

PONTIFÍCIA UNIVERSIDADE CATÓLICA
DO RIO DE JANEIRO



Rodrigo Pinto Moreira

**Modelo de Regressão Logística com Transição Suave
Estruturado por Árvore (STLR-Tree)**

Dissertação de Mestrado

Dissertação apresentada como requisito parcial para
obtenção do título de Mestre pelo Programa de Pós-
Graduação em Engenharia Elétrica da PUC-Rio.

Orientador: Prof. Dr. Álvaro Veiga

Rio de Janeiro
Abril de 2008



Rodrigo Pinto Moreira

**Modelo de Regressão Logística com Transição Suave
Estruturado por Árvore (STLR-Tree)**

Dissertação apresentada como requisito parcial para obtenção do grau de Mestre pelo Programa de Pós-Graduação em Engenharia Elétrica do Departamento de Engenharia Elétrica do Centro Técnico Científico da PUC-Rio. Aprovada pela Comissão Examinadora abaixo assinada.

Prof. Álvaro de Lima Veiga Filho
Orientador

Departamento de Engenharia Elétrica – PUC-Rio

Prof. Marcelo Cunha Medeiros
Departamento de Economia – PUC-Rio

Prof. Joel Maurício Corrêa da Rosa
UFPR

Prof. José Eugenio Leal
Coordenador Setorial do Centro
Técnico Científico

Rio de Janeiro, 11 de abril de 2008

Todos os direitos reservados. É proibida a reprodução total ou parcial do trabalho sem autorização da universidade, da autora e do orientador.

Rodrigo Pinto Moreira

Graduou-se em Ciência Estatísticas na Escola Nacional de Ciências Estatísticas - ENCE (Rio de Janeiro, Brasil). Durante o mestrado em Engenharia Elétrica trabalhou com técnicas de análise estatística multivariada, modelagem linear e não-linear e em projetos na área de seguros colaborando com seu orientador no desenvolvimento de modelos internos para seguradoras.

Ficha Catalográfica

Moreira, Rodrigo Pinto

Modelo de regressão logística com transição suave estruturado por árvore (STLR-Tree) / Rodrigo Pinto Moreira ; orientador: Álvaro de Lima Veiga Filho. – 2008.

81 f. : il. ; 30 cm

Dissertação (Mestrado em Engenharia Elétrica)– Pontifícia Universidade Católica do Rio de Janeiro, Rio de Janeiro, 2008.

Inclui bibliografia

1. Engenharia elétrica – Teses. 2. Modelos não-lineares estruturados por árvore. 3. Classificação. 4. Regressão logística. Árvore de classificação e regressão (CART). I. Veiga Filho, Álvaro de Lima. III. Pontifícia Universidade Católica do Rio de Janeiro. Departamento de Engenharia Elétrica. IV. Título.

CDD: 621.3

Agradecimentos

Ao meu orientador Álvaro Veiga Lima Filho, pelo apoio e incentivo a este trabalho.

À FAPERJ, CNPq, CAPES e à PUC-Rio, pelos auxílios financeiros concedidos, sem os quais este trabalho não poderia ter sido realizado.

Aos meus pais e minha irmã. Família que me atura e ajuda diariamente.

A todos os meus outros familiares, principalmente minhas avós e meu avô, que por muitos anos ainda me darão força.

Ao meu avô João (*in memoriam*), que certamente está me ajudando de um bom lugar.

À minha futura esposa, Suene, e à sua família. Ela foi a pessoa que mais me aturou no decorrer deste trabalho.

Aos meus queridos amigos da TDP, ENCE e todos os demais.

Aos meus companheiros da PUC-Rio, principalmente aos freqüentadores da sala L604, na favelinha.

Aos professores da ENCE, Kaizô Beltrão e Sandra Canton.

Aos professores Cristiano Fernandes, Marcelo Medeiros e Joel Corrêa da Rosa.

Ao mestrando do ICA, Gustavo Victor C. Ortega, pela ajuda com a aplicação de Redes Neurais e na obtenção dos dados para a mesma.

Ao pessoal da secretaria e do suporte do departamento de Engenharia Elétrica.

Enfim, a todos aqueles que contribuíram de forma direta ou indireta na realização deste feito.

Resumo

Moreira, Rodrigo Pinto; Veiga, Álvaro. **Modelo de Regressão Logística com Transição Suave Estruturado por Árvore (STLR-Tree)**. Rio de Janeiro, 2008. 82p. Dissertação de Mestrado — Departamento de Engenharia Elétrica, Pontifícia Universidade Católica do Rio de Janeiro.

Este trabalho tem como objetivo principal adaptar o modelo STR-Tree, o qual é a combinação de um modelo *Smooth Transition Regression* com *Classification and Regression Tree (CART)*, a fim de utilizá-lo em Classificação. Para isto algumas alterações foram realizadas em sua forma estrutural e na estimação. Devido ao fato de estarmos fazendo classificação de variáveis dependentes binárias, se faz necessária a utilização das técnicas empregadas em Regressão Logística, dessa forma a estimação dos parâmetros da parte linear passa a ser feita por Máxima Verossimilhança. Assim o modelo, que é paramétrico não-linear e estruturado por árvore de decisão, onde cada nó terminal representa um regime os quais têm seus parâmetros estimados da mesma forma que em uma Regressão Logística, é denominado *Smooth Transition Logistic Regression Tree (STLR-Tree)*. A inclusão dos regimes, determinada pela divisão dos nós da árvore, é feita baseada em testes do tipo Multiplicadores de Lagrange, que em sua forma para o caso Gaussiano utiliza a Soma dos Quadrados dos Resíduos em suas estatísticas de teste, aqui é substituída pela Função Desvio (*Deviance*), que é equivalente para o caso dos modelos não Gaussianos, cuja distribuição da variável dependente pertença à família exponencial. Na aplicação a dados reais selecionou-se dois conjuntos das variáveis explicativas de cada uma das duas bases utilizadas, que resultaram nas melhores taxas de acerto, verificadas através de Tabelas de Classificação (Matrizes de Confusão). Esses conjuntos de variáveis foram usados com outros métodos de classificação existentes, são eles: *Generalized Additive Models (GAM)*, *Regressão Logística*, *Redes Neurais*, *Análise Discriminante*, *k-Nearest Neighbor (k-NN)* e *Classification and Regression Trees (CART)*.

Palavras-chave

Modelos não-lineares estruturados por árvore. Classificação. Regressão Logística. Árvores de Classificação e Regressão (CART).

Abstract

Moreira, Rodrigo Pinto; Veiga, Álvaro. **Smooth Transition Logistic Regression Model Tree**. Rio de Janeiro, 2008. 82p. MsC Thesis — Department of Electric Engineering, Pontifícia Universidade Católica do Rio de Janeiro.

The main goal of this work is to adapt the STR-Tree model, which is the combination of a *Smooth Transition with Regression* model with *Classification and Regression Tree (CART)*, in order to use it in Classification. Some changes were made in its structural form and in the estimation. Due to the fact we are doing binary dependent variables classification, is necessary to use the techniques employed in Logistic Regression, so the estimation of the linear part will be made by Maximum Likelihood. Thus the model, which is nonlinear parametric and structured by a decision tree, where each terminal node represents a regime that have their parameters estimated in the same way as in a Logistic Regression, is called *Smooth Transition Logistic Regression Tree (STLR-Tree)*. The inclusion of the regimes, determined by the splitting of the tree's nodes, is based on Lagrange Multipliers tests, which for the Gaussian cases uses the Residual Sum-of-squares in their test statistic, here are replaced by the *Deviance function*, which is equivalent to the case of non-Gaussian models, that has the distribution of the dependent variable in the exponential family. After applying the model in two datasets chosen from the bibliography comparing with other methods of classification such as: *Generalized Additive Models (GAM)*, *Logistic Regression*, *Neural Networks*, *Discriminant Analyses*, *k-Nearest Neighbor (k-NN)* and *Classification and Regression Trees (CART)*. It can be seen, verifying in the Classification Tables (Confusion Matrices) that STLR-Tree showed the second best result for the overall rate of correct classification in three of the four applications shown, being in all of them, behind only from GAM.

Keywords

Tree structured nonlinear models. Classification. Logistic Regression. Classifications and Regression Trees (CART).

Sumário

1	Introdução	13
2	Regressão Logística	15
2.1	Revisão de Modelos Lineares Generalizados (MLG)	15
2.2	Dados binários (Regressão Logística)	18
3	Modelos e metodologias comparadas	29
3.1	Classification and Regression Trees (CART)	29
3.2	Generalized Additive Models (GAM)	33
3.3	k-Nearest Neighbor	36
3.4	Análise Discriminante	37
4	Modelo de Regressão Logística com Transição Suave Estruturado por Árvore (STLR-Tree)	40
4.1	Revisão do STR-Tree	40
4.2	Especificação do STLR-Tree	43
4.3	Estimação do STLR-Tree	47
4.4	Avaliação do STLR-Tree	49
4.5	Ciclo de Modelagem	49
5	Aplicação e comparação dos métodos	53
5.1	Bases de dados	53
6	Conclusão	66
	Referências Bibliográficas	68
A	Alguns Modelos Não-lineares	71
A.1	Threshold Auto Regressive (TAR)	71
A.2	Self-Exiting Threshold Auto Regressive (SETAR)	72
A.3	Smooth Transition Autoregression (STAR)	72
A.4	Logistic Smooth Transition Autoregression (LSTAR)	73
A.5	Exponencial Smooth Transition Autoregression (ESTAR)	74
A.6	Multiple Regime Smooth Transition Autoregression (MRSTAR)	74
A.7	Neural Coefficient Smooth Transition Autoregressive (NCSTAR)	74
A.8	Smooth Transition Regression (STR)	75
B	Comando do programa R 2.6.2	76
B.1	Comandos para GAM	76
B.2	Comandos para CART	76
B.3	Comandos para k-NN	77
B.4	Comandos para Regressão Logística	77
C	Estatísticas Descritivas	78
C.1	E-mail/Spam	78

C.2	Doenças Cardíacas na África do Sul	78
C.3	Fraude/Irregularidade no Consumo de Energia Elétrica	79
D	Estimativas dos Coeficientes	80
D.1	E-mail/Spam	80
D.2	Doenças Cardíacas na África do Sul	80
D.3	Fraude/Irregularidade no Consumo de Energia Elétrica	81
D.4	Coeficientes dos parâmetros não-lineares	82

Lista de figuras

3.1	Estrutura do modelo. Exemplo em (7)	31
3.2	Divisão do espaço das covariáveis	31
3.3	Hierarquia dos modelos	33
3.4	Exemplo: k-Nearest Neighbor	37
4.1	Exemplo Árvore 1	42
4.2	Exemplo Árvore 2	43
4.3	Dados gerados ($c_0 = 4.3$ e $s_0 = 1$ e $\gamma = 0.5$)	52
4.4	Dados gerados ($c_0 = 4.3$ e $s_0 = 1$ e $\gamma = 50$)	52
5.1	Estrutura do modelo - Spam	55
5.2	Estrutura do modelo - DCAS	58
5.3	Estrutura do modelo - Consumo de Energia	61
5.4	Gráfico das taxas de erro total - E-mail/Spam	65
5.5	Gráfico das taxas de erro total - DCAS	65
5.6	Gráfico das taxas de erro total - Fraude no Consumo de Energia Elétrica	65

Lista de tabelas

2.1	Tabela de Classificação	27
2.2	Qualidade do ajuste - ROC	28
3.1	Divisão do espaço das covariáveis	32
3.2	Matriz de Confusão	38
5.1	Tabela de Classificação (<i>in sample</i>) - Spam	56
5.2	Comparação das Taxas de Acerto (<i>in sample</i>) - Spam	56
5.3	Métodos de Classificação Ordenados por Taxas de Acerto (<i>in sample</i>) - Spam	56
5.4	Comparação das Taxas de Acerto (<i>out-of-sample</i>) - Spam	57
5.5	Métodos de Classificação Ordenados por Taxas de Acerto (<i>out-of-sample</i>) - Spam	57
5.6	Tabela de Classificação (<i>in sample</i>) - DCAS	58
5.7	Comparação das Taxas de Acerto (<i>in sample</i>) - DCAS	59
5.8	Métodos de Classificação Ordenados por Taxas de Acerto (<i>in sample</i>) - DCAS	59
5.9	Comparação das Taxas de Acerto (<i>out-of-sample</i>) - DCAS	59
5.10	Métodos de Classificação Ordenados por Taxas de Acerto (<i>out-of-sample</i>) - DCAS	59
5.11	Tabela de Classificação (<i>in sample</i>) - Consumo de Energia	62
5.12	Comparação das Taxas de Acerto (<i>in sample</i>) - Fraude no Consumo de Energia Elétrica	63
5.13	Métodos de Classificação Ordenados por Taxas de Acerto (<i>in sample</i>) - Fraude no Consumo de Energia	63
5.14	Comparação das Taxas de Acerto (<i>out-of-sample</i>) - Fraude no Consumo de Energia Elétrica	63
5.15	Métodos de Classificação Ordenados por Taxas de Acerto (<i>out-of-sample</i>) - Fraude no Consumo de Energia Elétrica	64
C.1	Estatísticas Descritivas - Spam	78
C.2	Estatísticas Descritivas - DCAS	78
C.3	Estatísticas Descritivas - Fraude no Consumo de Energia	79
D.1	Coeficientes - Spam	80
D.2	Coeficientes - DCAS	80
D.3	Coeficientes - Fraude no Consumo de Energia	81
D.4	Pesos Redes Neurais - Fraude no Consumo de Energia	81
D.5	Coeficientes Não-lineares	82

Jamais considere seus estudos como uma obrigação, mas como uma oportunidade invejável para aprender a conhecer a influência libertadora da beleza do reino do espírito, para seu próprio prazer pessoal e para proveito da comunidade à qual seu futuro trabalho pertencer.

Albert Einstein, *Notas Autobiográficas Albert Einstein*.
