.¹ Estes modelos, dependendo do delineamento do estudo, permitem também calcular a estatística *odds ratio* (OR) ou a probabilidade de ocorrência de um evento.¹

Existem vários modelos ordinais tais como o modelo de *odds* proporcional, modelo de *odds* proporcionais parciais, modelo de razão contínua e modelo estereótipo. Apesar dessa diversidade e da grande variedade de estudos^{1-4,7-15} sobre o assunto, a sua utilização na área de saúde pública ainda é escassa.¹ Esta constatação pode ser atribuída não só à complexidade, mas especialmente à dificuldade da validação de seus pressupostos.¹¹ Outro fator que poderia estar relacionado à pouca utilização destes modelos são as reduzidas opções de modelagem oferecidas em pacotes estatísticos comerciais utilizados na área de saúde pública, entre eles o SPSS e Minitab. Mesmo utilizando outros mais complexos tais como SAS e Stata, são freqüentes as dificuldades em selecionar os comandos apropriados e na interpretação dos resultados.³ Adiciona-se a isso o custo elevado, devido as licenças altamente restritas da maioria dos pacotes estatísticos comerciais.

Um pacote estatístico que vem se tornando cada vez mais popular é o *software* livre R,^{6a} distribuído sob licença pública geral. O pacote possui uma variedade de técnicas estatísticas, incluindo vários modelos de regressão logística ordinal, permitindo que eles sejam testados e o ajuste comparado.

O presente artigo teve como objetivo analisar o ajuste e a adequação dos principais modelos de regressão ordinal, mostrando os comandos de execução no *software* R. Complementando essa avaliação/análise, foi feito o ajuste do modelo de *odds* proporcionais parciais por meio do *software* Stata, pois ainda não se encontra implementado no *software* R.

Para exemplificar os métodos analisados, será utilizada parte dos dados do conhecido levantamento nacional de saúde e nutrição, intitulado *Second National Health and Nutrition Examination Survey* (NHANES II), disponíveis na internet^b e amplamente utilizadas como exemplos em estudos estatísticos e epidemiológicos. Este inquérito inclui informações demográficas,

antropométricas, história nutricional e de saúde e hematologia. Foram excluídas as informações referentes às crianças, totalizando 10.337 entrevistados com idades entre 20 e 74 anos no banco de dados.

ANÁLISE UNIVARIADA

Como em qualquer procedimento analítico que utiliza modelos de regressão, a análise múltipla por meio dos modelos ordinais deve sempre ser precedida pelo cruzamento de cada covariável com o evento de interesse. Por meio dessa análise, conhecida como univariada, é possível selecionar os fatores que serão introduzidos no modelo de regressão.

O qui-quadrado de tendência é um dos testes apropriados para seleção dos efeitos principais, já que considera o caráter ordinal da variável resposta. Normalmente, utiliza-se um nível de significância conservador (geralmente entre 10% e 25%) para entrada das covariáveis no modelo.¹⁰

Além disso, pode-se estimar OR, considerando uma categoria da variável-resposta como referência e comparando-a com as demais ou agrupando as categorias maiores e comparando-as às categorias menores.

MODELOS DE REGRESSÃO ORDINAL

Após a análise univariada, deve-se construir o modelo final de regressão múltipla, para controlar possíveis fatores de confusão. Como o evento de interesse é ordinal, deve-se utilizar um modelo de regressão logística ordinal.

Seja Y a variável resposta com k categorias codificadas em $1, 2, \dots, k$ e $\underline{x} = (x_1, x_2, \dots, x_p)$ o vetor de variáveis explicativas ou covariáveis. As k categorias de Y condicionalmente aos valores de \underline{x} ocorrem com probabilidades p_1, p_2, \dots, p_k , isto é, $p_j = P(Y=j)$, para $j=1, 2, \dots, k$. O termo α refere-se ao intercepto do modelo e β corresponde aos efeitos das covariáveis na variável-resposta. Na Tabela 1 são apresentadas as formas dos principais modelos, indicação de uso e comandos para execução desses nos pacotes R e Stata.

Modelo de *Odds* Proporcionais (MOP)

No MOP são considerados $(k - 1)$ pontos de corte das categorias sendo que o j -ésimo ($j=1, \dots, k-1$) ponto de corte é baseado na comparação de probabilidades acumuladas, como mostrado na Tabela 1.

^a The R Project for Statistical Computing [internet]. Viena: Viena University of Economics and Business Administration, [s.d.] [citado 2008 set 7]. Disponível em: <http://www.r-project.org/about.htm>

^b National Center for Health Statistics. Publications and Information Products: NHANES II public-use and data files. [citado 2008 nov 11] Disponível em: http://www.cdc.gov/nchs/products/elec_prods/subject/nhanesii.htm

Tabela 1. Informações sobre os principais modelos de regressão logística ordinal, indicação de uso e comandos para execução desses nos pacotes R e Stata, exemplificando com o estudo NHANES II.

Modelo	Forma funcional do modelo	Indicação de uso	Comandos R*	Comandos Stata
MOP	$\lambda_j(\underline{x}) = \ln \left\{ \frac{\Pr(Y = 1 \underline{x}) + \dots + \Pr(Y = j \underline{x})}{\Pr(Y = j + 1 \underline{x}) + \dots + \Pr(Y = k \underline{x})} \right\} = \ln \left\{ \frac{\sum_{j=1}^k \Pr(Y = j \underline{x})}{\sum_{j=1}^k \Pr(Y = j \underline{x})} \right\}$ $\lambda_j(\underline{x}) = \alpha_j + (\beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p), \quad j = 1, \dots, k-1$	Variável resposta original contínua e posteriormente agrupada e a suposição de chances proporcionais é válida	<pre>> mcp = lrm(saude~ idade + diabetes + cor_ pele, nhanes, x=TRUE, y=TRUE) > mcp > par(mfrow=c(1,5)) > residuals(mcp, type = 'score.binary', pl=TRUE) > residuals(mcp, type= 'partial', pl=TRUE)</pre>	<pre>. ologit saude idade diabetes cor_pele</pre>
MOPP-NR	$\lambda_j(\underline{x}) = \ln \left\{ \frac{\Pr(Y = 1 \underline{x}) + \dots + \Pr(Y = j \underline{x})}{\Pr(Y = j + 1 \underline{x}) + \dots + \Pr(Y = k \underline{x})} \right\} = \ln \left\{ \frac{\sum_{j=1}^k \Pr(Y = j \underline{x})}{\sum_{j=1}^k \Pr(Y = j \underline{x})} \right\}$ $\lambda_j(\underline{x}) = \alpha_j + [(\beta_1 + \gamma_1)x_1 + \dots + (\beta_q + \gamma_q)x_q + (\beta_{q+1}x_{q+1}) + \dots + (\beta_p x_p)], j = 1, \dots, k-1$	Quando a suposição de chances proporcionais não é válida	-	<pre>. gologit2 saude idade cor_pele diabetes, autoit lforce</pre>
MOPP-R	$\lambda_j(\underline{x}) = \ln \left\{ \frac{\Pr(Y = 1 \underline{x}) + \dots + \Pr(Y = j \underline{x})}{\Pr(Y = j + 1 \underline{x}) + \dots + \Pr(Y = k \underline{x})} \right\} = \ln \left\{ \frac{\sum_{j=1}^k \Pr(Y = j \underline{x})}{\sum_{j=1}^k \Pr(Y = j \underline{x})} \right\}$ $\lambda_j(\underline{x}) = \alpha_j + \tau_j [(\beta_1 + \gamma_1)x_1 + \dots + (\beta_q + \gamma_q)x_q] + (\beta_{q+1}x_{q+1}) + \dots + (\beta_p x_p), j = 1, \dots, k-1$	Suposição de chances proporcionais não é válida e existe relação linear entre OR de uma covariável e a variável resposta	-	-
MRC	$\lambda_j(\underline{x}) = \ln \left\{ \frac{\Pr(Y = j \underline{x})}{\Pr(Y = j + 1 \underline{x}) + \dots + \Pr(Y = k \underline{x})} \right\} = \ln \left\{ \frac{\Pr(Y = j \underline{x})}{\sum_{j=1}^k \Pr(Y = j \underline{x})} \right\}$ $\lambda_j(\underline{x}) = \alpha_j + (\beta_{j1}x_1 + \beta_{j2}x_2 + \dots + \beta_{jp}x_p), \quad j = 1, \dots, k$	Há um interesse intrínseco em uma categoria específica da variável resposta	<pre>> cr=cr.setup(nhanes\$saude) > sau.cr=cr\$y > cohort=cr\$cohort > ida.cr=idade[cr\$subs] > diab.cr = diabetes [cr\$subs] > cor.cr =cor_pele[cr\$subs] > mrc=lrm(sau.cr~ida.cr+ diab.cr+cor.cr+cohort) > mrc</pre>	<pre>. ocratio saude idade diabetes cor_pele</pre>
ME	$\lambda_j(\underline{x}) = \ln \left\{ \frac{\Pr(Y = j \underline{x})}{\Pr(Y = 0 \underline{x})} \right\}$ $\lambda_j(\underline{x}) = \alpha_j + \omega_j (\beta_1 x_1 + \dots + \beta_p x_p), \quad j = 1, \dots, k$	Variável resposta ordinal discreta que não provem de alguma variável contínua agrupada	<pre>> install.packages ("VGAM", repos="http://www.stat.auckland. ac.nz/~yee/") > library(VGAM) > ms = rrvglm(saude ~ idade + diabetes + cor_pele, nhanes, multinomial) > summary(ms)</pre>	<pre>. slogit saude idade diabetes cor_pele</pre>

MOP=ModeloOddsProporcionais;MOPP-NR=ModeloOddsProporcionaisParciaisNãoRestrito;MOPP-NR=ModeloOddsProporcionaisParciaisRestrito;MRC=ModelodeRazãoContínua;ME=ModeloEstereótipo
 *Estudo NHANES - variável resposta: saude; covariáveis: idade, diabetes e cor_pele
 NHANES: Second National Health and Nutrition Examination Survey

O termo α_j varia para cada uma das k categorias e cada β não depende do índice j , implicando que a relação entre X e Y é independente da categoria.

Logo, o modelo possui uma suposição de *odds* proporcionais acerca dos $(k-1)$ pontos de corte, também chamada de suposição de regressão paralela, que é assumida para cada covariável incluída no modelo. Essa suposição deve ser testada para cada covariável separadamente e no modelo final, utilizando, por exemplo, o teste escore.¹⁰

Esse modelo é apropriado para analisar variáveis ordinais, provenientes de uma variável contínua que foi, por sua vez, agrupada.

Modelo de Odds Proporcionais Parciais (MOPP)

Como a suposição de *odds* proporcionais é difícil de ser alcançada na prática, alternativamente pode ser utilizado o MOPP.¹³ Este modelo permite que algumas covariáveis possam ser modeladas com a suposição de *odds* proporcional mas, para as variáveis em que essa suposição não é satisfeita, é incrementado por um coeficiente (γ) que é o efeito associado com cada j -ésimo logito cumulativo, ajustado pelas demais covariáveis.¹⁰ A forma geral do modelo é a mesma anterior, agora com coeficientes associados a cada categoria da variável resposta.

Normalmente espera-se que haja um tipo de tendência linear entre cada OR dos pontos de corte específicos e a variável-resposta.¹ Nesse caso, um conjunto de restrições (γ_{kl}) podem ser incluídas no modelo, para esclarecer essa linearidade (Tabela 1). Quando essas restrições são incorporadas, esse modelo é chamado de modelo de *odds* proporcionais parciais restrito.

Os parâmetros τ_j são escalares fixos que tomam a forma de restrições alocadas nos parâmetros. Nesse caso, para uma dada covariável X_m , γ_m não depende dos pontos de corte, mas é multiplicado por τ_j para cada j -ésimo logito.¹¹

Modelo de Razão Contínua (MRC)

Esse modelo permite comparar a probabilidade de uma resposta igual à categoria com determinado escore, digamos y_j , $Y=j$, com a probabilidade de uma resposta maior, $Y > y_j$, como indicado na Tabela 1.

Esse modelo possui diferentes interceptos e coeficientes para cada comparação e pode ser ajustado por k modelos de regressão logística binária.¹¹ É mais apropriado quando há um interesse intrínseco em uma categoria específica da variável-resposta.¹

Modelo Estereótipo (ME)

O ME pode ser considerado uma extensão do modelo de regressão multinomial¹⁰ e compara cada categoria

da variável-resposta com uma categoria de referência, normalmente a primeira categoria ou a última. Entretanto, devido ao caráter ordinal dos dados é imposta uma estrutura linear aos β_{jl} ($j=1, \dots, k$ e $l=1, \dots, p$), ou seja, são atribuídos pesos (ω_j) aos coeficientes.¹¹

Os pesos (ω_j) do modelo são diretamente relacionados com o efeito das covariáveis. Por isso, OR formada terá uma tendência de crescimento, já que os pesos normalmente são construídos com ordenação

$$(0 = \omega_1 \leq \omega_2 \leq \dots \omega_k).$$

Esse modelo deve ser utilizado quando a variável-resposta é uma variável ordinal com categorias discretas.

Em todos os modelos ordinais mencionados a significância dos coeficientes deve ser testada por meio do teste de Wald.¹⁰ Para o exercício apresentado, ele foi calculado utilizando a aproximação pela distribuição normal padronizada.

VERIFICAÇÃO DA QUALIDADE DO AJUSTE DOS MODELOS ORDINAIS

Assim como em qualquer tipo de análise de regressão, é importante avaliar a qualidade do ajuste dos modelos de regressão logística ordinal, pois a falta de ajuste pode, por exemplo, levar a viés de estimação de efeitos. A avaliação do ajuste pode detectar: covariáveis importantes; interações omitidas; casos em que a função de ligação (logito) não foi apropriada; casos em que a forma funcional da modelagem das covariáveis não está correta; e, finalmente, casos em que a suposição de *odds* proporcional foi violada.⁴

Embora muitos métodos tenham sido desenvolvidos para avaliar o ajuste de modelos de regressão logística binária, poucos desses métodos foram estendidos para dados de resposta ordinal.¹⁰ Normalmente, a qualidade do ajuste dos modelos ordinais é verificada usando os testes de Pearson ou *deviance*. Esses testes envolvem a criação de uma tabela de contingência na qual as linhas consistem de todas as possíveis configurações das covariáveis do modelo e as colunas são as categorias da resposta ordinal.¹⁴

As contagens esperadas dessa tabela são expressas por

$E_{lj} = \sum_{i=1}^{N_l} \hat{p}_{lj}$, onde N_l é o número total de indivíduos classificados na linha l e \hat{p}_{lj} representa a probabilidade de um indivíduo na linha l ter a resposta j calculada a partir do modelo adotado.¹⁴ O teste de Pearson para avaliar a adequação do ajuste compara essas contagens esperadas com as observadas pela fórmula:

$$\chi^2 = \sum_{l=1}^L \sum_{j=1}^k \frac{(O_{lj} - E_{lj})^2}{E_{lj}} \quad (1)$$

A estatística *deviance* também compara contagens observadas e esperadas, mas pela fórmula:

$$D^2 = 2 \sum_{i=1}^L \sum_{j=1}^k O_{ij} \log \frac{O_{ij}}{E_{ij}} \quad (2)$$

Os testes para avaliar a qualidade do ajuste do modelo são baseadas nas aproximações das estatísticas (1) e (2) para a distribuição qui-quadrado com $(L-1)(k-1)p$ graus de liberdade (L é o número de linhas, k o número de colunas da tabela de contingência e p o número de covariáveis do modelo). Um valor-p significativo leva à conclusão de falta de ajuste do modelo aos dados estudados.¹⁴

Pulkstenis & Robinson¹⁴ (2004) relatam que as estatísticas (1) e (2) não fornecem uma boa aproximação da distribuição qui-quadrado quando são ajustadas covariáveis contínuas e propõem pequenas modificações neste caso.

No presente estudo, em todos os modelos considerados, foram utilizados os testes de qualidade do ajuste de Pearson ou *deviance*, uma vez que se encontram implementados nos pacotes estatísticos usuais.

Como a literatura sobre ferramentas de diagnóstico ou avaliação de ajuste para modelos ordinais é relativamente escassa, Hosmer & Lemeshow¹⁰ (2000) sugerem a utilização de regressões binárias separadas para cada ponto de corte para que sejam criadas as estatísticas de diagnóstico para os modelos ordinais. Normalmente são construídos gráficos dos resíduos para os modelos de *odds* proporcionais usando o ajuste desses modelos para prever uma série de eventos binários $Y > j, j=1, 2, \dots, k$. Assim, para a variável indicadora $[Y \geq j]$, o resíduo escore, para o caso i e covariável p é dado por:⁹

$$U_{ip} = X_{ip} ([Y_i \geq j] - \hat{P}_{ij}),$$

$$\hat{P}_{ij} = \frac{1}{1 + \exp[-(\hat{\alpha}_j + \underline{x}\hat{\beta})]} \quad (3)$$

Nos gráficos do resíduo escore são colocados no eixo vertical a média $\bar{U}_{.p}$ e os respectivos intervalos de confiança, com as categorias da variável resposta no eixo horizontal. Se a suposição de *odds* proporcionais for válida, para cada covariável, os intervalos de confiança para cada categoria da variável-resposta devem ter uma aparência semelhante.⁹

Existem ainda os resíduos parciais que são muito usados para verificar se todas as covariáveis do modelo se comportam de forma linear. No contexto de regressão ordinal, é necessário calcular modelos de regressão logística binária para todos os pontos de corte da variável resposta Y , sendo o resíduo parcial para cada caso i e covariável p definido da seguinte forma:⁹

$$r_{ip} = \hat{\beta}_p X_{ip} + \frac{[Y_i \geq j] - \hat{P}_{ij}}{\hat{P}_{ij}(1 - \hat{P}_{ij})}$$

$$\hat{P}_{ij} = \frac{1}{1 + \exp[-(\hat{\alpha}_j + \underline{x}\hat{\beta})]} \quad (4)$$

Os gráficos de resíduos parciais geram estimativas de como cada covariável (x) se relaciona com cada categoria da variável resposta (Y).⁹

Assim, os resíduos parciais são usados para verificar a necessidade de transformações nas covariáveis (linearidade) ou mesmo a validade da suposição de *odds* proporcionais (paralelismo das curvas).⁹

COMANDOS PARA EXECUÇÃO DOS MODELOS

Nessa seção serão mostrados os passos para ajuste dos modelos mostrados na Seção 2 nos *softwares* R ou Stata que foram sumarizados no Quadro 1.

Os comandos foram ilustrados com os dados do NHANES II, e as variáveis foram assim nomeadas:

- Variável-resposta: *saúde*
- Covariáveis: *idade, diabetes, cor_pele*

Ajuste dos Modelos no Software R

A) Modelo de Odds Proporcionais

No *software* R, o MOP pode ser ajustado por meio do comando *lrm*, desenvolvido por Harrell e integrante do pacote *Design* (Tabela 1). Esse comando ajusta modelos binários e ordinal de *odds* proporcionais usando o método de máxima verossimilhança ou alternativamente, máxima verossimilhança penalizada.⁹

Os argumentos utilizados são: fórmula, ou seja, os termos a serem incluídos no modelo (variável resposta e covariáveis), nome do arquivo de dados a ser utilizado, entre outros. As saídas mostradas após execução dos comandos são: expressão utilizada, tabela de frequências para a resposta, vetor com algumas estatísticas importantes, estimativas dos coeficientes, vetor das primeiras derivadas do log da função de verossimilhança e *deviance* do modelo.

B) Modelo de Razão Contínua

O MRC pode ser implementado no *software* R por meio de uma reestruturação dos dados que é feita com o comando *cr.setup*, integrante do pacote *Design*. Esse comando possibilita criar várias novas variáveis a partir da variável resposta y , que serão usadas no ajuste do modelo de razão contínua.¹

Quatro novas variáveis são adicionadas por esse comando:

1. *y* – nova variável binária que será usada como resposta no ajuste do modelo de regressão logística binária;
2. *cohort* – vetor indicando qual ponto de corte (das comparações do MRC) foi aplicado;
3. *subs* – vetor utilizado para replicar as outras variáveis (explicativas) da mesma forma que *y* foi replicado;
4. *reps* – variável que especifica quantas vezes cada observação original foi replicada.

O modelo é obtido ajustando uma regressão logística binária nos dados reestruturados com uma nova resposta dicotômica (*y*) como variável dependente, incluindo a covariável criada (*cohort*), que indica o nível do ponto de corte, e reestruturando as covariáveis pelo vetor (*subs*), como mostrado na Tabela 1.¹

A suposição de heterogeneidade dos pontos de corte pode ser testada incluindo no modelo um termo de interação entre a exposição de interesse e a variável indicadora do ponto de corte (*cohort*), o que é denominado de modelo saturado, e comparando o valor do log da função de verossimilhança dos modelos com e sem o termo de interação.

C) Modelo estereótipo

O ME pode ser ajustado por meio de modelos lineares generalizados que possuem matrizes de restrições estimadas. Os pesos (restrições) são estimados como parâmetros adicionais do modelo, utilizando a família multinomial para o ajuste.

No *software* R, o ME pode ser ajustado por meio do comando *rrvglm* integrante do pacote VGAM.¹⁷

Análise de Resíduos

A função *residuals.lrm* do pacote *Design* é usada para construir gráficos dos resíduos após ajuste do MOP no *software* R, como apresentado na Tabela 1.

No gráfico do resíduo score (*score.binary*), se a suposição de *odds* proporcionais for válida é esperado que, para cada covariável, a tendência em torno das categorias da variável resposta tenha um comportamento horizontal constante. Já no gráfico do resíduo parcial (*partial*), para um modelo bem ajustado espera-se que as curvas tenham um caráter linear e sejam paralelas.⁹

Software Stata

Ajuste do Modelo de Odds Proporcionais Parciais

Até o momento, o MOPP não se encontra disponível no *software* R, mas pôde ser ajustado no Stata 9.0 utilizando o comando *gologit2* desenvolvido por Williams¹⁶

(2006). Esse comando possibilitou testar a suposição de *odds* proporcionais por meio da opção *autofit* e ajustar coeficientes para as várias categorias das variáveis em que essa suposição é violada.

EXEMPLO DE APLICAÇÃO – NHANES II

Considerou-se como variável dependente a condição de saúde, classificada em cinco categorias (1=ruim, 2=regular, 3=média, 4=boa e 5=excelente). Didaticamente, foi construído um modelo considerando três variáveis explicativas: uma variável quantitativa – idade (em anos), uma variável categórica binária – diabetes (não; sim) e uma variável com mais de duas categorias – cor da pele (branca; não-branca; outras). Para a variável cor da pele foram criadas variáveis indicadoras, considerando como referência a cor não branca.

A Tabela 2 foi construída como forma didática de mostrar como pode ser obtido OR, considerando uma categoria como referência ou agrupando as categorias. No primeiro cálculo, a condição de saúde excelente foi considerada como referência e cada uma das subsequentes categorias foi comparada a ela separadamente, como é feito no ME. Observou-se que o valor de OR aumenta à medida que a condição de saúde piora.

No segundo cálculo, OR foi calculada de acordo com a equação do MOP, na qual se comparam valores menores ou iguais a uma dada categoria a valores maiores (Tabela 1). Quando comparada à condição de saúde excelente, OR das condições boa a ruim (OR=6,3) foi igual àquela quando comparamos condição de saúde excelente e boa às condições razoável, média e ruim. Este caso exemplifica o modelo de *odds* proporcional, ou seja, de OR semelhantes para todas as categorias comparadas, suposição principal do MOP.

Os resultados do MOP estão apresentados na Tabela 3. O teste score sugere violação da suposição de *odds* proporcional para as variáveis cor da pele e idade, isoladamente, e também para o modelo múltiplo. Além disso, o teste de *deviance* indicou falta de ajuste do modelo.

Os gráficos de resíduos (score e parcial) para avaliar a adequação do MOP estão apresentados nas Figuras 1 e 2, respectivamente. Na Figura 1 (resíduo score), os resultados reforçam a conclusão do teste score, pois as curvas para a variável diabetes mostraram um formato horizontal próximo ao zero. Entretanto, a variável idade apresentou uma oscilação no comportamento para as categorias de condição de saúde boa e média estando bem abaixo da linha do resíduo zero. O mesmo ocorreu para a covariável cor da pele. Porém, nesse caso, a oscilação maior foi observada nas categorias de condição de saúde razoável e ruim para a classificação outras e na categoria média para a cor branca.

Tabela 2. Associação entre diabetes e condições de saúde de acordo, com estimação da OR e intervalo de confiança 95%

Variável	Condição de saúde				
	Excelente	Boa	Média	Razoável	Ruim
Diabetes					
Não	2383 (24%)	2546 (26%)	2805 (29%)	1508 (15%)	594 (6%)
Sim	24 (5%)	45 (9%)	133 (27%)	162 (32%)	135 (27%)
OR [IC 95%]*	1,0	1,8 [1,1;2,9]	4,7 [3,0;7,3]	10,7 [6,9;16,5]	22,6 [14,5;35,2]
Diabetes					
Não	2383	2546	2805	1508	594
Sim	24	45	133	162	135
OR [IC 95%**	1,0		6,3 [4,2;9,6]		
Diabetes					
Não	2383	2546	2805	1508	594
Sim	24	45	133	162	135
OR [IC 95%**		1,0		6,3 [4,8;8,1]	
Diabetes					
Não	2383	2546	2805	1508	594
Sim	24	45	133	162	135
OR [IC 95%**		1,0		5,4 [4,7;7,2]	
Diabetes					
Não	2383	2546	2805	1508	594
Sim	24	45	133	162	135
OR [IC 95%**			1,0		5,8 [4,7;7,2]

* OR considerando uma categoria como referência (onde Y=condição de saúde, Y₀=categoria excelente, x(a)= sem diabetes,

$$OR_j = \frac{P(Y = Y_0 | x^{(A)})}{P(Y > Y_0 | x^{(A)})} \bigg/ \frac{P(Y = Y_0 | x^{(B)})}{P(Y > Y_0 | x^{(B)})}$$

x(b)=com diabetes) ** OR para dados ordinais - (onde Y =condição de saúde, j=cada categoria de Y, x(a)= sem diabetes, x(b)=com diabetes)

$$OR_j = \frac{P(Y \leq Y_j | x^{(A)})}{P(Y > Y_j | x^{(A)})} \bigg/ \frac{P(Y \leq Y_j | x^{(B)})}{P(Y > Y_j | x^{(B)})}$$

Nos gráficos do resíduo parcial (Figura 2), a suposição de regressão paralela pareceu bastante aceitável devido a seu aspecto linear e às retas aproximadamente paralelas para a variável diabetes. Já no gráfico para a categoria outras da variável cor da pele, apesar do comportamento linear, observou-se cruzamento das curvas, violando, portanto, a suposição do paralelismo. A covariável idade não apresentou comportamento linear, o que poderia contribuir para a falta de ajuste do modelo. Mesmo incluindo termos de graus mais elevados para a idade, o teste de *deviance* continuou mostrando um ajuste ruim.

O ajustamento do MOPP está apresentado na Tabela 4. Nesse modelo, os efeitos foram significativos para as quatro comparações, e os coeficientes não variaram para a variável diabetes, indicando que um indivíduo com pior condição de saúde tem 3,39 vezes mais chance de ser diabético, quando comparado a um indivíduo com melhores condições de saúde.

Comparada à cor da pele “não-branco”, não houve variação dos coeficientes nas diversas comparações

para a cor de pele “branca”, enquanto para a cor de pele “outra” houve variação. À medida que se avaliou piores condições de saúde, houve aumento no efeito de proteção, ou seja, diminuição do valor absoluto de OR. Para a covariável idade observou-se também mudança de magnitude do OR, nas diversas categorias de comparações das condições de saúde. Com o aumento de um ano na idade, a chance de apresentar uma condição de saúde de boa a ruim foi 1,03 vezes maior se comparada à condição de saúde excelente. Essa chance pode chegar a 1,05 quando se compara condição de saúde ruim com saúde de excelente à razoável.

Quanto ao ME, em todas as comparações, o efeito das covariáveis foi significativo (valor-p<0,01), e o teste *deviance* indicou um bom ajuste do modelo (Tabela 5), ou seja, pessoas com saúde ruim têm 10 vezes mais chance de ser portadoras de diabetes do que aquelas com saúde excelente. A magnitude desta associação diminuiu à medida que a condição de saúde se aproxima da excelente, chegando a 1,43 na comparação de saúde boa a excelente.

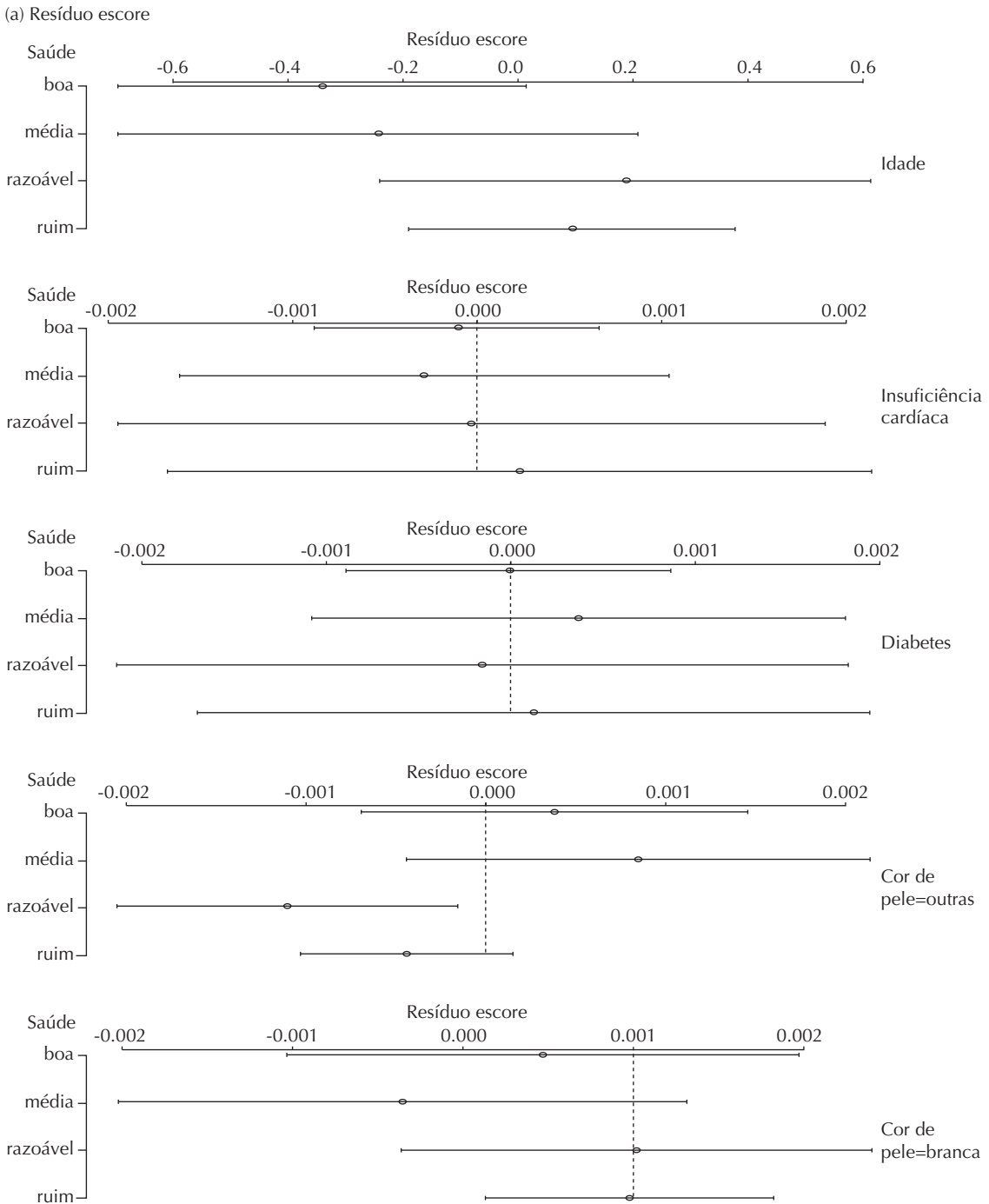


Figura 1. Gráficos de resíduos (escore e parcial) para as covariáveis incluídas no modelo tendo como resposta a condição de saúde.

COMENTÁRIOS

No exemplo apresentado no presente artigo, o MOP não apresentou bom ajuste e os gráficos de resíduos mostraram retas não paralelas para algumas covariáveis, indicando violação da premissa principal. Considerando que inferências baseadas nesse modelo podem não ser corretas foi alternativamente apresentado o

MOPP, com estimativa de OR para cada uma das comparações. Entretanto, o ME foi o que melhor se ajustou aos dados analisados, de acordo com os resultados do teste *deviance*.

De maneira geral, os modelos de regressão logística ordinal são recomendados para análise de dados ordinais.^{1,3,8,10,11} Ananth & Kleinbaum¹ (1997) relataram que

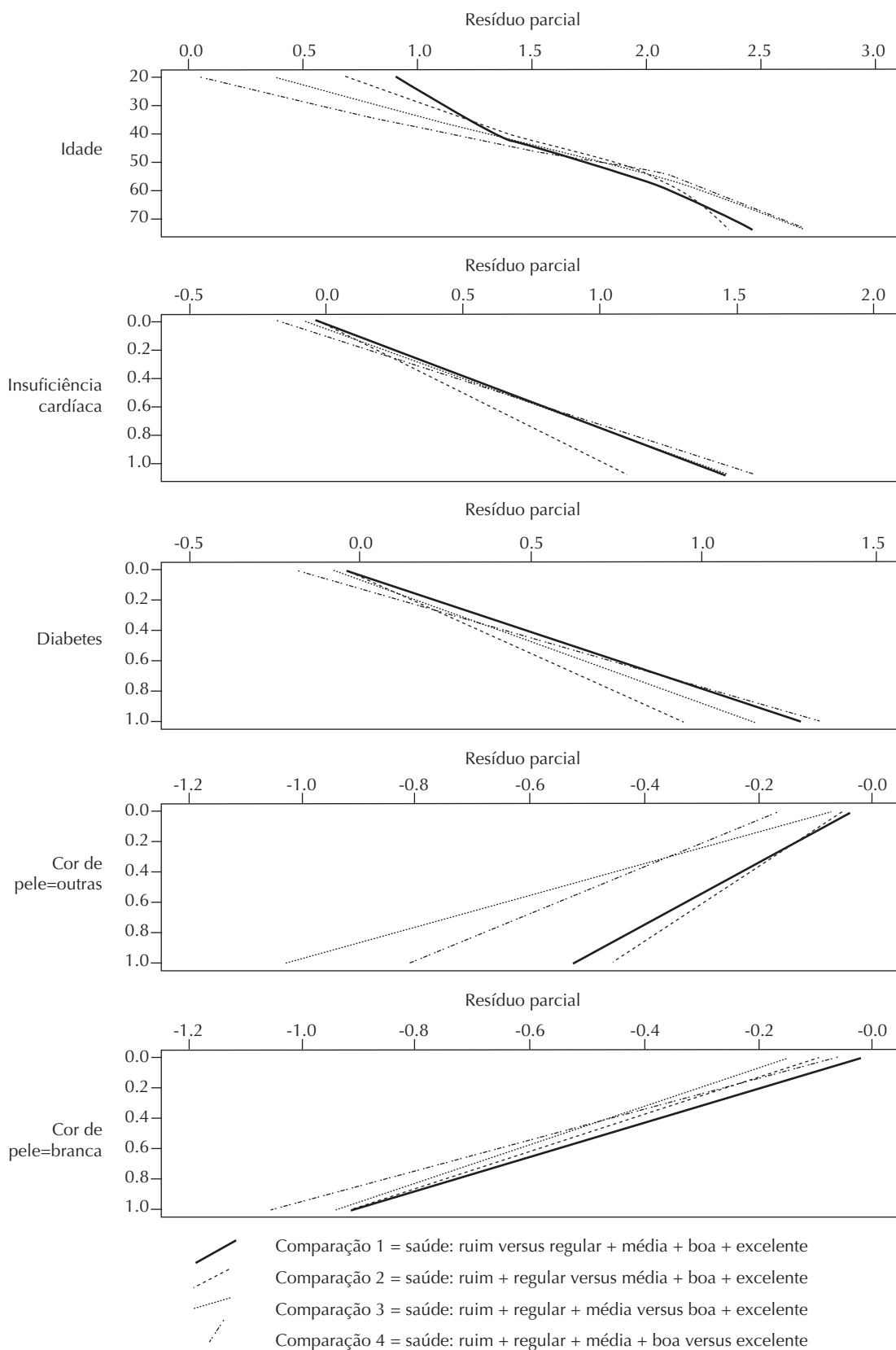


Figura 2. Gráficos de resíduos (escore e parcial) para as covariáveis incluídas no modelo tendo como resposta a condição de saúde.

Tabela 3. Resultados* do modelo de *odds* proporcional** final tendo como resposta a condição de saúde (NHANES II).

Covariável	β	EP(β)	OR	Valor-p Teste escore***
Diabetes				
Não			1,00	0,60
Sim	1,26	0,09	3,35	
Cor de pele				
Não branca			1,00	0,01
Branca	-0,86	0,14	0,42	
Outras	-0,64	0,06	0,53	
Idade (em anos)	0,04	0,01	1,04	<0,01

* Em todos os modelos, todas as variáveis foram significativas ao nível de significância de 1%

** Teste escore modelo final (valor-p<0,001)

*** Teste de deviance (valor-p<0,001)

o MOP e o MRC são os mais utilizados em aplicações epidemiológicas e biomédicas em relação aos MCP e ME. Porém, esses modelos levam a fortes suposições que, se não forem válidas, podem gerar interpretações incorretas, como ocorreu no exemplo utilizado.¹

Outros autores¹¹ afirmam que o tipo de modelo a ser utilizado depende do caráter da variável resposta ordinal, ou seja, se essa variável foi ordenada a partir de uma variável contínua reagrupada ou se originou de uma variável discreta. No primeiro caso, o MOP é o mais indicado quando a premissa de retas paralelas não é violada. No segundo caso, o ME é o mais indicado, como no presente estudo, em que a variável resposta, definida como condição de saúde do estudo NHANES II, foi tratada com categorias discretas.

Mesmo na presença de uma resposta ordinal, outras opções de análises multivariadas devem ser consideradas. Um caminho é utilizar a metodologia de árvores de decisão,⁵ um método mais descritivo e que considera um número maior de variáveis no modelo final. Também podem ser utilizadas outras funções de ligação do

modelo, como probito e complementar log-log. Entretanto, a regressão ordinal é uma técnica paramétrica que, impondo uma estrutura rígida do modelo, além de ser mais conservadora e parcimoniosa, permite a quantificação de intervalos de confiança para os parâmetros, e portanto, apresenta facilidade de interpretação de OR. Abordagens alternativas merecem análise e discussão, mas que não serão tratadas neste artigo.

Na construção dos modelos ordinais, Hosmer & Lemeshow¹⁰ (2000) propõem estratégias como as adotadas no exemplo do presente artigo. Esses autores recomendam inicialmente fazer uma análise univariada para seleção dos efeitos principais e incluir no modelo apenas as variáveis significativas com um nível de significância pré-fixado. Em seguida, ajustar o modelo, verificar sua adequação por meio dos testes adequados e gráficos de resíduos e, por fim, interpretar o modelo por meio da estimativa de OR.

Entretanto, os métodos para verificação do ajuste dos modelos ordinais são escassos. Até o momento encontramos na literatura disponível nenhuma técnica

Tabela 4. Resultados* do modelo de *odds* proporcionais parciais final tendo como resposta a condição de saúde (NHANES II).

Covariável	Comparações							
	Excelente versus (boa+média+razoável+ruim)		(excelente+boa) versus (média+razoável+ruim)		(excelente+boa+média) versus (razoável+ruim)		(excelente+boa+média+razoável) versus (ruim)	
	$\hat{\beta}_1$	\hat{OR}_1	$\hat{\beta}_2$	\hat{OR}_2	$\hat{\beta}_3$	\hat{OR}_3	$\hat{\beta}_4$	\hat{OR}_4
Diabetes								
Não		1,00		1,00		1,00		1,00
Sim	1,22	3,39	1,22	3,39	1,22	3,39	1,22	3,39
Cor de pele								
Não branca		1,00		1,00		1,00		
Branca	-0,88	0,42	-0,88	0,42	-0,88	0,42	-0,88	0,42
Outras	-0,54	0,58	-0,44	0,64	-1,04	0,35	-1,14	0,32
Idade (em anos)	0,03	1,03	0,04	1,04	0,05	1,05	0,05	1,05

* Em todos os modelos, todas as variáveis foram significativas ao nível de significância de 1%

Tabela 5. Resultados* do modelo estereótipo** final tendo como resposta a condição de saúde (NHANES II).

Covariável	Comparações da condição de saúde							
	Ruim versus excelente		Regular versus excelente		Média versus excelente		Boa versus excelente	
	$\hat{\beta}_1$	\hat{OR}_1	$\hat{\beta}_2$	\hat{OR}_2	$\hat{\beta}_3$	\hat{OR}_3	$\hat{\beta}_4$	\hat{OR}_4
Diabetes								
Não		1,00		1,00		1,00		1,00
Sim	2,25	9,45	1,69	5,43	0,97	2,64	0,36	1,43
Cor de pele								
Não branca		1,00		1,00		1,00		1,00
Branca	-1,64	0,19	-1,23	0,29	-0,71	0,49	-0,26	0,77
Outras	-1,30	0,27	-0,98	0,38	-0,56	0,57	-0,21	0,81
Idade (em anos)	0,08	1,08	0,06	1,06	0,03	1,03	0,01	1,01

* Em todos os modelos, todas as variáveis foram significativas ao nível de significância de 1%

** Teste de deviance (valor-p = 0,100)

de verificação do ajuste do ME. As estatísticas de diagnóstico existentes, propostas por Harrel⁹ (2002) e aplicáveis ao MOP, são apenas gráficos traçados a partir de regressões binárias separadas para os pontos de corte da variável ordinal. Embora incompletas, até o momento, essas técnicas são de grande importância para se ter uma indicação da qualidade de ajuste dos modelos ordinais. A análise dos resíduos parciais, mesmo que gráfica, é considerada muito útil para modelos ordinais, pois eles simultaneamente checam linearidade, indicando eventuais transformações que devem ser utilizadas, bem como o pressuposto de *odds* proporcionais.

Por outro lado, deve-se ter cuidado na interpretação dos resíduos, principalmente considerando novamente a escassez de informações de como fazê-la. Os gráficos, por vezes, podem sugerir informações confusas e dificultar a tomada de decisões quanto à violação ou não da suposição de *odds* proporcionais. Uma alternativa

é utilizar a análise dos resíduos conjuntamente com o teste do escore, pois quando houver alguma dúvida quanto ao formato do gráfico, esse teste pode contribuir para a conclusão final.

Os modelos de regressão logística ordinal têm-se mostrado apropriados para análise de dados com resposta ordinal. A escolha do melhor modelo depende do caráter da variável ordinal, adequação do modelo às suposições, qualidade do ajuste e capacidade de boa explicação com reduzido número de parâmetros a serem estimados.

Finalmente, uma boa implementação computacional e o domínio dos comandos para execução dos modelos ordinais são essenciais, não só para a escolha do modelo mais adequado, mas também para comparações entre modelos. Para isso, o programa R torna-se uma ferramenta importante, trazendo vários modelos e alguns gráficos de diagnóstico.⁹

REFERÊNCIAS

1. Ananth CV, Kleinbaum DG. Regression models for ordinal responses: a review of methods and applications. *Int J Epidemiol.* 1997;26(6):1323-33. DOI: 10.1093/ije/26.6.1323
2. Anderson JA. Regression and ordered categorical variables. *J R Statist Soc B.* 1984;46(1):1-30.
3. Bender R, Benner A. Calculating ordinal regression models in SAS and S-Plus. *Biometrical J.* 2000;42(6):677-99. DOI: 10.1002/1521-4036(200010)42:6<677::AID-BIMJ677>3.0.CO;2-O
4. Brant R. Assessing proportionality in the proportional odds model for ordinal logistic regression. *Biometrics.* 1990;46(4):1171-8. DOI: 10.2307/2532457
5. Breiman L, Friedman J, Stone CJ, Olshen RA. Classification and regression trees. New York: Chapman & Hall; 1984.
6. Colin RB. Bioestatística usando R: apostila para biólogos. Bragança; 2004.
7. Fienberg SE. Fixed Margins and Logit Models. In: Fienberg SE. The analysis of cross-classified categorical data. Cambridge, MA: MIT Press; 1980. p.110-6.
8. Greenland S. Alternative models for ordinal logistic regression. *Stat Med.* 1994;13(16):1665-77. DOI: 10.1002/sim.4780131607
9. Harrell Jr FE. Regression modelling strategies: with applications to linear models, logistic regression, and survival analysis. New York: Springer; 2002.
10. Hosmer DW, Lemeshow S. Applied logistic regression. 2. ed. New York: John Wiley & Sons; 2000.
11. Lall R, Campbell MJ, Walters SJ, Morgan K. A review of ordinal regression models applied on health-related quality of life assessments. *Stat Methods Med Res.* 2002;11(1):49-67. DOI: 10.1191/0962280202sm271ra
12. McCullagh P. Regression models for ordinal data. *J R Statist Soc B.* 1980;42(2):109-42.
13. Peterson BL, Hanrel FE. Partial proportional odds models for ordinal response variables. *Appl Statistic.* 1990;39(2):205-17. DOI: 10.2307/2347760
14. Pulkstenis E, Robinson TJ. Goodness-of-fit tests for ordinal response regression models. *Stat Med.* 2004;23(6):999-1014. DOI: 10.1002/sim.1659
15. Walker SH, Duncan DB. Estimation of the probability of an event as a function of several independent variables. *Biometrika.* 1967;54(1):167-79.
16. Williams R. Generalized ordered logit/partial proportional odds models for ordinal dependent variables. *Stata J.* 2006;6(1):58-82.
17. Yee TW, Hastie TJ. Reduced-rank vector generalized linear models. *Statist Model.* 2003;3(1):15-41. DOI: 10.1191/1471082X03st045oa

Caiaffa WT é apoiada pelo Conselho Nacional de Desenvolvimento Científico e Tecnológico (CNPq; bolsa produtividade em pesquisa); Abreu MNS foi apoiada pela Coordenação de Aperfeiçoamento de Pessoal de Nível Superior (Capes; bolsa de mestrado).