

**UNIVERSIDADE FEDERAL DO RIO GRANDE DO SUL
FACULDADE DE MEDICINA
PROGRAMA DE PÓS-GRADUAÇÃO EM EPIDEMIOLOGIA**



DISSERTAÇÃO DE MESTRADO

Métodos para estimar prevalências ajustadas

Natália Bordin Barbieri

Orientador: Prof. Dr. Álvaro Vigo

Porto Alegre, fevereiro de 2016

UNIVERSIDADE FEDERAL DO RIO GRANDE DO SUL
FACULDADE DE MEDICINA
PROGRAMA DE PÓS-GRADUAÇÃO EM EPIDEMIOLOGIA



DISSERTAÇÃO DE MESTRADO

Métodos para estimar prevalências ajustadas

Natália Bordin Barbieri

Orientador: Prof. Dr. Álvaro Vigo

A apresentação desta dissertação é exigência do Programa de Pós-graduação em Epidemiologia, Universidade Federal do Rio Grande do Sul, para obtenção do título de Mestre.

Porto Alegre, Brasil.
2016

BANCA EXAMINADORA

Prof.^a Dra. Vivian Cristine Luft

Faculdade de Medicina

PPG em Epidemiologia

Universidade Federal do Rio Grande do Sul

Prof.^a Dra. Suzi Alves Camey

Instituto de Matemática e Estatística

PPG em Epidemiologia

Universidade Federal do Rio Grande do Sul

Prof. Dr. Cleber Bisognin

Instituto de Matemática e Estatística

Departamento de Estatística

Universidade Federal do Rio Grande do Sul

AGRADECIMENTOS

Agradeço aos meus pais, Jeferson e Marfisa, e aos meus irmãos, Renata e Tomás, que sempre me apoiaram, incentivaram e oportunizaram para que eu pudesse chegar até aqui. Eles são a base de tudo.

Ao Lázaro Ribeiro Luz pelo companheirismo e amor em nossas jornadas.

Ao Professor Álvaro Vigo, meu orientador, pelo exemplo de pessoa e profissional. Se um dia eu conseguir chegar perto de onde ele chegou, a caminhada valeu a pena.

Aos Professores Maria Inês Schmidt e Bruce Duncan, pelas oportunidades ao longo do tempo. É uma honra poder trabalhar com vocês.

Ao Professor Loyd Chambless pela ajuda e discussão sobre os métodos utilizados.

Aos professores Cleber Bisognin, Suzi Camey e Vivian Luft, pela oportunidade de tê-los como banca neste trabalho.

A Paula Sientchkovski pelo companheirismo desde os tempos da graduação.

E a querida grande Equipe de Estatística ELSA-Brasil, que oportuniza a cada dia o crescimento pessoal e profissional, e que fazem a palavra "equipe" valer muito a pena.

SUMÁRIO

ABREVIATURAS E SIGLAS	6
RESUMO	7
1. APRESENTAÇÃO	9
2. INTRODUÇÃO	10
3. REVISÃO DA LITERATURA	11
3.1 Estimação de prevalências ajustadas	11
3.1.1 Método de Predição Condicional.....	14
3.1.2 Método de Predição Marginal.....	15
3.1.3 Estimação por intervalo utilizando método Delta.....	16
3.2 Aspectos computacionais.....	19
3.2.1 SAS - Macro %ADJ_PROP	19
3.2.2 SUDAAN (SAS-Callable): CONDMARG e PREDMARG.....	19
3.2.3 Stata - Função <i>margins</i>	20
3.2.4 R – Função <i>margins</i>	20
4. OBJETIVOS	21
5. REFERÊNCIAS BIBLIOGRÁFICAS.....	22
6. ARTIGO ORIGINAL	25
7. CONCLUSÕES E CONSIDERAÇÕES FINAIS	61

ABREVIATURAS E SIGLAS

ARIC	The Atherosclerosis Risk in Communities Study
ELSA-Brasil	Estudo Longitudinal da Saúde do Adulto
ICHD-3	International Classification of Headache Disorders

RESUMO

Objetivo: Apresentar e discutir métodos para estimar prevalências ajustadas em pesquisas clínicas e epidemiológicas, bem como desenvolver rotinas computacionais em SAS e R.

Métodos: No contexto de estudo transversal, foi simulada uma amostra de 2.000 observações independentes, considerando o desfecho dicotômico diabetes, sexo como a variável de exposição e idade como variável de ajuste. As estimativas de prevalências ajustadas (IC 95%) foram estimadas pelos métodos de predição condicional e marginal, utilizando as rotinas desenvolvidas em SAS e R. O método Delta foi usado para construir os intervalos de confiança. Os resultados foram comparados com aqueles do SUDAAN (SAS-Callable), Stata e a macro %ADJ_PROP (SAS).

Resultados: No exemplo simulado, 68,2% são do sexo feminino e a idade média (DP) foi 57,6 (5,0) anos, sendo 54,2 (3,9) anos em homens e 59,2 (4,6) anos em mulheres. A estimativa da prevalência global do desfecho foi de 25,3% (IC 95%:23,4-27,3); sendo 13,8% (IC 95%:11,7-16,7) e 30,7% (IC 95%:28,3-33,2), respectivamente para homens e mulheres. As estimativas de prevalências ajustadas por idade, por meio do método de predição condicional, foram de 19,6% (IC 95%:16,2-23,6) para homens, e 23,6% (IC 95%:21,2-26,1) para mulheres. Pelo método de predição marginal, as estimativas foram de 22,4% (IC 95%:18,7-26,5) para homens, e 26,3% (IC 95%:24,1-28,6) para mulheres.

Conclusão: A discrepância entre as estimativas não ajustadas é devida ao confundimento pela idade. Estimativas livres de confundimento podem ser obtidas por meio das prevalências ajustadas pela idade. No entanto, a estimativa pelo método de predição condicional não engloba a prevalência global. Em virtude disso, o método de predição marginal é, geralmente, mais adequado. A rotina desenvolvida na versão para R é uma alternativa aos softwares comerciais.

Palavras-chave: Prevalências ajustadas, Método de predição condicional, Método de predição marginal, Método Delta.

ABSTRACT

Objective: To present and discuss methods to estimate adjusted prevalences for clinical and epidemiological research, and develop computational routines in SAS and R.

Methods: In the context of cross-sectional study, it was simulated a sample of 2,000 independent observations, considering the dichotomous outcome diabetes, sex as the exposure variable and age as an adjustment variable. Adjusted prevalences were estimated by the conditional and marginal methods, using routines developed in SAS and R. Confidence intervals were constructed using the Delta method. The results were compared with those of the SUDAAN (SAS-callable), Stata and macro %ADJ_PROP (SAS).

Results: In simulated example, 68.2% are female and the mean (SD) age was 57.6 (5.00) years old, being that 54.2 (3.94) years for men and 59.2 (4.60) years in women. The estimated global prevalence of outcome was 25.3% (CI 95%: 23.4-27.3) and 13.8% (CI 95%: 11.7-16.7) and 30.7% (CI 95%: 28.3-33.2), respectively for men and women. Estimates of adjusted prevalence for age, through the conditional method, were 19.6% (CI 95%: 16.2-23.6) for men, and 23.6% (CI 95%: 21.2-26.1) for women. For marginal method, the estimates were 22.4% (CI 95%: 18.7-26.5) for men and 26.3% (CI 95%: 24.1-28.6) for women.

Conclusion: The observed discrepancy in estimates by sex, unadjusted, can be attributed to confounding due to difference in age distribution between sexes. Comparable estimates (without confounding) of the prevalences can be obtained through prevalence adjusted for age. However, the estimate for the conditional method does not comprise the global prevalence. As a result, the marginal method is in general more suitable. The developed routines can be useful for estimating adjusted prevalences, particularly the R version (an alternative to commercial software).

Keywords: Adjusted Prevalence, Conditional prediction method, Marginal prediction method, Delta method.

1. APRESENTAÇÃO

Este trabalho consiste na dissertação de mestrado intitulada “**Métodos para estimar prevalências ajustadas**”, apresentada ao Programa de Pós-Graduação em Epidemiologia da Universidade Federal do Rio Grande do Sul, em 23 de fevereiro de 2016. O trabalho é apresentado em quatro partes, na ordem que segue:

1. Introdução, Revisão da Literatura e Objetivo.
2. Exemplo
3. Artigo
4. Conclusões e Considerações finais.

Documentos de apoio estão apresentados no apêndice.

2. INTRODUÇÃO

Estimativas de prevalências são reportadas com grande frequência em estudos clínicos ou epidemiológicos. Para que essas estimativas sejam comparáveis em diferentes populações, é necessário levar em conta os desequilíbrios nas distribuições de variáveis que possam alterar a frequência do evento. Em outras palavras, na presença de confundimento, estimativas não ajustadas podem não ser adequadas, e podem dificultar a compreensão do fenômeno.

Reportar prevalências ajustadas possibilita ao leitor um julgamento da frequência do evento livre de possíveis confundimentos, uma vez que a magnitude da associação entre desfecho e o fator em estudo pode variar entre os diferentes níveis do fator. Estratificação e modelos de ajuste multivariável são ferramentas comuns utilizadas para controlar os efeitos de confundimento (Szklo e Nieto, 2007).

Em uma revisão de programas de análise de dados foi constatado que somente o SUDAAN disponibiliza, de forma direta, estimativas de prevalências ajustadas (e intervalos de confiança) por meio dos métodos marginal e condicional. O Stata fornece apenas estimativas pelo método condicional. Para os programas SAS e R, é necessário o uso de rotinas específicas (macros ou funções).

Essa dissertação apresenta uma revisão dos diferentes métodos para estimar prevalências ajustadas, bem como aspectos computacionais disponíveis. Também são apresentadas rotinas computacionais em SAS e R, desenvolvidas para esta finalidade. Um exemplo é utilizado para explorar e discutir os métodos estatísticos e computacionais.

3. REVISÃO DA LITERATURA

Nas seções a seguir foram detalhados os métodos de estimação de prevalências ajustadas e também as rotinas computacionais disponíveis.

3.1 Estimação de prevalências ajustadas

Em estudos clínicos ou epidemiológicos muitas vezes é necessário estimar prevalências de um evento/doença. Para que essas estimativas sejam comparáveis em diferentes populações, é necessário levar em conta os desequilíbrios das distribuições das variáveis/fatores importantes que possam alterar a frequência da doença. Reportar prevalências ajustadas possibilita ao leitor um julgamento da frequência do evento livre de possíveis confundimentos, uma vez que a magnitude da associação entre desfecho e o fator em estudo pode variar entre os diferentes níveis do fator. Estratificação e modelos de ajuste multivariável são ferramentas comuns utilizadas para controlar os efeitos de confundimento (Szklo e Nieto, 2007).

O uso de prevalências ajustadas é frequente na literatura. Como exemplo, An (2015) utilizou prevalências ajustadas para descrever o diagnóstico de diabetes e leitura de rótulos nutricionais entre adultos norte-americanos para suas escolhas alimentares. Para obter as prevalências ajustadas, foi utilizado o modelo de regressão logística, porém não foi especificado o método utilizado.

Shon *et al.* (2015) também utilizaram prevalências ajustadas para comparar e analisar a prevalência de infecção pelo vírus da hepatite C por região na República da Coreia durante os anos de 2005-2012. O método utilizado para a prevalência ajustada foi o método direto de padronização (*direct standardization method*), utilizando como referência para o ajuste de idade a população de 2010. As prevalências foram ajustadas para sexo, idade e região, e foi possível identificar regiões onde a doença tinha uma maior prevalência.

Lebedeva *et al.* (2015) estimaram a prevalência de transtornos de dor de cabeça primária, diagnosticados de acordo com a ICHD-3 beta (*International Classification of*

Headache Disorders), em três grupos sociais diferentes. Para obter as prevalências ajustadas, utilizaram o método direto de padronização.

Esses são alguns exemplos que ilustram o uso de prevalências ajustadas, tendo sido estimadas por diferentes métodos. Entre as abordagens utilizadas estão padronização direta (*direct adjustment*), padronização indireta (*indirect adjustment*) e modelos multivariáveis.

A padronização direta utiliza as taxas observadas na população em estudo (que pode ser afetada por instabilidade, ocasionada por tamanho amostral pequeno), enquanto que a padronização indireta utiliza taxas de uma população de referência, e que por isso tem sua aplicação limitada a essa comparação de referência. A vantagem da padronização direta é que as taxas estimadas podem ser comparadas entre grupos (por exemplo, comparar o grupo feminino *versus* masculino), enquanto a padronização indireta permite comparar cada grupo estimado apenas com sua taxa utilizada da população de referência - por exemplo, comparar a taxa do sexo masculino estimada com a taxa do sexo masculino da população de referência (Woodward, 2014).

Em estudos ecológicos, o método de padronização direta é preferível uma vez que a informação existente é sobre a população. Sempre que tivermos dados no nível individual, é preferível utilizar ajuste multivariável. (Szklo e Nieto, 2007).

Wilcosky e Chambless (1985) afirmam que as principais vantagens do método de ajuste direto são a simplicidade computacional e as poucas suposições estatísticas. Entretanto, com a maior disponibilidade e diversidade de recursos computacionais, esses aspectos atualmente não são mais relevantes, além de ter como desvantagem a necessidade de categorização de variáveis quantitativas. Os autores também discutiram os métodos de predição condicional ("*conditional prediction method*") e marginal ("*marginal prediction method*") para estimar prevalências ajustadas. Esses métodos utilizam o modelo de regressão logística multivariável e são mais convenientes para testes de interações e de diferenças entre grupos, permitem explorar a natureza da relação funcional entre as variáveis de controle e o desfecho, e não exigem a

categorização de variáveis quantitativas. Lane e Nelder (1982) também apresentaram a abordagem de ajuste multivariável em modelos lineares generalizados para estimar proporções ajustadas por meio dos métodos de predição marginal e condicional.

Os métodos propostos por Wilcosky e Chambless (1985) têm sido bastante utilizados. Para estimar a associação entre incidência da doença arterial coronariana e espessura da parede da artéria carótida, Chambless *et al.* (1997) apresentaram prevalências dos principais fatores de risco, ajustadas por idade, centro de investigação e raça. No estudo sobre lipoproteína fosfolipase A₂ associada e alta sensibilidade da proteína c-reativa, esse método foi utilizado por Nambi *et al.* (2009) para melhorar a estratificação de risco de acidente vascular isquêmico, tendo sido estimadas prevalências de características (como hipertensão, diabetes, e uso de aspirina, estatina, etc.) na linha de base, ajustadas por idade, raça e sexo.

Chor *et al.* (2015) utilizaram o método de predição marginal para estimar a prevalência de pressão arterial elevada, ajustado por idade e sexo, com dados do ELSA-Brasil.

Ohira *et al.* (2010) apresentaram prevalências de diferentes variáveis de linha de base do estudo *ARIC (The Atherosclerosis Risk in Communities Study)*, ajustadas por idade e raça, em participantes com e sem risco de tromboembolismo venoso. Não foi explicitado se foi utilizado o método de predição marginal ou condicional.

Para estudar fatores associados a patologias nas cordas vocais em professores, Souza *et al.* (2011), utilizaram o método de predição condicional para obter prevalências ajustadas para cada grupo de interesse.

Epstein *et al.* (2013a, 2013b) utilizaram o método de predição marginal para estimar a prevalência de uso de medicamentos antipsicóticos e anticonvulsantes e de analgésicos opióides durante a gestação, ao longo do tempo. As prevalências em cada ano, expressadas por mil gestações, foram ajustadas para diferentes variáveis maternas.

Os métodos de predição condicional e marginal são descritos nas seções 3.1.1 e

3.1.2, respectivamente.

3.1.1 Método de Predição Condicional

Para descrever estes métodos será considerado o contexto de um estudo transversal em que a variável dependente (Y) representa presença ($Y=1$) ou ausência ($Y=0$) de diabetes, a variável dicotômica (x_1) representa sexo (1=masculino, 0=feminino), e a variável x_2 representa idade (quantitativa). Naturalmente, pode ser generalizado para mais preditores. Dada uma amostra aleatória com n indivíduos, obtida por meio de um estudo transversal, deseja-se estimar a prevalência do desfecho para cada categoria da exposição, ajustada pela variável x_2 . O modelo logístico definido na equação [1] pode ser usado para essa finalidade, ou seja,

$$\log \frac{P[Y=1|\mathbf{x}]}{1-P[Y=1|\mathbf{x}]} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 \quad [1]$$

em que $\mathbf{x} = (x_1, x_2)^t$ representa o vetor de covariáveis do modelo. O método de estimação de máxima verossimilhança pode ser usado para estimar os parâmetros do modelo, obtendo-se o modelo estimado descrito a seguir

$$\log \frac{P[Y=1|\mathbf{x}]}{1-P[Y=1|\mathbf{x}]} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 \quad [2].$$

Detalhes sobre o método de estimação podem ser obtidos em Mcculagh e Nelder (1989). As probabilidades estimadas pelo modelo em [2] são obtidas por

$$\hat{p}_{ij}(x_{1i}, x_{2i}) = \hat{P}[Y=1 | x_{1i} = j, x_{2i}] = \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 \times j + \hat{\beta}_2 x_{2i})}} \quad [3]$$

em que $j = 0,1$ e $i = 1,2,\dots,n$.

Para obter estimativas de prevalências ajustadas pelo método de predição condicional (\hat{p}_{Ci}), a média amostral do preditor x_2 , denotada por \bar{x}_2 , é usada na equação [3] para cada categoria da variável x_1 , como descrito nas equações [4] e [5]. Detalhes do método podem ser encontrados no manual do programa SUDAAN (SUDAAN Language Manual, 2012).

$$\hat{p}_{C_0}(\mathbf{x}) = \hat{P}[Y = 1 | x_1 = 0, x_2 = \bar{x}_2] = \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 \times 0 + \hat{\beta}_2 \bar{x}_2)}} = \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_2 \bar{x}_2)}}, \quad [4]$$

e

$$\hat{p}_{C_1}(\mathbf{x}) = \hat{P}[Y = 1 | x_1 = 1, x_2 = \bar{x}_2] = \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 \times 1 + \hat{\beta}_2 \bar{x}_2)}} = \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 + \hat{\beta}_2 \bar{x}_2)}}. \quad [5]$$

3.1.2 Método de Predição Marginal

Adotando a notação da seção anterior, as prevalências ajustadas pelo método de predição marginal do desfecho para as categorias da exposição são estimadas pelas equações [6] e [7]:

$$\hat{p}_{M_0}(\mathbf{x}) = \hat{P}[Y = 1 | x_1 = 0, x_{2i}] = \frac{1}{n} \sum_{i=1}^n \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 \times 0 + \hat{\beta}_2 x_{2i})}} = \frac{1}{n} \sum_{i=1}^n \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_2 x_{2i})}}, \quad [6]$$

$$\hat{p}_{M_1}(\mathbf{x}) = \hat{P}[Y = 1 | x_1 = 1, x_{2i}] = \frac{1}{n} \sum_{i=1}^n \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 \times 1 + \hat{\beta}_2 x_{2i})}} = \frac{1}{n} \sum_{i=1}^n \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 + \hat{\beta}_2 x_{2i})}}. \quad [7]$$

As estimativas do intervalo de confiança para as prevalências ajustadas dependem das variâncias e covariâncias das estimativas dos parâmetros do modelo estimado. Aspectos essenciais do processo de construção desses intervalos de confiança são descritos na próxima seção.

3.1.3 Estimação por intervalo utilizando método Delta

As probabilidades estimadas pelo modelo são funções dos estimadores de máxima verossimilhança, podendo ser reescritas como

$$h(\hat{\boldsymbol{\beta}}) = \hat{p}(\mathbf{x}) = \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2)}}, \quad [8]$$

em que $\mathbf{x} = (x_1, x_2)^t$ e $\hat{\boldsymbol{\beta}} = (\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2)^t$ representa os estimadores de máxima verossimilhança no modelo definido em [2]. A variância de $h(\hat{\boldsymbol{\beta}})$ também é função dos estimadores de máxima verossimilhança, não sendo possível estimá-la de forma analítica. Nestas situações, o método Delta é um procedimento geral e flexível para construção de intervalos de confiança. Este método está extensamente detalhado na literatura, podendo-se destacar Xu e Long (2005) para sua aplicação na predição de probabilidades.

A função $h(\hat{\boldsymbol{\beta}})$ pode ser expandida em série de Taylor até o termo de primeira ordem, mostrada abaixo, que geralmente produz estimativas relativamente acuradas para a variância:

$$h(\hat{\boldsymbol{\beta}}) \approx h(\boldsymbol{\beta}) + (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^t h'(\boldsymbol{\beta}), \quad [9]$$

em que $h'(\boldsymbol{\beta}) = \frac{\partial h(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}}$. Assim,

$$\sqrt{n} [h(\hat{\boldsymbol{\beta}}) - h(\boldsymbol{\beta})] \approx \sqrt{n} \left[(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^t \frac{\partial h(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \right], \quad [10]$$

tal que,

$$h(\hat{\boldsymbol{\beta}}) \rightarrow N \left(h(\boldsymbol{\beta}), \frac{\partial h(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}^t} \text{Var}(\hat{\boldsymbol{\beta}}) \frac{\partial h(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \right). \quad [11]$$

Para estimar a variância de $h(\hat{\boldsymbol{\beta}})$, as derivadas parciais $\frac{\partial h(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}'}$ são avaliadas em

$\boldsymbol{\beta} = \hat{\boldsymbol{\beta}}$, tal que

$$\text{Var} \left(h(\hat{\boldsymbol{\beta}}) \right) = \left(\frac{\partial h(\hat{\boldsymbol{\beta}})}{\partial \hat{\boldsymbol{\beta}}'} \right) \text{Var}(\hat{\boldsymbol{\beta}}) \left(\frac{\partial h(\hat{\boldsymbol{\beta}})}{\partial \hat{\boldsymbol{\beta}}} \right) = (\nabla h(\hat{\boldsymbol{\beta}}))' \hat{\boldsymbol{\Sigma}} (\nabla h(\hat{\boldsymbol{\beta}})), \quad [12]$$

em que $\hat{\boldsymbol{\Sigma}}$ é a matriz de estimativas de variâncias e covariâncias dos parâmetros do modelo. As derivadas parciais de primeira ordem da equação [9] são mostradas abaixo:

$$\frac{\partial h(\hat{\boldsymbol{\beta}})}{\partial \beta_0} = \frac{e^{-\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2}}{\left[1 + e^{-\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2} \right]^2} = \frac{1}{1 + e^{-\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2}} \times \frac{e^{-\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2}}{1 + e^{-\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2}} = \hat{p}(\mathbf{x}) \times (1 - \hat{p}(\mathbf{x})). \quad [13]$$

Similarmente,

$$\frac{\partial h(\hat{\boldsymbol{\beta}})}{\partial \beta_1} = x_1 \times \frac{1}{1 + e^{-\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2}} \times \frac{e^{-\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2}}{1 + e^{-\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2}} = x_1 \times \hat{p}(\mathbf{x}) \times (1 - \hat{p}(\mathbf{x})) \quad [14]$$

e

$$\frac{\partial h(\hat{\boldsymbol{\beta}})}{\partial \beta_2} = x_2 \times \frac{1}{1 + e^{-\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2}} \times \frac{e^{-\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2}}{1 + e^{-\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2}} = x_2 \times \hat{p}(\mathbf{x}) \times (1 - \hat{p}(\mathbf{x})). \quad [15]$$

Assim, $\nabla h(\hat{\boldsymbol{\beta}})$ descrito na equação [12] é definido como

$$\nabla h(\hat{\boldsymbol{\beta}}) = \begin{bmatrix} \frac{1}{1 + e^{-\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2}} \times \frac{e^{-\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2}}{1 + e^{-\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2}} \\ x_1 \times \frac{1}{1 + e^{-\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2}} \times \frac{e^{-\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2}}{1 + e^{-\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2}} \\ x_2 \times \frac{1}{1 + e^{-\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2}} \times \frac{e^{-\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2}}{1 + e^{-\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2}} \end{bmatrix} = \begin{bmatrix} 1 \times \hat{p}(\mathbf{x}) \times (1 - \hat{p}(\mathbf{x})) \\ x_1 \times \hat{p}(\mathbf{x}) \times (1 - \hat{p}(\mathbf{x})) \\ x_2 \times \hat{p}(\mathbf{x}) \times (1 - \hat{p}(\mathbf{x})) \end{bmatrix}. \quad [16]$$

No contexto deste trabalho deseja-se obter as estimativas de prevalência para cada categoria do preditor x_1 , ajustado por x_2 . Assim, para cada um dos métodos de predição, condicional e marginal, o termo $\hat{p}(\mathbf{x})$ deve ser substituído de forma adequada (Xu e Long, 2005).

Para o método de predição condicional, quando $x_1 = 0$

$$\hat{p}_{c_0}(\mathbf{x}) = \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 \times 0 + \hat{\beta}_2 \bar{x}_2)}}, \quad [17]$$

e quando, $x_1 = 1$,

$$\hat{p}_{c_1}(\mathbf{x}) = \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 \times 1 + \hat{\beta}_2 \bar{x}_2)}}. \quad [18]$$

De maneira similar, para o método de predição marginal, para $x_1 = 0$ (não exposto),

$$\hat{p}_{M_0}(\mathbf{x}) = \frac{1}{n} \times \sum_{i=1}^n \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 \times 0 + \hat{\beta}_2 x_{2i})}}, \quad [19]$$

enquanto que, para $x_1 = 1$ (exposto),

$$\hat{p}_{M_1}(\mathbf{x}) = \frac{1}{n} \times \sum_{i=1}^n \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 \times 1 + \hat{\beta}_2 x_{2i})}}. \quad [20]$$

O método Delta, descrito nessa seção, foi implementado em rotinas computacionais SAS e R para obter as variâncias de $h(\hat{\beta})$ e, assim, estimar os intervalos de confiança para prevalências ajustadas. Alguns programas, como por exemplo, o SUDAAN, calculam estimativas de prevalências ajustadas pelos métodos de predição marginal e condicional e intervalos de confiança. A seção 3.2 descreve os métodos disponíveis, além do SUDAAN, para os programas SAS, Stata e R.

3.2 Aspectos computacionais

Rotinas computacionais (macros ou funções) para estimar prevalências ajustadas por meio de modelos multivariáveis estão disponíveis em programas usuais, tais como SAS, SUDAAN, R e STATA. Nesta seção é apresentada uma breve descrição das rotinas computacionais destes programas.

3.2.1 SAS - Macro %ADJ_PROP

O programa SAS não disponibiliza, de forma direta, procedimentos para estimar prevalências ajustadas e respectivos intervalos de confiança. Para esta finalidade, Dingyi Zhao (1985) criou a macro %ADJ_PROP escrita na linguagem SAS, originalmente para a versão 6.1, utilizando os métodos de predição marginal e condicional. Esta macro utiliza basicamente os comandos *PROC LOGISTIC* para a obtenção de prevalências ajustadas e o procedimento *PROC IML* para a obtenção dos respectivos intervalos de confiança. A macro %ADJ_PROP pode ser utilizada tanto para o cálculo de prevalências ajustadas, por meio de regressão logística, bem como para o cálculo de médias ajustadas, por meio de regressão linear. Por ter sido criada em uma versão antiga do programa SAS, foram realizadas modificações para que fosse possível obter as prevalências ajustadas.

Em 1998, Dingyi Zhao também criou a macro %ADJWPROP, para obter prevalências ajustadas no contexto de dados de planos amostrais com ponderação, mas que também pode ser utilizada para prevalências ajustadas sem o uso de pesos. Catellier (1998) realizou modificações na macro %ADJWPROP, renomeando-a para %ADJ. Face à documentação incompleta as macros %ADJWPROP e %ADJ não serão abordadas no trabalho.

3.2.2 SUDAAN (SAS-Callable): CONDMARG e PREDMARG

O SUDAAN versão SAS-Callable, é um programa para análise estatística de dados executado em conjunto com o programa SAS. No procedimento PROC

LOGISTIC do SUDAAN (*alias* PROC RLOGIST) estão disponíveis os comandos CONDMARG e PREDMARG, por meio dos quais é possível obter estimativas de prevalências ajustadas e intervalos de confiança para os métodos condicional e marginal, respectivamente. Neste trabalho os resultados do programa SUDAAN serão considerados como referência para as comparações com os resultados dos demais programas.

3.2.3 Stata - Função *margins*

Williams (2012) apresentou e discutiu diferentes métodos para estimação de prevalências ajustadas e efeitos marginais por meio do programa Stata, com destaque para o comando *margins* disponível a partir da versão 11. Com o comando *margins* é possível obter prevalências ajustadas e seus respectivos intervalos de confiança por meio do método condicional. O Stata não disponibiliza, de forma direta, estimativas de prevalências ajustadas para o método marginal.

3.2.4 R – Função *margins*

Thomas Leeper (2014) adaptou no programa R os mesmos procedimentos disponíveis para o Stata, apresentados por Williams (2012). Sendo assim, é possível obter apenas as prevalências ajustadas e seus respectivos intervalos de confiança por meio do método de predição condicional. Esta abordagem não será utilizada neste trabalho.

Rotinas computacionais para SAS e R foram desenvolvidas para estimar prevalências ajustadas por meio dos métodos descritos nas seções 3.1.1 e 3.1.2. Os respectivos intervalos de confiança foram derivados utilizando o método Delta, descrito na seção 3.1.3. Estes métodos, bem como aqueles descritos na revisão de literatura, serão mostrados em detalhes por meio de um exemplo.

4. OBJETIVOS

Objetivo geral

Apresentar e discutir métodos para estimar prevalências ajustadas em pesquisas clínicas e epidemiológicas e desenvolver rotinas computacionais em SAS e R para os métodos de predição marginal e condicional.

Objetivos específicos

- a) Caracterizar a importância de utilizar prevalências ajustadas na presença de confundimento.
- b) Revisar os métodos existentes para obtenção de prevalências ajustadas.
- c) Revisar aspectos computacionais disponíveis para obtenção de prevalências ajustadas.
- d) Desenvolver rotinas computacionais em SAS e R para estimar prevalências ajustadas e respectivos intervalos de confiança.

5. REFERÊNCIAS BIBLIOGRÁFICAS

An R. Diabetes diagnosis and nutrition facts label use among US adults, 2005–2010. *Public Health Nutr.* 2015 outubro 20: 1-8.

Chambless LE, Heiss G, Folsom AR, Rosamond W, Szklo M, Sharrett AR, et al. Association of coronary heart disease incidence with carotid arterial wall thickness and major risk factors: the Atherosclerosis Risk in Communities (ARIC) Study, 1987-1993. *Am J Epidemiol.* 1997 setembro 15; 146(6): 483-94.

Chor D, Ribeiro ALP, Carvalho MS, BB Duncan, Lotufo PA, Nobre AA, et al. Prevalence, Awareness, Treatment and Influence of Socioeconomic Variables on Control of High Blood Pressure: Results of the ELSA-Brasil Study. *PLOS ONE.* 2015 junho 23; 10(6): e0127382.

Epstein RA, Bobo WV, Martin PR, Morrow JA, Wang W, Chandrasekhar R, et al. Increasing pregnancy-related use of prescribed opioid analgesics. *Ann Epidemiol.* 2013 agosto. 23(8): 498-503.

Epstein RA, Bobo WV, Shelton RC, Arbogast PG, Morrow JA, Wang W, et al. Increasing use of atypical antipsychotics and anticonvulsants during pregnancy. *Pharmacoepidem Dr S.* 2013 julho. 22(7): 794-801.

Lane PW, Nelder JA. Analysis of Covariance and Standardization as Instances of Prediction. *Biometrics.* 1982; 38(3):613-21.

Lebedeva ER, Kobzeva NR, Gilev D, Olesen J. Prevalence of primary headache disorders diagnosed according to ICHD-3 beta in three different social groups. *Cephalalgia.* 2015 outubro 6.

Leeper T. Margins. Disponível de: <https://github.com/leeper/margins>. Acessado setembro 20, 2015.

McCullagh P, Nelder JA. Generalized Linear Models. 2th ed. 1989.

Nambi V, Hoogeveen RC, Chambless L, Hu Y, Bang H, Coresh J, et al. Lipoprotein-associated phospholipase A2 and high-sensitivity C-reactive protein improve the stratification of ischemic stroke risk in the Atherosclerosis Risk in Communities (ARIC) study. *Stroke*. 2009 fevereiro; 40(2): 376-81.

Ohira T, Folsom AR, Cushman M, White RH, Hannan PJ, Rosamond WD. Reproductive History, Hormone Replacement, and Incidence of Venous Thromboembolism: The Longitudinal Investigation of Thromboembolism Etiology. *Br J Haematol*. 2010 maio; 149(4): 606-612

Research Triangle Institute. SUDAAN Language Manual, Volumes 1 and 2, Release 11. 2012.

Schoenbach VJ, Rosamond WD. Understanding the Fundamentals of Epidemiology - an evolving text. 2000, Fall Edition.

Shon H-S, Choi HY, Kim JR, Ryu SY, Lee Y-J, Lee MJ, et al. Comparison and analysis of the prevalence of hepatitis C virus infection by region in the Republic of Korea during 2005-2012. *Clin Mol Hepatol*. 2015 setembro; 21(3): 249-56.

Smith AK, Cenzer IS, John Boscardin W, Ritchie CS, Wallhagen ML, Covinsky KE. Increase in Disability Prevalence Before Hip Fracture. *J Am Geriatr Soc*. 2015 outubro; 63(10): 2029-35.

Souza CL, Carvalho FM, Araújo TM, Reis EJM, Lima VCM, Porto LA. Factors associated with vocal fold pathologies in teachers. Rev Saúde Pública. 2011; 45(5).

Szklo M, Nieto FJ. Epidemiology Beyond the Basics. 2th ed. 2007.

Wilcosky TC, Chambless LE. A comparison of Direct Adjustment and Regression Adjustment of Epidemiologic Measures. J Chron Dis. 1985; 38: 849-856.

Williams R. Using margins command. The Stata Journal. 2012; 12(2): 308-331.

Woodward, M. Epidemiology Study Design and Data Analysis. 3th ed. 2014.

Xu J, Long JS. Using the Delta Method to Construct Confidence Intervals for Predicted Probabilities, Rates, and Discrete Changes Indiana University. 2005 agosto 22.

Zhao D. Logistic Regression Adjustment of Proportions and its Macro Procedure Disponível de: <http://www2.sas.com/proceedings/sugi22/POSTERS/PAPER227.PDF>. Acessado 7 de setembro de 2015.

6. ARTIGO ORIGINAL

Métodos para estimar prevalências ajustadas

Methods to estimate adjusted prevalences

Natália Bordin Barbieri

Programa de Pós-Graduação em Epidemiologia

Universidade Federal do Rio Grande do Sul

RESUMO

Objetivo: Apresentar e discutir métodos para estimar prevalências ajustadas em pesquisas clínicas e epidemiológicas, bem como desenvolver rotinas computacionais em SAS e R.

Métodos: No contexto de estudo transversal, foi simulada uma amostra de 2.000 observações independentes, considerando o desfecho dicotômico diabetes, sexo como a variável de exposição e idade como variável de ajuste. As estimativas de prevalências ajustadas (IC 95%) foram estimadas pelos métodos de predição condicional e marginal, utilizando as rotinas desenvolvidas em SAS e R. O método Delta foi usado para construir os intervalos de confiança. Os resultados foram comparados com aqueles do SUDAAN (SAS-Callable), Stata e a macro %ADJ_PROP (SAS).

Resultados: No exemplo simulado, 68,2% são do sexo feminino e a idade média (DP) foi 57,6 (5,0) anos, sendo 54,2 (3,9) anos em homens e 59,2 (4,6) anos em mulheres. A estimativa da prevalência global do desfecho foi de 25,3% (IC 95%:23,4-27,3); sendo 13,8% (IC 95%:11,7-16,7) e 30,7% (IC 95%:28,3-33,2), respectivamente para homens e mulheres. As estimativas de prevalências ajustadas por idade, por meio do método de predição condicional, foram de 19,6% (IC 95%:16,2-23,6) para homens, e 23,6% (IC 95%:21,2-26,1) para mulheres. Pelo método de predição marginal, as estimativas foram de 22,4% (IC 95%:18,7-26,5) para homens, e 26,3% (IC 95%:24,1-28,6) para mulheres.

Conclusão: A discrepância entre as estimativas não ajustadas é devida ao confundimento pela da idade. Estimativas comparáveis (livre de confundimento) podem ser obtidas por meio das prevalências ajustadas pela idade. No entanto, a estimativa pelo método de predição condicional não engloba a prevalência global. Em virtude disso, o método de predição marginal é, geralmente, mais adequado. As rotinas desenvolvidas podem ser úteis para estimação de prevalências ajustadas, sendo a versão para R uma alternativa aos softwares comerciais.

Palavras-chave: Prevalências ajustadas, Método de predição condicional, Método de predição marginal, Método Delta.

ABSTRACT

Objective: To present and discuss methods to estimate adjusted prevalences for clinical and epidemiological research, and develop computational routines in SAS and R.

Methods: In the context of cross-sectional study, it was simulated a sample of 2,000 independent observations, considering the dichotomous outcome diabetes, sex as the exposure variable and age as an adjustment variable. Adjusted prevalences were estimated by the conditional and marginal methods, using routines developed in SAS and R. Confidence intervals were constructed using the Delta method. The results were compared with those of the SUDAAN (SAS-callable), Stata and macro %ADJ_PROP (SAS).

Results: In simulated example, 68.2% are female and the mean (SD) age was 57.6 (4.97) years old, being that 54.2 (3.94) years for men and 59.2 (4.60) years in women. The estimated global prevalence of outcome was 25.3% (CI 95%: 23.4-27.3) and 13.8% (CI 95%: 11.7-16.7) and 30.7% (CI 95%: 28.3-33.2), respectively for men and women. Estimates of adjusted prevalence for age, through the conditional method, were 19.6% (CI 95%: 16.2-23.6) for men, and 23.6% (CI 95%: 21.2-26.1) for women. For marginal method, the estimates were 22.4% (CI 95%: 18.7-26.5) for men and 26.3% (CI 95%: 24.1-28.6) for women.

Conclusion: The observed discrepancy in estimates by sex, unadjusted, can be attributed to confounding due to difference in age distribution between sexes. Comparable estimates (without confounding) of the prevalences can be obtained through prevalence adjusted for age. However, the estimate for the conditional method does not comprise the global prevalence. As a result, the marginal method is in general more suitable. The developed routines can be useful for estimating adjusted prevalences, particularly the R version (an alternative to commercial software).

Keywords: Adjusted Prevalence, Conditional prediction method, Marginal prediction method, Delta method.

INTRODUÇÃO

Estimativas de prevalências são frequentemente reportadas em estudos clínicos ou epidemiológicos¹⁻⁶. Para que essas estimativas sejam comparáveis em diferentes populações, é necessário levar em conta os desequilíbrios nas distribuições de variáveis que possam alterar a frequência do evento. Em outras palavras, na presença de confundimento, estimativas não ajustadas podem não ser adequadas, e podem dificultar a compreensão do fenômeno.

Reportar prevalências ajustadas possibilita ao leitor um julgamento da frequência do evento livre de possíveis confundimentos, uma vez que a magnitude da associação entre desfecho e o fator em estudo pode variar entre os diferentes níveis do fator. Estratificação e modelos de ajuste multivariável são ferramentas comuns utilizadas para controlar os efeitos de confundimento⁷.

Outras abordagens utilizadas para obter estimativas ajustadas são os métodos de padronização direta (*direct adjustment*) ou indireta (*indirect adjustment*). Estes métodos são frequentemente utilizados em dados de estudos ecológicos, como por exemplo, para a comparação de taxas de mortalidade, ajustadas por sexo ou idade⁸⁻¹⁰.

Wilcosky e Chambless¹⁰ afirmam que as principais vantagens do método de ajuste direto são a simplicidade computacional e as poucas suposições estatísticas. Entretanto, com a maior disponibilidade e diversidade de recursos computacionais, esses aspectos atualmente não são mais relevantes, além de ter como desvantagem a necessidade de categorização de variáveis quantitativas. Os autores também discutiram os métodos de predição condicional ("*the conditional prediction method*") e marginal ("*the marginal prediction method*") para estimar prevalências ajustadas. Esses métodos utilizam o modelo de regressão logística multivariável e são mais convenientes para testes de interações ou de diferenças entre grupos, permitem explorar a natureza da relação funcional entre as variáveis de controle e o desfecho, e não exigem a categorização de variáveis quantitativas. Lane e Nelder⁷ também apresentaram a abordagem de ajuste multivariável em modelos lineares generalizados para estimar

proporções ajustadas por meio dos métodos de predição marginal e condicional.

Em uma revisão de programas de análise de dados foi constatado que somente o SUDAAN disponibiliza, de forma direta, estimativas de prevalências ajustadas (e intervalos de confiança) por meio dos métodos marginal e condicional. O Stata fornece apenas estimativas pelo método condicional. Para os programas SAS e R, é necessário o uso de rotinas específicas (macros ou funções).

Esse artigo apresenta uma revisão dos diferentes métodos para estimar prevalências ajustadas, bem como aspectos computacionais disponíveis. Também são apresentadas rotinas computacionais em SAS e R, desenvolvidas para estimar prevalências ajustadas por meio dos métodos condicional e marginal. Um exemplo simulado foi usado para ilustrar a aplicação e discutir vantagens e desvantagens dos métodos.

MÉTODOS

Para descrever estes métodos será considerado o contexto de um estudo transversal em que a variável dependente (Y) representa presença ($Y = 1$) ou ausência ($Y = 0$) de diabetes, a variável dicotômica (x_1) representa sexo ($1 = masculino, 0 = feminino$), e a variável x_2 representa idade (quantitativa), podendo ser generalizado para mais preditores. Dada uma amostra aleatória com n indivíduos, obtida por meio de um estudo transversal, deseja-se estimar a prevalência do desfecho para cada categoria da exposição, ajustada pela variável x_2 . O modelo logístico definido na equação [1] pode ser usado para essa finalidade, ou seja,

$$\log \frac{P[Y = 1 | \mathbf{x}]}{1 - P[Y = 1 | \mathbf{x}]} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 \quad [1]$$

em que $\mathbf{x} = (x_1, x_2)^t$ representa o vetor de covariáveis do modelo. O método de estimação de máxima verossimilhança pode ser usado para estimar os parâmetros do modelo, obtendo-se o modelo estimado descrito na equação [2]

$$\log \frac{P[Y = 1 | \mathbf{x}]}{1 - P[Y = 1 | \mathbf{x}]} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2. \quad [2]$$

Detalhes sobre o método de estimação podem ser obtidos em Mcculagh e Nelder¹². As probabilidades estimadas pelo modelo da equação [2] são obtidas por

$$\hat{p}_{ij}(x_{1i}, x_{2i}) = \hat{P}[Y = 1 | x_{1i} = j, x_{2i}] = \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 \times j + \hat{\beta}_2 x_{2i})}}, \quad [3]$$

em que $j = 0, 1$ e $i = 1, 2, \dots, n$.

Método de predição condicional

Para obter estimativas de prevalências ajustadas pelo método de predição condicional (\hat{p}_{C_i}), a média amostral do preditor x_2 , denotada por \bar{x}_2 , é usada na equação [3] para cada categoria da variável x_1 , como descrito nas equações [4] e [5]. Detalhes do método podem ser encontrados no manual do programa SUDAAN¹³.

$$\hat{p}_{C_0}(\mathbf{x}) = \hat{P}[Y = 1 | x_1 = 0, x_2 = \bar{x}_2] = \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 \times 0 + \hat{\beta}_2 \bar{x}_2)}} = \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_2 \bar{x}_2)}}, \quad [4]$$

e

$$\hat{p}_{C_1}(\mathbf{x}) = \hat{P}[Y = 1 | x_1 = 1, x_2 = \bar{x}_2] = \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 \times 1 + \hat{\beta}_2 \bar{x}_2)}} = \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 + \hat{\beta}_2 \bar{x}_2)}}. \quad [5]$$

Método de predição marginal

Pelo método marginal, as prevalências ajustadas para as categorias da exposição são estimadas pelas equações [6] e [7]:

$$\hat{p}_{M_0}(\mathbf{x}) = \hat{P}[Y = 1 | x_1 = 0, x_{2i}] = \frac{1}{n} \sum_{i=1}^n \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 \times 0 + \hat{\beta}_2 x_{2i})}} = \frac{1}{n} \sum_{i=1}^n \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_2 x_{2i})}}, \quad [6]$$

e

$$\hat{p}_{M_1}(\mathbf{x}) = \hat{P}[Y = 1 | x_1 = 1, x_{2i}] = \frac{1}{n} \sum_{i=1}^n \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 \times 1 + \hat{\beta}_2 x_{2i})}} = \frac{1}{n} \sum_{i=1}^n \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 + \hat{\beta}_2 x_{2i})}}. \quad [7]$$

As estimativas do intervalo de confiança para as prevalências ajustadas dependem das variâncias e covariâncias das estimativas dos parâmetros do modelo estimado¹⁴. Aspectos essenciais do processo de construção desses intervalos de confiança, pelo método Delta, são descritos a seguir.

Estimação por intervalo utilizando método Delta

Como descrito na equação [3], as probabilidades estimadas pelo modelo são funções dos estimadores de máxima verossimilhança, podendo ser reescritas como

$$h(\hat{\boldsymbol{\beta}}) = \hat{p}(\mathbf{x}) = \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2)}}, \quad [8]$$

em que $\mathbf{x} = (x_1, x_2)^t$ e $\hat{\boldsymbol{\beta}} = (\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2)^t$ representa os estimadores de máxima verossimilhança dos parâmetros do modelo definido em [2]. A variância de $h(\hat{\boldsymbol{\beta}})$ também é função dos estimadores de máxima verossimilhança, não sendo possível estimar de forma analítica. Nestas situações, o método Delta é um procedimento geral e flexível para construção de intervalos de confiança. O método está extensamente detalhado na literatura, podendo-se destacar em Xu e Long¹⁴.

A função $h(\hat{\boldsymbol{\beta}})$ pode ser expandida em série de Taylor até o termo de primeira ordem, mostrada abaixo, que geralmente produz estimativas relativamente acuradas para a variância:

$$h(\hat{\boldsymbol{\beta}}) \approx h(\boldsymbol{\beta}) + (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^t h'(\boldsymbol{\beta}), \quad [9]$$

em que $h'(\boldsymbol{\beta}) = \frac{\partial h(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}}$. Assim,

$$\sqrt{n} [h(\hat{\boldsymbol{\beta}}) - h(\boldsymbol{\beta})] \approx \sqrt{n} \left[(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})' \frac{\partial h(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \right], \quad [10]$$

tal que,

$$h(\hat{\boldsymbol{\beta}}) \rightarrow N \left(h(\boldsymbol{\beta}), \frac{\partial h(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}'} \text{Var}(\hat{\boldsymbol{\beta}}) \frac{\partial h(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \right). \quad [11]$$

Para estimar a variância de $h(\hat{\boldsymbol{\beta}})$, as derivadas parciais $\frac{\partial h(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}'}$ são avaliadas em

$\boldsymbol{\beta} = \hat{\boldsymbol{\beta}}$, tal que

$$\text{Var} (h(\hat{\boldsymbol{\beta}})) = \left(\frac{\partial h(\hat{\boldsymbol{\beta}})}{\partial \hat{\boldsymbol{\beta}}'} \text{Var}(\hat{\boldsymbol{\beta}}) \frac{\partial h(\hat{\boldsymbol{\beta}})}{\partial \hat{\boldsymbol{\beta}}} \right) = [\nabla h(\hat{\boldsymbol{\beta}})]' \hat{\boldsymbol{\Sigma}} [\nabla h(\hat{\boldsymbol{\beta}})] \quad [12]$$

em que $\hat{\boldsymbol{\Sigma}}$ é a matriz de estimativas de variâncias e covariâncias dos parâmetros do modelo. As derivadas parciais de primeira ordem da equação [9] são mostradas abaixo:

$$\frac{\partial}{\partial \beta_0} h(\hat{\boldsymbol{\beta}}) = \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2)}} \times \frac{e^{-(\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2)}}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2)}} = \hat{p}(\mathbf{x}) \times (1 - \hat{p}(\mathbf{x})) \quad [13]$$

Similarmente,

$$\frac{\partial}{\partial \beta_1} h(\hat{\boldsymbol{\beta}}) = x_1 \times \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2)}} \times \frac{e^{-(\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2)}}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2)}} = x_1 \times \hat{p}(\mathbf{x}) \times (1 - \hat{p}(\mathbf{x})) \quad [14]$$

e

$$\frac{\partial}{\partial \beta_2} h(\hat{\boldsymbol{\beta}}) = x_2 \times \frac{1}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2)}} \times \frac{e^{-(\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2)}}{1 + e^{-(\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2)}} = x_2 \times \hat{p}(\mathbf{x}) \times (1 - \hat{p}(\mathbf{x})). \quad [15]$$

Assim,

$$\nabla h(\hat{\boldsymbol{\beta}}) = \begin{bmatrix} \frac{1}{1+e^{-(\hat{\beta}_0+\hat{\beta}_1x_1+\hat{\beta}_2x_2)}} \times \frac{e^{-(\hat{\beta}_0+\hat{\beta}_1x_1+\hat{\beta}_2x_2)}}{1+e^{-(\hat{\beta}_0+\hat{\beta}_1x_1+\hat{\beta}_2x_2)}} \\ x_1 \times \frac{1}{1+e^{-(\hat{\beta}_0+\hat{\beta}_1x_1+\hat{\beta}_2x_2)}} \times \frac{e^{-(\hat{\beta}_0+\hat{\beta}_1x_1+\hat{\beta}_2x_2)}}{1+e^{-(\hat{\beta}_0+\hat{\beta}_1x_1+\hat{\beta}_2x_2)}} \\ x_2 \times \frac{1}{1+e^{-(\hat{\beta}_0+\hat{\beta}_1x_1+\hat{\beta}_2x_2)}} \times \frac{e^{-(\hat{\beta}_0+\hat{\beta}_1x_1+\hat{\beta}_2x_2)}}{1+e^{-(\hat{\beta}_0+\hat{\beta}_1x_1+\hat{\beta}_2x_2)}} \end{bmatrix} = \begin{bmatrix} 1 \times \hat{p}(\mathbf{x}) \times (1 - \hat{p}(\mathbf{x})) \\ x_1 \times \hat{p}(\mathbf{x}) \times (1 - \hat{p}(\mathbf{x})) \\ x_2 \times \hat{p}(\mathbf{x}) \times (1 - \hat{p}(\mathbf{x})) \end{bmatrix}. \quad [16]$$

No contexto deste trabalho deseja-se obter estimativas de prevalência ajustada para cada categoria do preditor x_1 ajustado por x_2 . Assim, para cada um dos métodos de predição, condicional e marginal, o termo $\hat{p}(\mathbf{x})$ deve ser substituído de forma adequada¹⁴, como descrito a seguir.

Para o método de predição condicional, quando $x_1 = 0$

$$\hat{p}_{C_0}(\mathbf{x}) = \frac{1}{1+e^{-(\hat{\beta}_0+\hat{\beta}_1 \times 0 + \hat{\beta}_2 \bar{x}_2)}}, \quad [17]$$

e quando, $x_1 = 1$,

$$\hat{p}_{C_1}(\mathbf{x}) = \frac{1}{1+e^{-(\hat{\beta}_0+\hat{\beta}_1 \times 1 + \hat{\beta}_2 \bar{x}_2)}}. \quad [18]$$

De maneira similar, para o método de predição marginal, quando $x_1 = 0$

$$\hat{p}_{M_0}(\mathbf{x}) = \frac{1}{n} \times \sum_{i=1}^n \frac{1}{1+e^{-(\hat{\beta}_0+\hat{\beta}_1 \times 0 + \hat{\beta}_2 x_{2i})}}, \quad [19]$$

e quando, $x_1 = 1$,

$$\hat{p}_{M_1}(\mathbf{x}) = \frac{1}{n} \times \sum_{i=1}^n \frac{1}{1+e^{-(\hat{\beta}_0+\hat{\beta}_1 \times 1 + \hat{\beta}_2 x_{2i})}}. \quad [20]$$

O método Delta foi implementado em rotinas computacionais SAS e R para obter a variância de $h(\hat{\boldsymbol{\beta}})$ e, assim, estimar os intervalos de confiança para prevalências

ajustadas. Alguns programas, como por exemplo, o SUDAAN, calculam estimativas de prevalências ajustadas pelos métodos de predição marginal e condicional, incluindo intervalos de confiança. A seguir são descritos os métodos computacionais disponíveis, além do SUDAAN, para os programas SAS, Stata e R.

Aspectos computacionais

Métodos para estimar prevalências ajustadas por meio de modelos estão disponíveis em programas usuais, tais como SAS, SUDAAN, R e STATA. A seguir é apresentada uma breve descrição das rotinas computacionais encontradas na revisão de métodos disponíveis.

O programa SAS não disponibiliza de forma direta procedimentos para estimar prevalências ajustadas e respectivos intervalos de confiança. Para esta finalidade, Dingyi Zhao¹⁵ criou a macro %ADJ_PROP para o programa SAS na versão 6.1, utilizando os métodos de predição marginal e condicional. Esta macro utiliza basicamente os comandos *PROC LOGISTIC* para a obtenção de prevalências ajustadas e o procedimento *PROC IML* para a obtenção dos respectivos intervalos de confiança.

A macro %ADJ_PROP pode ser utilizada tanto para o cálculo de prevalências ajustadas, por meio de regressão logística, quanto para o cálculo de médias ajustadas, por meio de regressão linear. Por ter sido criada em uma versão antiga do programa SAS, foram realizadas modificações para que fosse possível obter as prevalências ajustadas.

O SUDAAN SAS-Callable é um programa para análise estatística de dados executado em conjunto com o programa SAS. No procedimento *PROC LOGISTIC* (alias *PROC RLOGIST*) do SUDAAN¹³ estão disponíveis os comandos *PREDMARG* e *CONDMARG*, por meio dos quais é possível obter estimativas de prevalências ajustadas e seus respectivos intervalos de confiança, utilizando os métodos condicional e marginal. Neste artigo os resultados do programa SUDAAN serão considerados como

referência nas comparações com os resultados dos demais programas.

Williams¹⁶ aborda diferentes métodos para estimação de prevalências ajustadas e efeitos marginais por meio do programa Stata, com destaque para o comando *margins* disponível a partir da versão 11. Com o comando *margins* é possível obter prevalências ajustadas e respectivos intervalos de confiança por meio do método condicional.

Thomas Leeper¹⁷ adaptou no programa R os mesmos procedimentos disponíveis para o Stata, apresentados por Williams¹⁶. Sendo assim, é possível obter apenas as prevalências ajustadas e seus respectivos intervalos de confiança por meio do método de predição condicional.

Rotinas computacionais para SAS e R foram desenvolvidas para estimar prevalências ajustadas por meio dos métodos marginal e condicional, disponibilizadas nos Suplementos I-III. Para exemplificar os métodos, foi simulada uma amostra de tamanho $n=2000$, considerando um contexto com desfecho e exposição dicotômicos e uma variável de confundimento quantitativa, como descrito na Tabela 1. Detalhes sobre a geração da amostra e dos aspectos computacionais para obtenção dos resultados das análises podem ser obtidos no trabalho original, disponível no Suplemento IV. Estimativas de prevalências ajustadas foram obtidas por meio dos programas SUDAAN, R, SAS e STATA.

RESULTADOS

A amostra é composta por 636 (31,8%) homens e 1364 (68,2%) mulheres. A idade média (DP) foi igual a 57,6 (4,97) anos, sendo 54,2 (3,94) anos e 59,2 (4,60) anos, respectivamente para homens e mulheres. Das 2000 observações, 507 apresentaram o desfecho, correspondendo a uma estimativa de prevalência de 25,3% (IC 95% 23,4-27,3). Estimativas de prevalências não ajustadas foram 13,8% (IC 95% 11,7-16,7) e 30,7% (IC 95% 28,3-33,2), respectivamente para homens e mulheres. A

discrepância observada nas estimativas por sexo pode ser atribuída ao confundimento devido à diferença na distribuição da idade entre sexos.

A Tabela 2 mostra prevalências do desfecho para cada sexo, ajustadas pela idade, por meio dos métodos marginal e condicional. Os resultados obtidos pelo SUDAAN serão considerados como referência para comparações com os resultados de outros programas.

As estimativas pontuais de prevalências ajustadas obtidas pelos demais programas para os métodos marginal e condicional foram iguais àsquelas do programa SUDAAN, com exceção das estimativas geradas pela macro %ADJ_PROP pelo método marginal. As estimativas por intervalo mostraram pequenas discrepâncias em relação às do SUDAAN, em geral na terceira casa decimal.

Na Tabela 3 são apresentadas as médias ponderadas (proporcional às frações amostrais) das prevalências ajustadas para cada categoria, obtidas por meio do método de predição condicional e marginal em comparação com a proporção global observada. É importante notar que, para homens e mulheres, as prevalências ajustadas estimadas pelo método condicional são menores do que a estimativa de prevalência global não ajustada, mas isso não ocorre no método marginal. O método de predição condicional tem a média ponderada (22,3%) muito distante da prevalência global observada (25,3%). Esta é uma desvantagem do método condicional, que muitas vezes pode confundir o leitor. O método marginal não é suscetível a este problema, exceto talvez por arredondamentos.

DISCUSSÃO

Prevalências ajustadas são frequentemente reportadas em estudos epidemiológicos¹⁻⁶, e são importantes para o leitor fazer uma análise do fenômeno livre de confundimento. Atualmente os métodos de estimação marginal e condicional são as abordagens mais utilizadas.

As rotinas SAS e R desenvolvidas no trabalho podem ser úteis para estimar prevalências ajustadas, produzindo resultados muito similares aos do SUDAAN. Em especial, a versão para o R é uma alternativa aos softwares comerciais, uma vez que esse programa é livre.

As estimativas pontuais das rotinas desenvolvidas para o SAS e R são idênticas aos do SUDAAN, e as discrepâncias entre as estimativas por intervalo não parecem ser relevantes na prática.

Os resultados relativamente discrepantes gerados pela macro %ADJ_PROP pelo método marginal podem ser explicados pela diferença dos métodos utilizados. O método condicional utiliza uma constante de aproximação k , produzindo resultados similares aos do SUDAAN. Outra desvantagem é a falta de documentação dos códigos, dificultando sua compreensão para descrição dos métodos ou eventuais modificações.

A média ponderada obtida por meio do método marginal foi mais próxima da prevalência global no exemplo estudado, o que faz com que o método marginal tenha preferência sobre o método condicional.

Em suma, prevalências ajustadas são importantes para uma análise mais realista do fenômeno em estudo. A discussão apresentada é importante para entender os métodos para estimação de prevalências ajustadas e as rotinas computacionais desenvolvidas podem ser úteis para sua utilização.

TABELAS

Tabela 1 - Descrição das variáveis simuladas.

Variável	Tipo	Natureza	Valores Possíveis
Y	Desfecho (diabetes)	Dicotômica	0=Ausente 1=Presente
Sexo	Exposição	Dicotômico	0=Masculino 1=Feminino
Idade	Confundidor	Quantitativo	35 a 75 anos

Tabela 2 - Estimativas de prevalências (IC 95%) por sexo, ajustados pela idade, obtidas pelos métodos de predição condicional e marginal.

Programa	Marginal		Condicional	
	Feminino	Masculino	Feminino	Masculino
SUDAAN	0,2627 (0,241-0,286)	0,2240 (0,187-0,265)	0,2358 (0,212-0,261)	0,1961 (0,162-0,236)
%ADJ_PROP (SAS)	0,2678 (0,239-0,298)	0,2228 (0,184-0,268)	0,2358 (0,21 -0,262)	0,1961 (0,162-0,236)
SAS*	0,2627 (0,241-0,284)	0,2240 (0,185-0,263)	0,2358 (0,210-0,261)	0,1961 (0,159-0,233)
R*	0,2627 (0,241-0,284)	0,2240 (0,185-0,263)	0,2358 (0,202-0,269)	0,1961 (0,150-0,242)
STATA	-	-	0,2358 (0,210-0,261)	0,1961 (0,159-0,233)

Estimativas dos coeficientes (EP) do modelo de regressão logística:

Intercepto= -10,8999 (0,7797), $\hat{\beta}_{sexo}$ =0,2349 (0,1464), $\hat{\beta}_{idade}$ =0,1646 (0,0138)

EP - Erro padrão

* Intervalos de confiança estimados pelo método Delta

Tabela 3 - Comparação dos métodos de predição condicional e marginal

Método	Proporção mulheres	Prevalência ajustada mulheres	Proporção homens	Prevalência ajustada homens	Média ponderada*	Prevalência global (%)
Condicional	0,682	23,58	0,318	19,61	22,32	25,35
Marginal	0,682	26,27	0,318	22,40	25,04	

* Média ponderada = (Proporção de mulheres x Prevalência ajustada em mulheres) + (Proporção homens x Prevalência ajustada em homens).

REFERÊNCIAS

1. An R. Diabetes diagnosis and nutrition facts label use among US adults, 2005–2010. *Public Health Nutr.* 2015 outubro 20: 1-8.
2. Chambless LE, Heiss G, Folsom AR, Rosamond W, Szklo M, Sharrett AR, et al. Association of coronary heart disease incidence with carotid arterial wall thickness and major risk factors: the Atherosclerosis Risk in Communities (ARIC) Study, 1987-1993. *Am J Epidemiol.* 1997 setembro 15; 146(6): 483-94.
3. Lebedeva ER, Kobzeva NR, Gilev D, Olesen J. Prevalence of primary headache disorders diagnosed according to ICHD-3 beta in three different social groups. *Cephalalgia.* 2015 outubro 6.
4. Nambi V, Hoogeveen RC, Chambless L, Hu Y, Bang H, Coresh J, et al. Lipoprotein-associated phospholipase A2 and high-sensitivity C-reactive protein improve the stratification of ischemic stroke risk in the Atherosclerosis Risk in Communities (ARIC) study. *Stroke.* 2009 fevereiro; 40(2): 376-81.
5. Shon H-S, Choi HY, Kim JR, Ryu SY, Lee Y-J, Lee MJ, et al. Comparison and analysis of the prevalence of hepatitis C virus infection by region in the Republic of Korea during 2005-2012. *Clin Mol Hepatol.* 2015 setembro; 21(3): 249-56.
6. Smith AK, Cenzer IS, John Boscardin W, Ritchie CS, Wallhagen ML, Covinsky KE. Increase in Disability Prevalence Before Hip Fracture. *J Am Geriatr Soc.* 2015 outubro; 63(10): 2029-35.
7. Lane PW, Nelder JA. Analysis of Covariance and Standardization as Instances of Prediction. *Biometrics.* 1982; 38(3):613-21.

8. Szklo M, Nieto FJ. Epidemiology Beyond the Basics. 2th ed. 2007.
9. Schoenbach VJ, Rosamond WD. Understanding the Fundamentals of Epidemiology - an envolving text. 2000, Fall Edition.
10. Woodward M. Epidemiology Study Design and Data Analysis. 2014, Third Edition.
11. Wilcosky TC, Chambless LE. A comparision of Direct Adjustment and Regression Adjustment of Epidemiologic Measures. J Chron Dis. 1985; 38: 849-856.
12. McCullagh P, Nelder JA. Generalized Linear Models. 2th ed. 1989.
13. Research Triangle Institute . SUDAAN Language Manual, Volumes 1 and 2, Release 11. 2012.
14. Xu J, Long JS. Using the Delta Method to Construct Confidence Intervals for Predicted Probabilities, Rates, and Discrete Changes Indiana University. 2005 agosto 22.
15. Zhao D. Logistic Regression Adjustment of Proportions and its Macro Procedure
Disponível de:
<http://www2.sas.com/proceedings/sugi22/POSTERS/PAPER227.PDF>.
Acessado 7 de setembro de 2015.
16. Williams R. Using margins command. The Stata Journal. 2012; 12(2): 308-331.

17. Leeper T. Margins. Disponível de: <https://github.com/leeper/margins>. Acessado setembro 20, 2014.

SUPLEMENTOS

Suplemento I - Rotina SAS para utilizar a macro %ADJ_PROP

```
options ls=120;
libname L1 'C:\Users\NataliaBarbieri\Documents\ - UFRGS - Pos
Graduacao\Projeto Dissertacao\Simulacao\Banco simulado dissertacao';

data SIMULADO;
    set L1.base_simulada_201115;
run;
proc sort data=SIMULADO; by SEXO; run;

%include 'C:\Users\NataliaBarbieri\Documents\ - UFRGS - Pos
Graduacao\Projeto
Dissertacao\Simulacao\Sintaxes_Dissertacao\ADJ_PROP_Dingyi_Zhao.sas';

%ADJ_PROP(SIMULADO, LOGISTIC, , Y, SEXO, IDADE, , , YES, DES, OUTPROP);
```

Suplemento II - Rotina SAS para estimar prevalências ajustadas

```
%MACRO adjmodel(database,desfecho,exposicao,confundidor);

title1 "Modelo logístico";
proc logistic data=&database descending out=Betas covout;
    model &desfecho = &exposicao &confundidor / rl;
run;

/* Matriz de covariancias*/
data Cov;
    set Betas;
    if _TYPE_ = "COV";
    keep Intercept &exposicao &confundidor;
run;

data Betas1;
    set Betas;
    keep Intercept &exposicao &confundidor;
    rename &exposicao = Beta_&exposicao;
    rename &confundidor = Beta_&confundidor;
run;

data &database;
    set &database;
    if _n_ = 1 then set betas1;

    XBETA = Intercept + Beta_&exposicao*&exposicao +
            Beta_&confundidor*&confundidor;
    P = logistic(XBETA);
run;

proc means data=&database mean noprint;
    var &confundidor;
    output out=mean&confundidor mean=mean&confundidor;
run;

data mean&confundidor;
    set mean&confundidor;
    keep mean&confundidor;
run;

data &database;
    set &database;
    if _n_ = 1 then set mean&confundidor;
run;

data &database;
    set &database;

* Marginal - PREDMARG results;
    P_Marg = 1/(1 + exp(-(Intercept +Beta_&exposicao*&exposicao
    + Beta_&confundidor*&confundidor)));

*Proporcao em expostos = mulheres;
    PE_Marg = 1/(1 + exp(-(Intercept + Beta_&exposicao*1 +
    Beta_&confundidor*&confundidor)));

    label PE_Marg = "Predicted Marginal (Y=1|Tratamento)";
*PE - Proporcao em expostos;
```

```

PE_m=PE_Marg*(1- PE_Marg);
    PE_m0=PE_m*1; *b0;
    PE_m1=PE_m*1; *b1;
    PE_m2=PE_m*&confundidor; *b2;

*Proporcao em nao expostos;
PNE_Marg = 1/(1 + exp(-(Intercept + Beta_&exposicao*0 +
Beta_&confundidor*&confundidor)));
label PNE_Marg = "Predicted Marginal (Y=1|Controle)";

*p(x)*(1-p(x)); *PNE - Proporcao em nao expostos;
PNE_m=PNE_Marg*(1- PNE_Marg);
    PNE_m0=PNE_m*1;
    PNE_m1=PNE_m*0;
    PNE_m2=PNE_m*&confundidor;

* Condicional - CONDMARG results;
P_Cond = 1/(1 + exp(-(Intercept + Beta_&exposicao*&exposicao +
Beta_&confundidor*mean&confundidor)));

*Proporcao em expostos = mulheres;
PE_Cond = 1/(1 + exp(-(Intercept + Beta_&exposicao*1 +
Beta_&confundidor*mean&confundidor)));
label PE_COND = "Conditional Marginal (Y=1|Tratamento)";

PE_c=PE_Cond*(1- PE_COND); *PE - Proporcao em expostos;
    PE_c0=PE_c*1; *b0;
    PE_c1=PE_c*1; *b1;
    PE_c2=PE_c*&confundidor; *b2;

*Proporcao em nao expostos = homens;
PNE_Cond = 1/(1 + exp(-(Intercept + Beta_&exposicao*0 +
Beta_&confundidor*mean&confundidor)));
label PNE_COND = "Conditional Marginal (Y=1|Controle)";

*p(x)*(1-p(x)); *PNE - Proporcao em nao expostos;
PNE_c=PNE_Cond*(1- PNE_COND);
    PNE_c0=PNE_c*1;
    PNE_c1=PNE_c*0;
    PNE_c2=PNE_c*&confundidor;

run;

* Proporcoes ajustadas, por grupo;
proc means data=&database mean; output out=adj_prev mean=Pm
    P_Margm PE_Margm PNE_Margm PE_m0m PE_m1m PE_m2m PNE_m0m
    PNE_m1m PNE_m2m          PE_CONDM PNE_CONDM PE_c0m PE_c1m PE_c2m
    PNE_c0m PNE_c1m PNE_c2m;
var P P_Marg PE_Marg PNE_Marg PE_m0 PE_m1 PE_m2 PNE_m0
    PNE_m1 PNE_m2 PE_COND PNE_COND PE_c0 PE_c1 PE_c2 PNE_c0
    PNE_c1 PNE_c2;
run;

/*Marginal*/
data delPE_marg;
    set adj_prev;
    keep PE_m0m PE_m1m PE_m2m;
run;
proc transpose data=delPE_marg out=delPE_marg_T ;run;

data delPE_marg_T;

```

```

        set delPE_marg_T;
        drop _NAME_;
run;

data delPNE_marg;
    set adj_prev;
    keep PNE_m0m PNE_m1m PNE_m2m;
run;
proc transpose data=delPNE_marg out=delPNE_marg_T ;run;

data delPNE_marg_T;
    set delPNE_marg_T;
    drop _NAME_;
run;

/*Condicional*/
data delPE_cond;
    set adj_prev;
    keep PE_c0m PE_c1m PE_c2m;
run;
proc transpose data=delPE_cond out=delPE_cond_T ;run;

data delPE_cond_T;
    set delPE_cond_T;
    drop _NAME_;
run;

data delPNE_cond;
    set adj_prev;
    keep PNE_c0m PNE_c1m PNE_c2m;
run;
proc transpose data=delPNE_cond out=delPNE_cond_T ;run;

data delPNE_cond_T;
    set delPNE_cond_T;
    drop _NAME_;
run;

proc iml;
/*    Marginal*/
    use delPE_marg_T;
    read  all into delPE_marg_T;
    close delPE_marg_T;

    use delPNE_marg_T;
    read all into delPNE_marg_T;
    close delPNE_marg_T;

    use COV;
    read  all into COV;
    close COV;

/*    Condicional*/
    use delPE_cond_T;
    read  all into delPE_cond_T;
    close delPE_cond_T;

    use delPNE_cond_T;
    read all into delPNE_cond_T;
    close delPNE_cond_T;

    use COV;
    read  all into COV;

```

```

close COV;

/* Marginal*/
VarPE_Margm = t(delPE_marg_T)*COV*delPE_marg_T;
VarPNE_Margm = t(delpNE_marg_T)*COV*delpNE_marg_T;

del_m=delPNE_marg_T||delPE_marg_T;
var_PNE_Marg = t(del_m)*COV*del_m;

print del_m;
print var_PNE_Marg;

print delPE_marg_T delpNE_marg_T COV;
print VarPE_Margm VarPNE_Margm;

VarName1 = ("VarPE_Margm");
create VarPE_Margm from VarPE_Margm [C=VarName1];
append from VarPE_Margm;
close VarPE_Margm;

VarName2 = ("VarPNE_Margm");
create VarPNE_Margm from VarPNE_Margm [C=VarName2];
append from VarPNE_Margm;
close VarPNE_Margm;

/* Condicional*/
VarPE_Condm = t(delPE_cond_T)*COV*delPE_cond_T;
VarPNE_Condm = t(delpNE_cond_T)*COV*delpNE_cond_T;

del_c=delPNE_cond_T||delPE_cond_T;
var_PNE_Condm = t(del_c)*COV*del_c;

print del_c;
print var_PNE_Condm;

print delPE_cond_T delpNE_cond_T COV;
print VarPE_Condm VarPNE_Condm;

VarName1_c = ("VarPE_Condm");
create VarPE_Condm from VarPE_Condm [C=VarName1_c];
append from VarPE_Condm;
close VarPE_Condm;

VarName2_c = ("VarPNE_Condm");
create VarPNE_Condm from VarPNE_Condm [C=VarName2_c];
append from VarPNE_Condm;
close VarPNE_Condm;

run;
quit;

/* Marginal*/
data Adj_prev_marg;
set Adj_prev;
if _N_=1 then set VarPE_Margm;
run;

data Adj_prev_marg;
set Adj_prev_marg;

if _N_=1 then set VarPNE_Margm;

*IC PE_Marg - tratamento;
LI_PE_Margm=PE_Margm-1.96*sqrt(VarPE_Margm);

```

```

LS_PE_Margm=PE_Margm+1.96*sqrt(VarPE_Margm);
label LI_PE_Margm = "Limite inferior";
label LS_PE_Margm = "Limite superior";

*IC PNE_Marg - controle;
LI_PNE_Margm=PNE_Margm-1.96*sqrt(VarPNE_Margm);
LS_PNE_Margm=PNE_Margm+1.96*sqrt(VarPNE_Margm);
label LI_PNE_Margm = "Limite inferior";
label LS_PNE_Margm = "Limite superior";

run;

data Adj_prev_cond;
set Adj_prev;
if _N_=1 then set VarPE_Condm;
run;
/* Condicional*/
data Adj_prev_cond;
set Adj_prev_cond;

if _N_=1 then set VarPNE_Condm;

*IC PE_Marg - tratamento;
LI_PE_Condm=PE_Condm-1.96*sqrt(VarPE_Condm);
LS_PE_Condm=PE_Condm+1.96*sqrt(VarPE_Condm);
label LI_PE_Condm = "Limite inferior";
label LS_PE_Condm = "Limite superior";

*IC PNE_Marg - controle;
LI_PNE_Condm=PNE_Condm-1.96*sqrt(VarPNE_Condm);
LS_PNE_Condm=PNE_Condm+1.96*sqrt(VarPNE_Condm);
label LI_PNE_Condm = "Limite inferior";
label LS_PNE_Condm = "Limite superior";

run;

/* Marginal*/
proc print data= Adj_prev_marg;
title1 "Prevalencia ajustada no grupo exposto (IC 95%) -
metodo marginal";
var PE_Margm LI_PE_Margm LS_PE_Margm;
run;

proc print data= Adj_prev_marg;
title1 "Prevalencia ajustada no grupo não exposto (IC 95%) -
metodo marginal";
var PNE_Margm LI_PNE_Margm LS_PNE_Margm;
run;

/* Condicional*/
proc print data= Adj_prev_cond;
title1 "Prevalencia ajustada no grupo exposto (IC 95%) -
metodo condicional";
var PE_Condm LI_PE_Condm LS_PE_Condm;

title1 "Prevalencia ajustada no grupo não exposto (IC 95%) -
metodo condicional";
var PNE_Condm LI_PNE_Condm LS_PNE_Condm;
run;

%MEND adjmodel;

```


Suplemento III - Rotina R para estimar prevalências ajustadas

```
# Function 'logistic_adj_model'
# Arguments - MARGINAL MODEL
# model : a logistic model object.
# groups : nominal variable name with group information.
# adjustvar : numeric variable name (adjustable).
# level : significance level for confidence intervals.
# method : model method - marginal or conditional.
#           Use "marginal" or "m" for a marginal model,
#           "conditional" or "c" for a conditional model.
logistic_adj_model <- function(model, groups, adjustvar, level = 0.95,
method = "marginal" ) {
  # handling errors
  if ( length(groups) != 1 ) stop("Please, set one nominal variable
name with group information.")
  if ( length(adjustvar) != 1 ) stop("Please, use a numeric variable
name (adjustable).")
  if ( ! grep(pattern = adjustvar, x = model$"formula") > 0 )
stop(paste("'",adjustvar,'" not present in the model formula.", sep =
""))
  if ( ! grep(pattern = groups, x = model$"formula") > 0 )
stop(paste("'",groups,'" not present in the model formula.", sep = ""))

  # covariance matrix
  varbetaAux = vcov(model)

  #X <- as.numeric(t(model$data[2]))
  # X data
  nX <- which(names(model$data) == adjustvar)
  X <- as.numeric(t(model$data[nX]))

  # Group levels
  nG <- which(names(model$data) == groups)
  GROUP <- as.numeric(names(table(model$data[nG])))

  COEFS <- model$"coefficients"
  meanX <- mean(X)

  if ( method == "marginal" || method == "m" ) {
    PE_Marg = 1/(1 + exp(-(COEFS[1] + COEFS[2]*GROUP[2] +
COEFS[3]*X)))
    PE_m=PE_Marg*(1- PE_Marg)
    PE_m0=PE_Marg*1
    PE_m1=PE_Marg*GROUP[2]
    PE_m2=PE_Marg*X

    PNE_Marg = 1/(1 + exp(-(COEFS[1] + COEFS[2]*GROUP[1] +
COEFS[3]*X)))
    PNE_m=PNE_Marg*(1- PNE_Marg)
    PNE_m0=PNE_Marg*1
    PNE_m1=PNE_Marg*GROUP[1]
    PNE_m2=PNE_Marg*X

    PE_Margm <- mean(PE_Marg)
    PE_m0m <- mean(PE_m0)
    PE_m1m <- mean(PE_m1)
    PE_m2m <- mean(PE_m2)
    delhbTT <- c(PE_m0m, PE_m1m, PE_m2m)

    PNE_Margm <- mean(PNE_Marg)
```

```

PNE_m0m <- mean(PNE_m0)
PNE_m1m <- mean(PNE_m1)
PNE_m2m <- mean(PNE_m2)
delhbCT <- c(PNE_m0m, PNE_m1m, PNE_m2m)

VarPE_Margm = t(delhbTT) %*% varbetaAux %*% t(t(delhbTT));
VarPNE_Margm = t(delhbCT) %*% varbetaAux %*% t(t(delhbCT));

# IC PE_Marg - tratamento; qnorm(c(.025,.975)) ;
LI_PE_Margm=PE_Margm+qnorm((1-level)/2)*sqrt(VarPE_Margm);
LS_PE_Margm=PE_Margm+qnorm(1-(1-level)/2)*sqrt(VarPE_Margm);

# IC PNE_Marg - controle;
LI_PNE_Margm=PNE_Margm+qnorm((1-
level)/2)*sqrt(VarPNE_Margm);
LS_PNE_Margm=PNE_Margm+qnorm(1-(1-
level)/2)*sqrt(VarPNE_Margm);

ans <- c(PE_Margm, LI_PE_Margm, LS_PE_Margm, PNE_Margm,
LI_PNE_Margm, LS_PNE_Margm)
names(ans) <- c("PE_Marg", "LI_PE_Marg", "LS_PE_Marg",
"PNE_Marg", "LI_PNE_Marg", "LS_PNE_Marg")
return(ans)
} else if ( method == "conditional" || method == "c" ) {
PE_Cond = 1/(1 + exp(-(COEFS[1] + COEFS[2]*GROUP[2] +
COEFS[3]*meanX)))
PE_c=PE_Cond*(1- PE_Cond)
PE_c0=PE_Cond*1
PE_c1=PE_Cond*GROUP[2]
PE_c2=PE_Cond*X

PNE_Cond = 1/(1 + exp(-(COEFS[1] + COEFS[2]*GROUP[1] +
COEFS[3]*meanX)))
PNE_c=PNE_Cond*(1- PNE_Cond)
PNE_c0=PNE_Cond*1
PNE_c1=PNE_Cond*GROUP[1]
PNE_c2=PNE_Cond*X

PE_Condm <- mean(PE_Cond)
PE_c0m <- mean(PE_c0)
PE_c1m <- mean(PE_c1)
PE_c2m <- mean(PE_c2)
delhbTT <- c(PE_c0m, PE_c1m, PE_c2m)

PNE_Condm <- mean(PNE_Cond)
PNE_c0m <- mean(PNE_c0)
PNE_c1m <- mean(PNE_c1)
PNE_c2m <- mean(PNE_c2)
delhbCT <- c(PNE_c0m, PNE_c1m, PNE_c2m)

VarPE_Condm = t(delhbTT) %*% varbetaAux %*% t(t(delhbTT));
VarPNE_Condm = t(delhbCT) %*% varbetaAux %*% t(t(delhbCT));

# IC PE_Cond - tratamento; qnorm(c(.025,.975)) ;
LI_PE_Condm=PE_Condm+qnorm((1-level)/2)*sqrt(VarPE_Condm);
LS_PE_Condm=PE_Condm+qnorm(1-(1-level)/2)*sqrt(VarPE_Condm);

# IC PNE_Cond - controle;
LI_PNE_Condm=PNE_Condm+qnorm((1-
level)/2)*sqrt(VarPNE_Condm);
LS_PNE_Condm=PNE_Condm+qnorm(1-(1-
level)/2)*sqrt(VarPNE_Condm);

```

```
      ans <- c(PE_Condm, LI_PE_Condm, LS_PE_Condm, PNE_Condm,
LI_PNE_Condm, LS_PNE_Condm)
      names(ans) <- c("PE_Cond", "LI_PE_Cond", "LS_PE_Cond",
"PNE_Cond", "LI_PNE_Cond", "LS_PNE_Cond")
      return(ans)

    } else {
      stop("Please, try a valid model method option.")
      return(ans)
    }
  }
}
```

Suplemento IV - Detalhes do exemplo simulado e execução das rotinas computacionais

O exemplo simulado mostra detalhes da geração dos dados e da execução das rotinas computacionais para a estimação de prevalências ajustadas. O contexto é de um estudo transversal em que o desfecho Y representa a presença ($Y = 1$) ou ausência ($Y = 0$) de diabetes. De forma intencional, a idade está associada com o desfecho e distribuição é diferente em cada sexo, isto é, a associação entre o desfecho e sexo está confundida pela idade. A categoria sexo feminino foi considerada como exposição e masculino como não exposição. Assim, é importante estimar prevalências do desfecho em cada sexo, ajustadas pela idade. Foi considerada uma amostra de $n = 2000$ indivíduos selecionados ao acaso da população em investigação, por meio de dados simulados.

Geração dos dados

Utilizando o programa SAS, versão 9.4, o conjunto de dados do exemplo foi gerado considerando a matriz de correlação entre as variáveis X_1 e X_2

$$\rho = \begin{bmatrix} 1 & 0,9 \\ 0,9 & 1 \end{bmatrix}$$

a partir da qual foram geradas 2.000 observações independentes, seguindo os passos:

- 1) Foi gerada a matriz Z com dimensão (2000×2) , a partir da distribuição normal multivariada com vetor de médias $(0,0)$, variâncias $(1,1)$ e matriz de correlações ρ ;

- 2) Foi definida a matriz $U (2000 \times 2)$ por meio da transformação $U_j(Z_j) = \Phi(Z_j)$, $j=1,2$, em que $\Phi(\cdot)$ é a função de distribuição da normal padrão. Desse modo, as colunas de U são variáveis aleatórias com distribuição Uniforme em $[0,1]$, porém não são independentes;
- 3) Foi criada a matriz $X (2000 \times 2)$, cujas colunas foram definidas por meio da transformação $X_j = F^{-1}(U_j)$, em que $F^{-1}(\cdot)$ é a inversa da função de distribuição especificada: para X_1 foi usada a distribuição Uniforme em $[0,1]$; para X_2 foi usada a distribuição normal com média 57 e desvio padrão 5;
- 4) Uma variável aleatória W , com distribuição uniforme em $[0,1]$ e independente de X_1 foi gerada para definir o preditor dicotômico $SEXO$, por meio da dicotomização de X_1 , tal que $SEXO=1$ (feminino) se $X_1 > (W - 0,2)$, e $SEXO=0$ (masculino) em caso contrário;
- 5) A variável $IDADE$ foi definida como o menor inteiro menor ou igual a X_2 ;
- 7) O preditor linear $XBETA = -12 + \ln(1,3) \times SEXO + \ln(1,2) \times IDADE$ foi utilizado para determinar a probabilidade de sucesso (presença do desfecho) definida por:

$$p(SEXO, IDADE) = \frac{1}{1 + e^{-XBETA}}.$$

Dessa forma, sexo ($RC = 1,3$) e idade ($RC = 1,2$) estão associadas com o desfecho (diabetes).

8) Finalmente, para definir a variável resposta Y (diabetes), para cada linha do banco de dados foi gerado um valor v com distribuição Uniforme $[0,1]$, usado da seguinte forma:

$$Y = \begin{cases} 1, v < p(SEXO, IDADE) \\ 0, v \geq p(SEXO, IDADE) \end{cases}$$

Em resumo, a variável Y representa o desfecho diabetes, sexo é a exposição e idade é um confundidor. A tabela abaixo descreve a amostra.

Descrição das variáveis na amostra.

Variável	n(%) ou média (DP*)
Desfecho (Y)	
Ausente	1493 (74,7)
Presente	507 (25,3)
Sexo	
Masculino	636 (31,8)
Feminino	1364 (68,2)
Idade (anos)	57,6 (5,0)

* DP = desvio padrão

A prevalência global do desfecho é de 25,3% (IC 95% 23,4 - 27,3). Em homens, a média (DP) de idade foi 54,2 (3,94) anos, enquanto que para mulheres foi 59,2 (4,60) anos. A seção 3.1 descreve diferentes métodos para estimar a prevalência do desfecho, por sexo, ajustada e não ajustada. As estimativas de prevalência do desfecho, não ajustadas, foram 13,8% (IC 95% 11,7 - 16,7) e 30,7% (IC 95% 28,3 - 33,2), respectivamente, para homens e mulheres. A discrepância observada nessas estimativas pode ser atribuída ao confundimento devido às diferenças nas distribuições da idade entre sexos. Estimativas livres de confundimento podem ser obtidas por meio das prevalências ajustadas pela idade, por meio dos métodos marginal e condicional. A

tabela abaixo apresenta as estimativas do modelo de regressão logística para os dados do exemplo simulado.

Estimativas dos parâmetros do modelo de regressão logística no exemplo simulado.

Variável	Coeficiente	Erro padrão	p-valor
Intercepto	-10,8899	0,7797	<0,0001
Sexo	0,2349	0,1464	0,1086
Idade	0,1646	0,0138	<0,0001

A seguir são apresentados os detalhes para obtenção de prevalências ajustadas, pelos métodos de predição marginal e condicional, utilizando as diferentes abordagens computacionais.

SUDAAN (SAS-Callable)

Neste trabalho foi utilizado o programa SUDAAN (SAS-Callable) versão 11, em conjunto com o programa SAS versão 9.4. A sintaxe abaixo é utilizada para obter as prevalências ajustadas e seus respectivos intervalos de confiança.

```
proc rlogist data=SIMULADO filetype=sas design=SRS;
  model Y = SEXO IDADE;
  setenv colspce=1 colwidth=12 decwidth=4;
  predmarg SEXO / SEXO=(0 1);
  condmarg SEXO / SEXO=(0 1);
run;
```

O modelo de regressão logística é ajustado por meio do comando *proc rlogist*, *model* especifica o modelo utilizado, em que *Y* é a variável desfecho (diabetes) e *SEXO* e *IDADE* são as variáveis preditoras. Os métodos marginal e condicional são obtidos por meio dos comandos *predmarg* e *condmarg*, respectivamente, especificando a variável para a qual se deseja estimar prevalências ajustadas e os valores das

categorias. A opção `design=SRS` especifica o delineamento amostral utilizando amostra aleatória simples (*SRS - Simple Random Sample*). Os resultados obtidos seguem abaixo:

Conditional Marginal	Conditional Marginal	SE	Lower 95% Limit	Upper 95% Limit	T:Marg=0	P-value
SEXO						
0	0.1961	0.0189	0.1617	0.2358	10.3922	0.0000
1	0.2358	0.0126	0.2119	0.2615	18.6646	0.0000

Predicted Marginal	Predicted Marginal	SE	Lower 95% Limit	Upper 95% Limit	T:Marg=0	P-value
SEXO						
0	0.2240	0.0199	0.1874	0.2654	11.2556	0.0000
1	0.2627	0.0114	0.2410	0.2857	23.0596	0.0000

SAS - Macro %ADJ_PROP

Zhao (1985) apresenta a versão original da macro `%ADJ_PROP`, mas foram necessárias pequenas alterações para sua execução no SAS versão 9.4. A chamada da macro segue abaixo:

```
%MACRO ADJ_PROP (INFILE, MODEL, BYVAR, YVARS, IX_VAR,
                 CONTVARS,GIVENVAR, GIVENVAL, PRINT, DES, OUT);
```

Os parâmetros especificados na chamada da macro são descritos por Zhao (1985) e também estão no início do código. De forma sucinta, são descritos a seguir: *INFILE* especifica o banco de dados utilizado, *MODEL* especifica se o modelo de regressão é linear ou logístico, *BYVAR* pode ser utilizado na existência de uma variável classificadora, *YVARS* é a variável dependente do modelo (desfecho), *IX_VAR* especifica as variáveis de exposição, *CONTVARS* define as variáveis de controle,

GIVENVAR especifica as variáveis de controle com valores de referência especificados, *GIVENVAL* define os valores de referência das variáveis especificadas na opção *CONTVARS* (se não forem especificados, utiliza a média das variáveis), *PRINT=YES* ativa a visualização dos resultados intermediários, *DES=DES* especifica o modelo logístico, assumindo que *Y=1* representa presença do desfecho de interesse, e *OUT* é o nome do banco de dados que será criado com os resultados gerados. Para os casos em que existirem mais de uma variável ou valores, os mesmos devem ser separados por espaços em branco.

Abaixo é apresentada a sintaxe para executar a macro para os dados do exemplo simulado:

```
%ADJ_PROP (SIMULADO, LOGISTIC, , Y, SEXO, IDADE, , , YES, DES, OUTPROP) ;
```

Os resultados obtidos são mostrados abaixo, sendo que *ADJVALUE* refere-se à prevalência ajustada obtida por meio do método de predição condicional e *CILOW* e *CIUP* são os limites inferior e superior do intervalo de confiança, respectivamente. *ADJUSTED* refere-se à prevalência ajustada obtida por meio do método de predição marginal, e *LOWCI* e *UPCI* são os limites inferior e superior do intervalo de confiança, respectivamente.

SEXO	CILOW	CIUP	ADJVALUE	XBETA	XNUMBER	NPREDIT	SUMY	OVERALL	K	ADJUSTED	LOWCI	UPCI
0	0.16172	0.23583	0.19614	-1.41059	636	124.745	507	446.423	1.13569	0.22276	0.18366	0.26783
1	0.21124	0.26234	0.23583	-1.17566	1364	321.678	507	446.423	1.13569	0.26784	0.23990	0.29794

SAS - Macro %adjmodel

A rotina SAS (macro %adjmodel) desenvolvida para obter as prevalências ajustadas é apresentada no Suplemento II. A chamada da macro é realizada pela

```
%adjmodel(database=SIMULADO, desfecho=Y, exposicao=SEXO,  
confundidor=IDADE);
```

sintaxe abaixo:

Os resultados são estimativas pontuais e respectivos intervalos de confiança. Assim, PNE_Marg é a estimativa da prevalência ajustada do desfecho em não expostos (homens), obtidos por meio do método marginal, e LI_PNE_Marg e LS_PNE_Marg são os limites inferior e superior do intervalo de confiança, respectivamente. Similarmente, PE_Marg, LI_PE_Marg e LS_PE_Marg são as estimativas de prevalência ajustada em expostos (mulheres), e os limites inferior e superior do intervalo de confiança, respectivamente.

PNE_Marg	LI_PNE_Marg	LS_PNE_Marg
0.22401	0.18486	0.26317
PE_Marg	LI_PE_Marg	LS_PE_Marg
0.26272	0.24112	0.28433

De forma similar, PNE_Cond é a estimativa pontual da prevalência ajustada do desfecho em não expostos (homens), obtidos por meio do método condicional, e LI_PNE_Cond e LS_PNE_Cond são os limites inferior e superior do intervalo de confiança, respectivamente. Do mesmo modo, PE_Cond, LI_PE_Cond e LS_PE_Cond são as estimativas de prevalência ajustada em expostos (mulheres), e os limites inferior e superior do intervalo de confiança, respectivamente.

PNE_COND	LI_PNE_Cond	LS_PNE_Cond
0.19614	0.15911	0.23318
PE_COND	LI_PE_Cond	LS_PE_Cond
0.23583	0.21028	0.26139

Stata - margins

No programa Stata (versão 13), é necessário primeiramente ajustar o modelo de regressão logística, usando o comando *logit*. O comando *margins*, em conjunto com a opção *atmeans*, estima as prevalências ajustadas por meio do método condicional para cada categoria da exposição (sexo), ajustado pela idade.

```
logit y i.sexo idade, nolog
margins sexo, atmeans
```

Os resultados do comando *margins* seguem abaixo:

Margin	Std. Err.	z	P> z	Delta-method [95% Conf. Interval]	

sexo					
0	.1961129	.0188946	10.38	0.000	.1590802 .2331456
1	.2358233	.0130401	18.08	0.000	.2102651 .2613815

R – Função *adjmodel*

O suplemento III apresenta a função *adjmodel* com os mesmos métodos da macro %*adjmodel* (SAS). O ajuste do modelo de regressão logística no programa R é realizado por meio da função *glm*, como descrito abaixo:

```
model<-glm(Y ~ SEXO + IDADE, data = simulado, family=binomial)
```

A obtenção das prevalências ajustadas pelo método condicional e marginal ocorre por meio da função *adjmodel*, em que *model* é o modelo de regressão logística

(que deve ser especificado a priori), *groups* é a variável preditora (dicotômica), *adjustvar* é a variável de ajuste (quantitativa), e *method* é o método de predição para estimar prevalências ajustadas (marginal ou condicional), conforme segue:

```
adjmodel(model, groups = "SEXO", adjustvar = "IDADE", level=0.95, method = "marginal" )
```

```
adjmodel(model, groups = "SEXO", adjustvar = "IDADE", level=0.95, method = "condicional" )
```

Os resultados obtidos são mostrados abaixo, utilizando a mesma nomenclatura da rotina SAS:

PE_Marg	LI_PE_Marg	LS_PE_Marg	PNE_Marg	LI_PNE_Marg	LS_PNE_Marg
0.2627174	0.2308656	0.2945693	0.2239920	0.1654427	0.2825412

PE_Cond	LI_PE_Cond	LS_PE_Cond	PNE_Cond	LI_PNE_Cond	LS_PNE_Cond
0.2358233	0.2023784	0.2692681	0.1961129	0.1500472	0.2421786

7. CONCLUSÕES E CONSIDERAÇÕES FINAIS

Prevalências ajustadas são frequentemente reportadas em estudos epidemiológicos e são importantes para o leitor fazer uma análise do fenômeno livre de confundimento. Atualmente os métodos de estimação marginal e condicional são as abordagens mais utilizadas.

As rotinas SAS e R desenvolvidas no trabalho podem ser úteis para estimar prevalências ajustadas, produzindo resultados muito similares aos do SUDAAN. Em especial, a versão para o R é uma alternativa aos *softwares* comerciais.

As estimativas pontuais das rotinas desenvolvidas para o SAS e R são idênticas aos do SUDAAN, e as discrepâncias entre as estimativas por intervalo não são relevantes na prática.

No método condicional, a média ponderada (proporcional às frações amostrais) das prevalências ajustadas para cada categoria pode ser diferente da proporção global observada. No exemplo, esta média ponderada foi 22,3% ($68,2 \times 23,58 + 31,8 \times 19,61$), muito distante da prevalência global 25,3%. Esta é uma desvantagem do método condicional, que muitas vezes pode confundir o leitor. O método marginal não é suscetível a este problema, exceto talvez por arredondamentos no cálculo da média ponderada, que no exemplo foi igual a 25,0% ($68,2 \times 26,27 + 31,8 \times 22,40$).

Os resultados relativamente discrepantes gerados macro %ADJ_PROP pelo método marginal podem ser explicados pela diferença dos métodos utilizados. O método condicional utiliza uma constante de aproximação k , produzindo resultados similares aos do SUDAAN. Outra desvantagem é a falta de documentação dos códigos, dificultando sua compreensão para descrição dos métodos ou eventuais modificações.

As rotinas computacionais ainda precisam ser estendidas para possibilitar a incorporação de número maior de variáveis de exposição e de controle, de forma

automática. Outra possibilidade é estender o estudo e as rotinas para o modelo log-binomial, comparando as estimativas de prevalências ajustadas com aquelas obtidas pelo modelo logístico.

Em suma, prevalências ajustadas são importantes para uma análise mais realista do fenômeno em estudo. A discussão apresentada é importante para entender os métodos para estimação de prevalências ajustadas, e as rotinas computacionais desenvolvidas podem ser úteis para sua utilização.