

Estimando valor de tempo de viagem com diferentes fontes de dados utilizando modelos logit

[Estimating travel time value by different data source using the logit model]

Francisco Gildemir Ferreira da Silva*, Sergio Aquino DeSouza

Agência Nacional de Transportes Terrestres, Universidade Federal do Ceará/CAEN

Submitted 3 Sep 2012; received in revised form 9 Jan 2013; accepted 20 Jan 2013

Resumo

A estimação de parâmetros de valor de tempo de viagem – VTV – pode ser feita utilizando metodologias diversas divididas em duas vertentes de especificação funcional a de equações estruturadas com base na teoria econômica e a essencialmente estatística com múltiplas equações. A primeira abordagem é a adotada neste trabalho, pois ela apresenta maior racionalidade e durante a evolução dos modelos de estimação do VTV, muitas dúvidas surgiram dos métodos de estimativas principalmente da fidedignidade dos parâmetros estimados relação ao parâmetro populacional. Dentre os métodos de estimação um que se propõe a melhorar os resultados é o de juntar fontes de dados distintas, mas que não foi utilizado para ver a melhoria na estimação do VTV. Neste sentido, este trabalho apresenta um estudo relativo a estimativas de modelos de escolha discreta utilizando uma combinação de duas fontes de dados e calculando os valores de tempo para diferentes modelos estruturais, conforme a teoria do valor de tempo. O trabalho amplia a abordagem de Morikawa (1989) formalizando as hipóteses necessárias a uma estimação com duas fontes de dados provenientes da pesquisa do Plano Diretor de Operação de Transporte do Ceará e indica que a combinação das fontes de dados de preferência declarada e revelada feita naquele estado gera resultados diferentes do modelo clássico logit e que isso dependerá da fonte geradora de dados.

Palavras-Chave: valor de tempo de viagem, fontes de dados, transporte de passageiros, modelos de escolha discreta.

Abstract

The estimation of travel time value - TTV - parameter can be done using some methodologies that can be match into two functional specification background: the first is the structured economic theory approach and the second is the statistical with multiple equations. The first approach is adopted in this study because it has a higher rationality and during the evolution of estimation models of TTV, many doubts arose related to the reliability of the estimation methods comparing the estimated parameter with the population parameter. Among the estimation methods one that aims to improve outcomes is to bring together different data sources, but that was not used to see improvement in the estimation of the TTV. This paper presents a study on the estimation of discrete choice models using a combination of two sources of data and calculating values of travel time for different structural models, according to the theory of time value. The work expands Morikawa (1989) approach formalizing the assumptions necessary for a estimation with two data sources using the data base from Ceará and it indicates that good data source combination can depend of the data source.

Key words: travel time value, data sources, passengers transport, discrete choice model.

* Email: gildemir@gmail.com.

Recommended Citation

DaSilva, F. G. F. and DeSouza, S. A. (2013) Estimando valor de tempo de viagem com diferentes fontes de dados utilizando modelos logit. Journal of Transport Literature, vol. 7, n. 4, pp. 107-129.

■ JTL|RELIT is a fully electronic, peer-reviewed, open access, international journal focused on emerging transport markets and published by BPTS - Brazilian Transport Planning Society. Website www.transport-literature.org. ISSN 2238-1031.

This paper is downloadable at www.transport-literature.org/open-access.

Introdução

A estimação de parâmetros de valor de tempo de viagem – VTV – são controversos, por tratar-se de um parâmetro subjetivo teoricamente representado. Uma vez estimado, enseja-se que os estimadores tenham no mínimo três propriedades: ausência de viés, eficiência e consistência. A ausência de viés implica na esperança dos estimadores serem igual ao parâmetro populacional. A eficiência implica na variância do estimador ser menor possível e ser tão próximo quanto possível do limite de Cramer-Rao (Train, 2003). A consistência é uma propriedade assintótica dos estimadores que implica na convergência em probabilidade do estimador para o parâmetro populacional.

Infelizmente, as propriedades somente serão atendidas se os dados forem robustos, propiciando o alcance assintótico das propriedades acima descritas. Isso não acontece por conta de erros de medidas mínimos que, às vezes, são inerentes a forma de aquisição dos dados. Os problemas inerentes aos dados podem ser solucionados, de tal sorte que as propriedades desejadas dos estimadores sejam atendidas total ou parcialmente. Utilizar todos os dados, mesmo com os problemas de medida indesejadas é válido, pois eles possuem também informações relevantes referentes ao processo de decisão dos indivíduos, ou seja, deve-se tentar de todas as formas extrair as informações e expurgar os erros inerentes aos dados.

A proposta de tentar retirar o melhor de cada fonte de dados é válida ao resultar em melhores estimadores favorecendo decisões públicas mais acertadas e ao propiciando simulações mais precisas. Neste sentido, este trabalho apresenta um estudo relativo a estimações de modelos de escolha modal utilizando uma combinação de duas fontes de dados tal como apresentado em Morikawa (1989) e Hensher, Louviere e Swait (2000) e calculando os valores de tempo para diferentes modelos estruturais, conforme a teoria do valor de tempo apresentado em DaSilva (2012). O resultado é que a combinação de fontes de dados pode gerar informações mais acertadas, contudo, deve-se avaliar a robustez e compatibilidade entre as fontes geradoras de dados. Além disso, o trabalho amplia a abordagem de Morikawa (1989) formalizando as hipóteses necessárias a uma estimação com duas fontes de dados e descrevendo os problemas inerentes as fontes de dados.

Para apresentar o estudado, o trabalho inicia indicando as características de dados de “Preferência Declarada – PD” e “Preferência Revelada – PR”, as formas de estimar utilizando fontes de dados distintas. Em seguida, estimam-se os modelos utilizando dados pesquisados no Ceará e faz-se uma comparação dos modelos com e sem enriquecimento de dados e finaliza-se com a conclusão.

1. Modelo Logit: histórico e teoria

Nesta seção será apresentado o histórico dos modelos de maximização de utilidade aleatória e formalizará o modelo Logit binário.

1.1 Surgimento, evolução e atualidade

Em McFadden (2000a) e McFadden (2000b) é feita uma revisão histórica dos modelos de escolha discreta, caracterizando a teoria dos seus primórdios a atualidade da agenda de pesquisa, resgatando as hipóteses básicas e sua vantagem quando associada a teoria econômico, principalmente por adequação da abordagem dada por Lancaster (1966).

Inicialmente a hipótese de distribuição IID foi relaxada e conseqüentemente reduziu-se o problema de IIA, e na possibilidade de uma matriz de variância e covariância que acomodasse heterogeneidade entre indivíduos ou entre grupos de indivíduos. Isso propiciou o surgimento dos modelos: Multinomial Logit, Nested Logit e Mixed Logit (Train, 2003).

Nos últimos anos, muitos trabalhos foram desenvolvidos, principalmente em pesquisas de demanda por viagem com o uso de dados desagregados tais como: Morikawa (1989), Ben-Akiva & Morikawa (1990), Morikawa, Ben-Akiva, & Yamada (1991), Hensher & Bradley (1993), Hensher, Louviere, & Swait (1999) e Brownstone & Train (1999). Entretanto, o objetivo do trabalho aqui não será ir além do que o modelo Logit pode auxiliar na mensuração do valor de tempo, assim, seguir-se-a o capítulo apresentado a formalização do modelo Logit e o processo de estimação dele.

1.2 Formalização do Modelo Logit

A função de utilidade que representa o grau de preferência de uma alternativa para certo indivíduo é definida em função dos valores dos atributos das alternativas e das características socioeconômicas do indivíduo, tal como na Equação 1.

$$U_{in} = U(z_{in}, S_n) \quad (1)$$

Em que: U_{in} : Utilidade de uma alternativa i para um indivíduo n ;

z_{in} : Vetor dos atributos relevantes da alternativa i ;

S_n : Vetor das características socioeconômicas do indivíduo n ;

Dado que a utilidade não pode ser prevista com total certeza, seguindo Ben-Akiva e Lerman (1985) trata-se a utilidade como uma variável aleatória, formada por uma componente determinística, também chamada de sistemática, e outra aleatória, que reflete as “irracionalidades” da escolha do indivíduo. Essa abordagem provem dos trabalhos de psicologia, inicialmente propostos por Thurstone em 1927 sobre utilidades aleatórias, presumindo que os indivíduos não escolhem com 100% de certeza as coisas, mas com uma idiosincrasia presente no termo erro. Dessa forma, a utilidade de uma alternativa i para um indivíduo n (U_{in}) pode ser representada pela seguinte expressão:

$$U_{in} = V_{in} + \varepsilon_{in} \quad (2)$$

Com: U_{in} : Utilidade global de uma alternativa i para um indivíduo n ;

V_{in} : Componente sistemática da utilidade de uma alternativa i para um indivíduo n ;

ε_{in} : Componente aleatória da utilidade de uma alternativa i para um indivíduo n .

A escolha de uma alternativa i em relação a outra j , se dá pela comparação entre duas utilidades tal que $U_{in} > U_{jn}, \forall i \neq j$. Como a utilidade possui um componente determinístico e outro aleatório, a escolha do indivíduo, conforme a probabilidade de escolha se dará conforme satisfeita a inequação acima relatada e atendendo a seguinte probabilidade:

$$\begin{aligned}
\Pr(y = j | \mathbf{z}, \mathbf{x}) &= \Pr[U_{in} > U_{jn} \forall i \neq j | \mathbf{z}, \mathbf{x}] \\
&= \Pr[V_{in} + \varepsilon_{in} > V_{jn} + \varepsilon_{jn} \forall i \neq j | \mathbf{z}, \mathbf{x}] \\
&= \Pr[\varepsilon_{in} - \varepsilon_{jn} > V_{jn} - V_{in} \forall i \neq j | \mathbf{z}, \mathbf{x}] \\
&= \Pr[\tilde{\varepsilon} > -\tilde{V} \forall i \neq j | \mathbf{z}, \mathbf{x}]
\end{aligned}$$

Assumindo a parte determinística como uma forma linear descrita na Equação 3 e a distribuição de probabilidade dos termos aleatórios assumido de valor extremo foi demonstrado por Luce (1959) que a probabilidade de escolha é igual a Equação 4.

$$V_{in} = \beta_0 + \beta_1 x_{in1} + \beta_2 x_{in2} + \beta_3 x_{in3} + \dots + \beta_k x_{ink} \quad (3)$$

Onde: x_{ink} : Atributo k da alternativa i para o indivíduo n;

β_0 : Constante específica da Alternativa e representa o efeito de escolha da alternativa que não está incluída nos efeitos dos atributos;

β_k : Peso relativo do atributo x_{ink} na composição da função utilidade.

$$P_n(i) = \frac{e^{V_{in}}}{\sum_{j \in A(h)} e^{V_{jn}}} \quad (4)$$

em que: $P_n(i)$: Probabilidade de escolha de uma alternativa i por um indivíduo n;

$A(n)$: Conjunto de alternativas j disponíveis para o indivíduo n;

V_{in} : Utilidade sistemática de uma alternativa i para um indivíduo n;

V_{jn} : Utilidade sistemática de uma alternativa j para um indivíduo n;

A estimação do modelo na forma binária é feita com a maximização da função 5. Considerando uma amostra com N observações e definindo $y_{in} = 1$ (se, na observação n, foi escolhida a alternativa i) ou $y_{in} = 0$ (caso contrário).

$$LL = \sum_{n=1}^N (y_{in}) \log(P_{ni}) + (1 - y_{in}) \log(1 - P_{ni}) \quad (5)$$

Onde: LL : logaritmo da função de verossimilhança;

$P_n(i)$: Probabilidade do indivíduo n escolher a alternativa i , expressa em (4);

N : tamanho da amostra.

As estatísticas para avaliar a qualidade de estimação dos modelos da família Logit, conforme Ben-Akiva & Lerman (1985), são os seguintes:

- $L(0)$: valor da função logarítmica de verossimilhança quando todos os parâmetros são zero;
- $L(c)$: é o valor da função logarítmica de verossimilhança quando somente a constante específica da alternativa é incluída. Isto corresponde ao caso onde a probabilidade de escolha é função apenas da fração de amostra que escolheu a determinada alternativa;
- $L(\beta)$: valor máximo da função logarítmica de verossimilhança;
- $-2(L(0) - L(\beta))$: estatística utilizada para testar a hipótese de que todos os parâmetros são zero; é assintoticamente distribuída como χ^2 com k graus de liberdade, onde k é igual ao número de parâmetros estimados;
- $-2(L(c) - L(\beta))$: estatística utilizada para testar a hipótese nula de que todos os parâmetros são zero; é assintoticamente distribuída como χ^2 com $k - J + 1$ graus de liberdade, onde J é o número de alternativas;
- ρ_{zero}^2 : mede a fração do valor de verossimilhança explicado pelo modelo, definido como $1 - \left(\frac{L(\beta)}{L(0)} \right)$. Os valores de ρ_{zero}^2 dependerão do tipo de modelo a ser construído.

Essa medida é mais adequada na comparação de duas especificações desenvolvidas com o mesmo conjunto de dados.

2. Fontes de dados para modelos de escolha discreta

Qualquer que seja a forma de modelar a escolha dos usuários de transporte, a sua concepção é fundamentada, principalmente, no planejamento de coleta de dados de preferência declarada (PD) e/ou preferência revelada (PR). Os dados de PR representam as escolhas já realizadas por usuários, portanto, representam o comportamento real de escolha. Por outro lado, os dados de PD fornecem informações a respeito da escolha do indivíduo diante de um conjunto de alternativas previamente definidas, com no mínimo uma alternativa hipotéticas. O viés da hipótese é uma crítica ao método, pois a hipótese de opções pode não ser realista, entretanto, mesmo que a hipótese não tenha viés, Morikawa (1989) discute a estabilidade e a validade dos dados para modelagem e divide os problemas dos dados nos seguintes itens:

1. Questionamentos referentes às decisões em dados de PD e PR: na confecção dos questionários o pesquisador recai em quatro possibilidades;
 - a. Importância Falaciosa de atributos ou “*Prominence Hypothesis*”: o respondente escolhe observando um ou poucos atributos no seu processo de avaliação da escolha, desconhecendo que existe o efeito compensado (*trade-off*) entre os atributos que podem ser incorporados na escolha;
 - b. Viés Política da Resposta ou “*Policy-Response Bias*”: o respondente acha que poderá se beneficiar da resposta dada (Morikawa (1985) indica a leitura de de Dooley (1981) e Suzuki *et. al.* (1986) para melhor compreensão desse comportamento);
 - c. **Inércia das preferências**: o respondente não consegue compreender como o novo produto influenciaria na sua vida (Morikawa (1985) recomenda a leitura de Nelson (1979) para aprofundamento no termo) ; e
 - d. Viés de Justificativa ou “*Justification Bias*”: o respondente pode querer justificar comportamentos passados na sua resposta.

2. Descrição imperfeita das alternativas:

- a. Omissão de variáveis para simplificação dos questionários;
- b. Associar imagens a situações reais pode viesar a resposta ao perder a percepção do respondente frente à situação vivida e a situação presente na imagem;

3. Omissão de situações restritas:

- a. O respondente pode, consciente ou inconscientemente, ignorar suas restrições pessoais, uma vez que o questionário se aplica a uma hipótese.

Como descrito acima, os resultados podem ser viesados pela forma como os questionários foram aplicados. Mas cada uma das fontes de dados possui suas vantagens e limitações. Os dados de PR possuem como característica principal refletirem a participação atual de cada alternativa incidente no mercado atual. No entanto, os dados de PR possuem sérios aspectos limitantes, podendo-se destacar (Ortúzar e Willumsem, 1994):

- A existência de altas correlações entre os atributos, que impede a estimação isolada dos efeitos dos atributos;
- A dificuldade de estimar variáveis qualitativas; e
- A não possibilidade de avaliar alternativas que ainda não atuam no mercado atual.

As lacunas formadas pelas limitações dos dados de PR podem ser preenchidas pelas vantagens advindas no uso de dados de PD, dentre as quais é possível destacar as seguintes:

- Permitem o controle dos valores dos atributos através de projetos ortogonais, que permitem a estimação dos efeitos de cada atributo isoladamente;
- Permitem a análise de variáveis qualitativas, com o alcance de resultados satisfatórios;
- Permitem a análise de alternativas que ainda não existem no mercado atual.

Por outro lado, a principal limitação dos dados de PD consiste no fato de não refletirem o comportamento atual do mercado, característica peculiar dos dados de PR, como já mencionado. Sempre que viável, a partir de dados de PD ou de PR, a modelagem/estimação

conjunta com dados de PR e de PD é recomendável, pois assim seria possível, ao mesmo tempo, unir as vantagens e diminuir as limitações de cada fonte de dados (Morikawa, 1989 e Hensher, Louviere e Swait, 2000).

Em casos de estudos de previsão de demanda, a modelagem conjunta de dados de PD e PR pode melhorar a previsibilidade de comportamentos. Em estudos de análise compensações entre custo, qualidade e tempo de viagem “*trade-off*”, que permitem a determinação da importância relativa dos atributos, Hensher, Louviere e Swait (2000) recomendam a utilização de dados de PD.

Os métodos de estimação com dados conjuntos partem do princípio de explorar, por um lado, as estimações dos efeitos isolados de cada atributo da função de utilidade, obtidas com os dados de PD, e por outro, informações acerca da divisão atual de mercado, obtidas com os dados de PR. Esses métodos também são chamados de “*métodos de enriquecimento de dados*” Morikawa (1989), Ben-Akiva e Morikawa (1990), Morikawa, Ben-Akiva e Yamada (1991) estudam a heterogeneidade entre diferentes amostras e propõem uma forma de juntar diferentes amostras (dados de PD e PR). Na seção que segue são apresentadas as estratégias de estimação.

3. Estratégias de Estimação conjunta com fontes de dados distintas

Existem algumas formas de estimar modelos com a junção de dados. Swait, Louviere e Williams (1994) e Hensher, Louviere e Swait (2000) fazem descrição deste métodos e indicam que o trabalho de Morikawa (1989) é o primeiro trabalho consistente neste sentido. Além de Morikawa (1989) outros métodos de estimação com várias fontes de dados são apresentados. Para a validade de qualquer método de estimação levantam-se as seguintes hipóteses:

1. O processo de geração de dados das fontes utilizadas possuem mesmos parâmetros para os atributos comuns (média e desvio padrão), modificados apenas por um parâmetro de escala diferente para cada fonte de dados;
2. O termo erro das duas amostras tem distribuição de Valor Extremo II e é IID;

3. Os atributos de tempo e custo, que podem ser utilizados para avaliar políticas compensatórias em amostras de PR são deficientes, mas atributos socioeconômicos ou de oferta se aplicam a capturar aspectos de equilíbrio de mercado, ou seja, inerentes a decisão dos indivíduos;
4. Os atributos de equilíbrio de mercado em amostras de PD são deficientes, mas os outros atributos se aplicam a capturar aspectos de “*Trade-offs*”, ou seja, inerentes ao objeto escolhido.

Tomando as hipóteses acima temos duas possibilidades: utilizar os atributos de “*Trade-off*” como ligação entre as duas amostras, mesmo sabendo da deficiência deste na amostra de PR ou utilizar os atributos de “*Trade-off*” dos dados de PD e os de equilíbrio de PR e estimar em uma função conjunta. Assim, tomam-se os parâmetros de escala como relativos e modela-se a partir da hipótese 2 por uma função *Logit*. Para os dois casos podem-se estimar sequencialmente ou simultaneamente os parâmetros, com a vantagem para o último caso de gerar um estimador com melhor eficiência. Hensher, Louviere e Swait (2000) propõem dois procedimentos para o caso sequencial: uma estimativa com uma procura do fator de escala e outro com um modelo aninhado *Logit*. O primeiro caso é resumido na Equação 6 cuja a matriz de dados tem a estrutura apresentada na Equação 7.

$$LL = \sum_{n \in PR \cup PD} \sum_{i \in C_n} (y_{in}) \log(P_{ni}(Q(\lambda^{PD}))/\tau) \quad (6)$$

$$Q(\lambda^{PD}) = \begin{bmatrix} I^{PR} & 0 & X^{PR} & Z & 0 \\ 0 & \lambda^{PD} I^{PD} & \lambda^{PD} X^{PD} & 0 & \lambda^{PD} W \end{bmatrix} \quad (7)$$

Onde I^{PR} e I^{PD} são matrizes identidade; X^{PR} e X^{PD} matrizes de atributos de “*Trade-off*”; Z e W atributos de equilíbrio. A forma funcional para $P_{ni}(\cdot)$ dependerá da estrutura do erro, no caso da hipótese que aponta para Valor Extremo II, então tem-se um modelo *Logit*, conforme apresentado no segundo capítulo. O parâmetro λ é um fator de escala entre os atributos da PD e do PR, os X são atributos de “*Trade-off*” e os Z e W são os outros atributos, o τ é o desvio padrão do parâmetro e y é a escolha efetuada.

- algoritmo para a estimação segue os seguintes passos:
 - Escolhe-se um espectro de valores de λ^{PD} (lista de λ^{PD});
 - Escolhe-se um valor de λ^{PD} ;
 - Constrói-se a matriz 2;
 - Maximiza-se a função de verossimilhança (Equação 1) e guarda-se o resultado;
 - Faz-se os passos 1 a 5 repetidamente para a lista de valores de λ^{PD} ;
 - Comparam-se os valores da função de verossimilhança maximizada e escolhe-se o modelo que apresentou o menor valor.

Para o caso de estimação via estrutura aninhada, define-se uma árvore virtual onde, em um ramo, colocam-se os dados de PR e, no outro, os de PD, Como se supõe que os parâmetros são escalonados, então se restringe o parâmetro de PR como unitário, e, dessa forma tem-se os parâmetros estimados por uma estrutura aninhada. Alternativamente, Morikawa (1989) propõe sua estrutura sequencial que consiste nos seguintes passos:

1. Estima-se o modelo com os dados de PD e obtêm-se $\mu\beta^*$ e $\mu\gamma^*$ onde μ é o parâmetro de escala;
2. Constrói-se a variável $V = \mu\beta^* X^{PR}$;
3. Estima-se o modelo com os dados de PR com a variável V para obter os parâmetros de escala e das variáveis W, conforme Equação 8.

$$U_{in} = \lambda V + \alpha W + \varepsilon_{in} \quad (8)$$

4. Calcula-se $\mu = 1/\lambda$, $\beta = \mu\beta^*/\mu$ e $\alpha = \mu\alpha^*/\mu$

Caso deseje-se melhorar a precisão, multiplica-se os dados de X e Z do PD por μ , empilham-se os dados de PR e PD e estima-se novamente.

Por fim, para o caso simultâneo o método proposto por Morikawa (1989) consiste em fazer uma estimação via máxima verossimilhança utilizando a Equação (4). Na proposta, existi uma variável explicativa indicando que a amostra é de uma fonte de dado ou de outra. Formalmente tem-se a função de verossimilhança da junção dos dados de PD e PR, atendendo as Hipótese 1 a 4.

- **Hipótese 1:** $G_{PD}(\mu_x(x, z/\theta, \gamma)) = G_{PR}((x, y/\theta, \delta))$
- **Hipótese 2:** A função $G(*)$ é $C(2)$.
- **Hipótese 3:** A função $f(*)$ é $C(2)$, dado que é uma função multiplicativa de duas funções $C(2)$.
- **Hipótese 4:** A função $f(*)$ admite uma transformação monotônica.

Assim, a função união dos dados será:

$$f((x, z, y/\theta, \gamma, \delta, \mu) = (G_{PD}(\mu_x(x, z/\theta, \gamma))^d) * (G_{PR}((x, y/\theta, \delta))^{(1-d)}) \quad (9)$$

$$\text{Onde: } d = \begin{cases} 1, & \text{se a amostra for de PD} \\ 0, & \text{caso contrário} \end{cases}$$

$x, z, y :=$ são variáveis explicativas;

$\theta, \gamma, \delta, \mu :=$ são parâmetros a estimar.

Note que a função união, trata-se de uma função contínua em partes, onde, a continuidade em todo o espaço é garantida no ponto de quebra dado pela variável *dummy*, onde os dados de PD são complementados pelos dados de PR. A função densidade de escolha, para o conjunto de escolha $E=\{0,1\}$ é dado pela seguinte função de verossimilhança.

$$f((x, z, y/\theta, \gamma, \delta, \mu) = [(G_{PD}(\mu_x(x, z/\theta, \gamma))^d) * (G_{PR}((x, y/\theta, \delta))^{(1-d)})]^{(E_1)} * \\ * [(G_{PD}(\mu_x(x, z/\theta, \gamma))^d) * (G_{PR}((x, y/\theta, \delta))^{(1-d)})]^{(1-E_1)}$$

Logaritimizando a função de verossimilhança tem-se:

$$\log(f((x, z, y / \theta, \gamma, \delta, \mu)) = E_1 * [d * \log((G_{PD}(\mu_x(x, z / \theta, \gamma))) + (1 - d) * \log((G_{PR}((x, y / \theta, \delta))))] + (1 - E_1) * [d * \log((G_{PD}(\mu_x(x, z / \theta, \gamma))) + (1 - d) * \log((G_{PR}((x, y / \theta, \delta))))]$$

Assumindo verdadeira a hipótese acima tem-se:

$$\log(f((x, z, y / \theta, \gamma, \delta, \mu)) = [d * \log((G_{PD}(\mu_x(x, z / \theta, \gamma))) + (1 - d) * \log((G_{PR}((x, y / \theta, \delta))))]$$

Sendo a função $G(*)$ contínua, diferenciável e monótona, então tem-se uma solução interior. As variáveis z e y são comuns ou não às amostras, e a retirada de uma das variáveis em uma das amostras pode ser um artifício para a correção de multicolinearidade. Assumindo $G(*)$ como *Logit*, então se tem a continuidade e solução interior, podendo fazer a maximização.

Deve-se ter em mente que dado as naturezas dos dados diferentes a presença de heteroscedasticidade será inevitável, portanto, deve-se tratar a heteroscedasticidade inerente a estimação. Por outro lado, para atentar à hipótese de que os dados têm a mesma função geradora dos dados, ou seja, os dados das fontes distintas que serão utilizadas na modelagem possuem a mesma distribuição de probabilidade podendo ser considerados semelhantes, teste-se a compatibilidade da função geradora dos dados das diferentes fontes. Para tanto, deve-se fazer o teste proposto em Hensher, Louviere e Swait (2000) e que consiste em fazer uma estimativa para cada fonte de dados e outra conjunta obter a estatística descrita na Equação 10 e testar para uma distribuição qui-quadrada com $(|\beta| - 1)$ graus de liberdade.

$$-2((L^{PR} + L^{PD}) - L^{Conjunto}) \quad (10)$$

Na sequência, aplica-se a proposta de Morikawa (1989) para a estimação dos modelos, dada sua maior robustez e dado que as outras estratégias de estimação derivaram dela.

4. Estimação dos Modelos

Nessa seção serão estimados vários modelos estruturais (conjunto de quatro) a partir da teoria do valor do tempo com um modelo Logit simples utilizando os dados de PD e com o modelo Logit estimado pelo procedimento de junção das fontes de dados. Sequencialmente, será executada uma comparação com os resultados de modelos de PD sem enriquecimento de dados e com enriquecimento e dos VTV encontrados para os diferentes modelos.

4.1 Descrição dos Dados

Os dados são provenientes da pesquisa de preferência declarada e de satisfação do usuário feitas – dados de PR - para definição do PDOTIP-CE (2005). O motivador da pesquisa foi que o transporte intermunicipal regular no Ceará sofreu decréscimo de demanda intimamente ligado a problemas institucionais e de operação, com isso o clandestino se apropriou do mercado problemático propiciando um transporte mais ágil e com flexibilidade em tempo, em preço e com melhor acessibilidade.

Segundo o relatório do PDOTIP (2005) no sistema existem problemas de mobilidade e acessibilidade dos usuários de transporte regular com 13% da população do Ceará, a época, sem acesso direto a rede de transportes e que esses necessitavam de um deslocamento médio, á pé, de 3km para acessar os pontos de parada. No mesmo relatório é indicado 10% da população do interior do estado, a época, dispunha de automóveis nas suas residências. Além disso, indica ineficiência de rede e dos operadores de transportes e uma baixa qualidade dos serviços prestados, fazendo um diagnóstico dos problemas institucionais existentes e dá destaque a baixa competitividade nas linhas liberadas para operação.

As pesquisas foram realizadas nos terminais rodoviários e embarcadas nos veículos. Os dados obtidos foram: dados socioeconômicos dos usuários, motivo e frequência de viagem, origem e destino da viagem, e a escolha do indivíduo com relação a duas alternativas (ônibus e van, van e trem e trem e ônibus).

Para a modelagem utilizou-se a definição que trem e ônibus são modos regularmente estabelecidos enquanto que van trata-se de um modo clandestino. Os locais de coleta dos dados foram: Fortaleza, Juazeiro e Iguatu e os entrevistados tinham residência e destino em

diversos municípios do estado e em estados vizinhos, tais como Pernambuco, Piauí e Rio Grande do Norte.

Para o exercício empírico serão utilizadas as variáveis: sexo, número de carros na residência, número de moto na residência, idade do entrevistado, tempo da viagem, custo da viagem. Foram retiradas observações onde não existisse resposta para uma dessas variáveis.

Na pesquisa de PD um indivíduo era perguntado sobre várias hipóteses. Nesse caso adotou-se cada resposta como se fosse de um indivíduo, procedimento coerente estatisticamente, conformem Ben-Akiva e Lerman (1984), por assumir não haver incoerência na escolha dos indivíduos quanto a hipótese de transitividade das escolhas.

O custo da viagem usada na amostra de PR é resultado de uma pesquisa complementar executada com as cooperativas de vans e das tarifas praticadas nos ônibus, mas na pesquisa de satisfação de usuário foram feitas perguntas na escala de satisfação do usuário com relação ao custo da viagem. A união das amostras foi feita com a adição de uma variável que indica 1 se a observação pertence a PD e 0 se pertence a PR.

O processo de filtragem serviu para ter uma base de dados completa sem valores perdidos.

Um resumo estatístico das amostras segue na Tabela 1. Observe que as variáveis têm pequenas mudanças da PD para a PR. Os números de motos e de carros por residência e o tempo de viagem são os que apresentam maior mudança. No caso do tempo de viagem, a pesquisa de PR aparenta ter sido feita sobre indivíduos que fazem deslocamentos de grandes distâncias e o de PD os que fazem deslocamentos menores (entre municípios vizinhos), portanto, tem-se um complemento informacional entre PD e PR quanto a diversidade da amostra, mas podendo ser indesejado caso haja grandes idiosincrasias entre os usuários de transporte de curta e longa distância. No caso do número de motos e carros, pode ter, também, uma complementaridade, pois os indivíduos que possuem poucos veículos em casa aparecem na pesquisa de PD e os que tem mais, na pesquisa de PR.

Tabela 1 - Estatística Descritiva dos Dados de PD e PR¹

| | Pesquisa Declarada | | Pesquisa Revelada | |
|---------------|--------------------|---------------|-------------------|---------------|
| | Média | Desvio Padrão | Média | Desvio Padrão |
| Sexo | 0.51 | 0.50 | 0.49 | 0.50 |
| Idade | 34.25 | 14.10 | 34.62 | 13.80 |
| Auto | 0.25 | 0.49 | 1.20 | 0.48 |
| Moto | 0.23 | 0.47 | 1.25 | 0.50 |
| Tempo_viagem1 | 30.36 | 15.88 | 118.52 | 101.17 |
| Custo_Tarifa1 | 1.15 | 0.09 | 11.41 | 9.74 |
| Tempo_viagem2 | 38.80 | 16.17 | 109.04 | 93.07 |
| Custo_Tarifa2 | 1.03 | 0.05 | 9.31 | 7.95 |
| Nicho_1 | 0.64 | 0.48 | 0.61 | 0.49 |
| Amostra | 1554.00 | | 3662.00 | |

4.2 Estratégia de Estimação

Adotaremos a proposta simultânea de Morikawa (1989), embora um tratamento computacional mais complexo. As funções de utilidade são conforme as equações 11 e 12:

$$U_{PD}(\text{Tempo de viagem}, \text{Custo de viagem}; \theta) \quad (11)$$

$$U_{PR}(\text{Sexo}, \text{Idade}, \# \text{Auto}, \# \text{Moto}, \text{Custo de viagem}; \theta) \quad (12)$$

O atributo passível de compensação na análise será o tempo de viagem e os atributos de equilíbrio são o custo de transportes e as variáveis socioeconômicas. Os modelos seguiram modelos estruturais conforme a seguinte estrutura: Equação (13) para o modelo estrutural de Becker (1965), a Equação (14) para o DeSerpa (1971) e a Equação (15) para o de Tuong e Hensher (1985b).

$$V_i = -\lambda x_i - \mu t_i \quad (13)$$

$$V_i = -\lambda x_i - (\mu - k_i) t_i \quad (14)$$

¹ Fonte: Sistema de Transporte Intermunicipal de Passageiros do Ceará) para os estudos do PDOTIP-CE Sexo=masculino 1, feminino 0; idade= idade dos entrevistados; Auto= número de automóveis na residência; Moto = número de motos na residências; Tempo_viagem1=tempo de viagem no modo regular (Ônibus ou Tren); Tempo_viagem2=tempo de viagem no modo irregular (Van); Custo_Tarifa1=tarifa da viagem no modo regular (Ônibus ou Tren); Custo_Tarifa2=tarifa da viagem no modo irregular (Van); Nicho_1=se pertencente a PD=1, pertencente a PR=0.

$$V_i = -\lambda x_i - (\mu - \bar{k})t_i + (\alpha + \beta t_i^2 + \gamma x_i t_i) \quad (15)$$

Onde: V_i := nível da utilidade indireta associada a escolha da viagem i ;

x_i, t_i := são o custo e o tempo associado à viagem i ;

λ, μ := parâmetros a serem estimados que retratem utilidades marginais da renda e do tempo respectivamente;

k_i := multiplicador de lagrange associado ao consumo de tempo restrito a tecnologia e que retrata a utilidade marginal do tempo poupado.

O modelo não estruturado, *a la* McFadden (1974), combina tempos de viagem com a renda dos indivíduos de uma forma multiplicativa tal como feito em Cherchi e Ortúzar (2002). A relação $\frac{\mu}{\lambda}$ obtido da Equação (8) representa o preço sombra associado ao tempo da viagem.

No segundo modelo, o ganho em tempo, específico de cada alternativa i , é representado pela relação $\frac{k_i}{\lambda}$. Alternativamente, dada a impossibilidade de identificação de μ em (15), Tuong e Hensher (1985b) assumem que $k_i = f(t_i, p_i)$ e fazendo uma aproximação por expansão de série de Taylor na Equação (15), chega-se à expressão (16) para a utilidade indireta.

$$V_i = -\lambda x_i - (\mu - \bar{k})t_i + (\alpha + \beta t_i^2 + \gamma x_i t_i) \quad (16)$$

$$\text{Com } \alpha = -\left(\frac{\partial k}{\partial t}\right)_i \bar{t} - \left(\frac{\partial k}{\partial x}\right)_i \bar{x}; \quad \beta = \left(\frac{\partial k}{\partial t}\right)_i; \quad \gamma = \left(\frac{\partial k}{\partial x}\right)_i$$

Na abordagem de Tuong e Hensher (1985b) o valor do tempo é calculado da seguinte forma:

$$VOT = \left(\frac{\partial V}{\partial t_i} \right) \bigg|_{V_i = const} \quad (17)$$

$$VOT = \frac{(\mu - \bar{k})}{\lambda} - \frac{(\gamma x_i + 2\beta t_i)}{\lambda} \quad (18)$$

Com esta abordagem, pode-se fazer uma distinção entre diferentes viagens, tempo andando, tempo parado, etc. e especificar o valor de tempo para cada uma das subaditividades.

4.3 Resultados da Estimação

Os resultados dos diferentes modelos são apresentados na Tabela 2. O modelo de DeSerpa (1972) admite diferenças na percepção do tempo por modo do indivíduo, e assim o modelo não é estimado com diferenças de tempo, mas com a variável tempo em nível.

Os modelos estimados com dados de PD e enriquecidos apresentam diferenças relativas aos parâmetros de custo e os parâmetros resultantes da estimação com enriquecimento apresentaram valores menores, mas, curiosamente, todos os parâmetros são significantes estatisticamente, a exceção do parâmetro tempo, o mais importante para nossa análise. A aderência e o poder de previsão dos modelos é menor no *Logit* com enriquecimento do que no *Logit* simples, conforme pode ser observado no pseudo R^2 . A estimação com enriquecimento para o modelo de Tuong e Hensher apresenta mudança de magnitude e de sinal nos parâmetros. A estimação para este modelo foi feita se imaginado a possibilidade de correção da multicolinearidade inerente a este modelo, entretanto, o resultado não corroborou a ideia. Os resultados dos modelos de McFadden e Becker possuem parâmetros de custo iguais tal como em DeSerpa e Tuong e Hensher.

Tabela 2 - Comparação de modelos simples e com enriquecimento dos dados²

| | Logit Simples | | | | Logit Com Enriquecimento | | | |
|---|-----------------------|-----------------|-----------------|-------------------------|--------------------------|-----------------|-----------------|-----------------------|
| | McFadden (Conceitual) | Becker (5.1) | DeSerpa (5.2) | Tuong & Hensher (5.3) | McFadden (Conceitual) | Becker (5.1) | DeSerpa (5.2) | Tuong & Hensher (5.3) |
| const | 1.789 0.230 | 1.866 0.235 | 1.618 0.242 | 1.839 0.260 | 1.777 0.229 | 1.827 0.235 | 1.613 0.242 | 1.522 0.291 |
| Temp*SR | -0.046 0.016 | - | - | - | -0.045 0.016 | - | - | - |
| Temp*1SL | -0.046 0.013 | - | - | - | -0.045 0.013 | - | - | - |
| Temp*1A2SL | -0.062 0.013 | - | - | - | -0.061 0.013 | - | - | - |
| Temp*2A5SL | -0.061 0.014 | - | - | - | -0.060 0.014 | - | - | - |
| Temp*5A10SL | -0.040 0.019 | - | - | - | -0.039 0.019 | - | - | - |
| Custo | -6.518 1.174 | -6.870 1.194 | -8.204 1.245 | -9.417 -1.336 | -6.453 1.171 | -6.661 1.196 | -8.187 1.243 | -8.392 1.528 |
| Tempo | - | -0.057 0.011 | - | 0.007704 ** 0.039 | - | -0.055 0.011 | - | 0,07396 ** 0.041 |
| Tempo^2 | - | - | - | 0.000 0.000 | - | - | - | -0.001 0.001 |
| Tempo*Custo | - | - | - | -0.078 0.024 | - | - | - | -0.092 0.023 |
| Tempo_Reg | - | - | -0.076 0.012 | - | - | - | -0.076 0.012 | - |
| Tempo_Irreg | - | - | -0.061 0.011 | - | - | - | -0.061 0.011 | - |
| μ | | | | | -0.009 | -0.017 | -0.010 | -0.025 |
| N | 1554.000 | 1554.000 | 1554.000 | 1554.000 | 1695.000 | 1695.000 | 1695.000 | 1695.000 |
| Pseudo R ² | 0.018 | 0.018 | 0.026 | 0.027 | 0.016 | 0.016 | 0.023 | 0.022 |
| lnL | -1007.000 | -1007.000 | -998.900 | -997.900 | -1089.690 | -1089.522 | -1081.974 | -1082.636 |
| Graus de Liberdade | | | | | 11 | 7 | 8 | 9 |
| Teste de Hensher, Louviere & Swait (2000) | | | | | -4192.862 | -4192.526 | -4161.230 | -4160.554 |

A Tabela 3 apresenta o resultado da estimação apenas com os dados de PR, a amostra é de 3769 e foram incorporadas variáveis de equilíbrio de mercado, ou seja, inerente aos indivíduos e não ao objeto escolhido.

² Fonte: Autor(es), Const=constante; Temp*SR=produto entre tempo e uma *dummy* para os indivíduos que não apresentaram renda fixa; Temp*1SL=produto entre tempo e uma *dummy* para os indivíduos que apresentaram um salário mínimo; Temp*1A2SL=produto entre tempo e uma *dummy* para os indivíduos que apresentaram de um a dois salários mínimos; Temp*2A5SL=produto entre tempo e uma *dummy* para os indivíduos que apresentaram dois a cinco salários mínimos; Temp*5A10SL=produto entre tempo e uma *dummy* para os indivíduos que apresentaram cinco a dez salários mínimos; Custo=custo da viagem; Tempo=tempo da viagem; Tempo^2=quadrado do tempo da viagem; Tempo*Custo=produto do tempo da viagem pelo custo; Tempo_Reg=tempo da viagem do regular; Tempo_Irreg=Tempo da viagem do irregular. A renda base para o modelo McFadden é acima de dez salários mínimos. Os valores abaixo do parâmetro representam os desvios padrões, em dividindo o valor do parâmetro pelo desvio padrão e obtendo valores maiores que 2, entende-se que o parâmetro é significativa a 5%. Os ** indicam que o parâmetro não é significativa estatisticamente.

Tabela 3 - Resultados do modelo com os dados de preferência revelada³

| Revelada | | |
|------------------------|-------------|---------------|
| | coeficiente | desvio padrão |
| const | -0.4169 | 0.1625 |
| Sexo | -0.0449 | 0.0707 |
| Idade | -0.0003 | 0.0026 |
| Auto | 0.0952 | 0.0763 |
| Moto | -0.0994 | 0.0709 |
| Custo | -0.4896 | 0.0284 |
| N | 3769 | |
| %Previsão | 63% | |
| Log da verossimilhança | -2320.97 | |

Ao testar se os dados possuem as mesmas funções geradoras, rejeita-se a hipótese nula, não podendo ser afirmado que os dados possuem a mesma função geradora dos dados. Assim, pode-se concluir que os resultados não foram melhores do que os da estimativa simples. Para verificar a robustez do teste, far-se-á um estudo dos valores de tempo obtido com os modelos enriquecidos. Os valores de tempo para os modelos são apresentados na Tabela 4. Em McFadden foi obtido da razão entre a mediana dos coeficientes estimados para a relação tempo renda pelo parâmetro de custo. Os outros modelos foram calculados conforme equações apresentadas acima, com destaque ao fato que, no modelo de Tuong e Hensher, utilizou-se o valor dos tempos e custos médios.

Tabela 4 - Valor de Tempo para os diferentes modelos em reais por minuto⁴

| | Simples | | | | Com Enriquecimento | | | |
|----------|----------|--------|---------|-----------------|--------------------|--------|---------|-----------------|
| | McFadden | Becker | DeSerpa | Tuong & Hensher | McFadden | Becker | DeSerpa | Tuong & Hensher |
| VTVReg | 0.422 | 0.497 | 0.554 | 0.691 | 0.467 | 0.499 | 0.553 | 0.507 |
| VTVIrreg | 0.422 | 0.497 | 0.446 | 0.662 | 0.467 | 0.499 | 0.446 | 0.510 |

Os resultados com enriquecimento aumentaram os valores de tempo de viagem estimados para os modelos de McFadden, Becker, o de DeSerpa foi imperceptível e o modelo Tuong e Hensher apresentou redução do valor do tempo de viagem. Note que a adoção de um enriquecimento aproximou as estimativas de VTL do modelo estruturado de Tuong e Hensher para os outros, possivelmente em decorrência de uma correção da multicolinearidade inerente

³ Fonte: autores.

⁴ Idem.

ao modelo de Tuong e Hensher. Além disso, pelo fato do aumento da amostra ter sido pequeno, as mudanças de valores de tempo de viagem não foram grandes quando se compara o modelo *Logit* simples com o *Logit* com enriquecimento.

Conclusão

Este trabalho objetivou comparar valores de tempo de viagem obtidos com modelos *Logit* simples e adotando a proposta de Morikawa (1989) para combinação de fontes de dados de PD e PR. O trabalho apresentou uma revisão exaustiva sobre modelos de escolha discreta e da possibilidade de extrair informações combinando fontes de dados distintas. Indicou, assim, que os métodos de estimação de modelos de escolha discreta evoluíram a partir do método de enriquecimento de dados proposto em Morikawa (1989), o qual pode gerar ganhos em estimação pontual dependendo se a função geradora dos dados das diferentes fontes ser compatível, conforme salientado em Hensher, Louviere e Swait (2000).

Como resultado, identificou-se que o enriquecimento de dados pode gerar informações falaciosas sobre a propensão a pagar e o valor de tempo para o usuário. É importante, portanto, compreender a função geradora dos dados para não incorrer em erros. Para resultarem em uma boa combinação e em informações mais acertadas, os dados devem ter propriedades semelhantes, ou seja, possuírem a mesma fonte geradora de dados. Assim, Swait *et. al.* (1999), estão corretos ao indicar que se deve utilizar a técnica de enriquecimento com parcimônia, pois a proposta de Morikawa (1989), conforme foi ilustrado, dependendo da especificação funcional possui limitações e pode gerar estimativas controversas. De outro modo, este trabalho amplia a abordagem de Morikawa (1989) por formalizar as hipóteses necessárias a uma estimação com duas fontes de dados e o e descreve os problemas inerentes às fontes de dados.

São várias as possibilidades futuras de pesquisa, tomando este artigo como ponto inicial. Pode-se investigar as melhorias de estimação ou solução dos problemas de multicolinearidade em modelos de escolha discreta, por meio da combinação de fontes distintas de dados. Pode-se também investigar uma estimação conjunta de dados utilizando uma abordagem bayesiana ou em simulação, objetivando melhorar as estimativas de parâmetros de valor de tempo de viagem ou estender para o conceito de propensão a pagar.

Referências

- Becker, G. S. (1965) A theory of the allocation of time. *Economic Journal*, vol. 75, n. 299, pp. 493-517.
- Ben-Akiva, M. e S. Lerman (1985) *Discrete choice analysis*. The MIT Press, Cambridge Massachusetts.
- Ben-Akiva, M. e Morikawa, T. (1990) Estimation of switching models from revealed preferences and stated intentions. *Transportation Research A*, vol. 24, pp. 485-495.
- Cherchi, E. e Ortúzar J. D. (2002) Mixed RP/SP models incorporating interaction effects: modeling new suburban train services in Cagliari. *Transportation*, vol. 29, n. 4, pp. 371-395
- DeSerpa, A. C., (1971) A Theory of the economics of time, *Economic Journal*, vol. 8, pp. 28-846.
- DeSerpa, A. C., (1972) Microeconomic theory and the valuation of travel time: some clarification. *Regional and Urban Economics*, vol. 2, n. 4, pp. 401-410.
- DaSilva, F. G. F. (2012) Valor de tempo de viagem e idiosincrasia dos usuários do transporte regular e clandestino no Ceará: um estudo empírico via estimativa bayesiana. *Journal of Transport Literature*, vol. 6, n. 1, pp. 71-92.
- Granjeiro, C. F. (2006) *Plano Diretor e Operacional do Transporte Intermunicipal de Passageiros do Estado do Ceará (PDOTIP-CE)*. Contrato DERT-CE/ASTEF.
- Hensher, D. e Bradley, M. (1993) Using stated response data to enrich revealed preference discrete choice models. *Marketing Letters*, vol. 4, pp. 39-152.
- Hensher, D., Louviere, J. e Swait, J. (1999) Combining sources of preference data. *Journal of Econometrics*, vol. 87, pp. 97-221.
- Jiang, M. e Morikawa, T. (2004) Theoretical analysis on the variation of value of travel time savings. *Transportation Research A*, vol. 38, pp. 551-571.
- Lenk, P.J., Desarbo, W.S., Green, P.E. e Young, M.R. (1996) Hierarchical Bayes conjoint analysis: recovery of partworth heterogeneity from reduced experimental designs. *Marketing Science*, vol. 15, pp. 173-191.
- McFadden, D. L., (1974) The measurement of urban travel demand, *Journal of Public Economics*, vol. 3, n. 4, pp. 303-328.
- Morikawa, T. (1989) Incorporating Stated Preference Data in Travel Demand Analysis, *Ph.D. Dissertation, Department of Civil Engineering, MIT*.
- Morikawa, T., Ben-Akiva, M. e Yamada, K. (1991) Forecasting Intercity Rail Ridership Using Revealed Preference and Stated Preference Data. *Transportation Research Record*, vol. 1328, pp. 30-35.
- Ortúzar, J. D. e Willumsen, L. G. (1997) *Modeling Transport*. Second Edition. New York: John Wiley & Sons.
- R Development Core Team (2009). *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria. Available at www.R-project.org.
- Swait, J., Louviere, J. e Williams, M. (1994) A Sequential Approach to Exploiting the Combined Strengths of SP and RP Data: Application to Freight Shipper Choice. *Transportation*, vol. 21, pp. 135-152.
- Train, K. E., (2003), *Discrete Choice Models with Simulation*, Cambridge: Cambridge Press.

Tuong, P. e Hensher D. A., (1985a) Measurement of Travel Time Values and Opportunity Cost from a Discrete-Choice Model. *Economic Journal*, vol. 95, n. 378, pp.438-451.

Tuong, P. e Hensher D. A., (1985b) A Valuation of Travel Time Savings A Direct Experimental Approach. *Journal of Transport Economics and Policy*, vol. 19, pp. 237-261.