



**Rodrigo Tosta Peres**

## **Seleção de Dados em LVQ**

### **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: Carlos Eduardo Pedreira

Rio de Janeiro, agosto de 2004



**Rodrigo Tosta Peres**

## **Seleção de Dados em LVQ**

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

**Carlos Eduardo Pedreira**

Orientador  
PUC-Rio

**Álvaro Veiga**

PUC-Rio

**Alexandre P. Alves da Silva**

UFRJ

**Felipe França**

UFRJ

**José Eugênio Leal**

Coordenador(a) Setorial do Centro Técnico Científico - PUC-Rio

Rio de Janeiro, 12 de agosto de 2004

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

**Rodrigo Tosta Peres**

Graduado em Licenciatura em Matemática pela UFF (Universidade Federal Fluminense), em 2000.

Ficha Catalográfica

Peres, Rodrigo Tosta

Seleção de dados em LVQ / Rodrigo Tosta Peres; orientador: Carlos Eduardo Pedreira. – Rio de Janeiro: PUC, Departamento de Engenharia Elétrica, 2004.

68 f. : il. ; 30 cm

Dissertação (mestrado) – Pontifícia Universidade Católica do Rio de Janeiro, Departamento de Engenharia Elétrica.

Inclui referências bibliográficas.

1. Engenharia elétrica – Teses. 2. Classificação de padrões. 3. Seleção de dados. 4. LVQ. 5. Aprendizado por quantização vetorial. 6. Active learning. 7. Aprendizado exaustivo. 8. Redes neurais. I. Pedreira, Carlos Eduardo. II. Pontifícia Universidade Católica do Rio de Janeiro. Departamento de Engenharia Elétrica. III. Título.

CDD: 621.3

Aos meus pais, Mário e Suely,  
pelo amor incondicional durante todos esses anos.

## Agradecimentos

A Deus, que está comigo durante todo o tempo.

Ao meu orientador, Professor Carlos Eduardo Pedreira, pelo apoio em todas as horas, pelo espírito científico que vem me ajudando a desenvolver, pela confiança e, acima de tudo, pela amizade.

A PUC-Rio, pelos auxílios concedidos, sem os quais esse trabalho não poderia ter sido realizado.

Aos meus pais e aos meus irmãos, Diego e Tiago, por todo o amor, carinho e paciência, além de terem me dado o suporte para que eu pudesse chegar até aqui.

A Professora Ana Pavani, por tudo que me ensinou durante os anos de trabalho, pela amizade e respeito que ficarão para sempre.

Ao Professor Álvaro Veiga pela força e pelo conselho de ir em frente mesmo nas horas difíceis.

Ao Professor Yaser Abu-Mostafa e a sua equipe pelo apoio técnico durante o desenvolvimento do trabalho, especialmente ao Hsuan-Tien Lin pela troca de idéias sempre interessante.

A todos os familiares, amigos e colegas que, direta ou indiretamente, contribuíram para que esse trabalho pudesse ser desenvolvido.

## Resumo

Peres, Rodrigo Tosta; Pedreira, Carlos Eduardo. **Seleção de Dados em LVQ**. Rio de Janeiro, 2004. 68p. Dissertação de Mestrado - Departamento de Engenharia Elétrica, Pontifícia Universidade Católica do Rio de Janeiro.

Nesta dissertação, propomos uma metodologia para seleção de dados em modelos de Aprendizado por Quantização Vetorial, referenciado amplamente na literatura pela sigla em inglês LVQ. Treinar um modelo (ajuste dentro-da-amostra) com um subconjunto selecionado a partir do conjunto de dados disponíveis para o aprendizado pode trazer grandes benefícios no resultado de generalização (fora-da-amostra). Neste sentido, é muito importante realizar uma busca para selecionar dados que, além de serem representativos de suas distribuições originais, não sejam ruído (no sentido definido ao longo desta dissertação). O método proposto procura encontrar os pontos relevantes do conjunto de entrada, tendo como base a correlação do erro de cada ponto com o erro do restante da distribuição. Procura-se, em geral, eliminar considerável parte do ruído mantendo os pontos que são relevantes para o ajuste do modelo (aprendizado). Assim, especificamente em LVQ, a atualização dos protótipos durante o aprendizado é realizada com um subconjunto do conjunto de treinamento originalmente disponível. Experimentos numéricos foram realizados com dados simulados e reais, e os resultados obtidos foram muito interessantes, mostrando claramente a potencialidade do método proposto.

## Palavras-chave

Classificação de Padrões; Seleção de Dados; LVQ; Aprendizado por Quantização Vetorial; Active Learning; Aprendizado Exaustivo; Redes Neurais.

## Abstract

Peres, Rodrigo Tosta; Pedreira, Carlos Eduardo. **Data Selection for LVQ**. Rio de Janeiro, 2004. 68p. MSc. Dissertation - Departamento de Engenharia Elétrica, Pontifícia Universidade Católica do Rio de Janeiro.

In this dissertation, we consider a methodology for selection of data in models of Learning Vector Quantization (LVQ). The generalization can be improved by using a subgroup selected from the available data set. We search the original distribution to select relevant data that aren't noise. The search aims at relevant points in the training set based on the correlation between the error of each point and the average of error of the remaining data. In general, it is desired to eliminate a considerable part of the noise, keeping the points that are relevant for the learning model. Thus, specifically in LVQ, the method updates the prototypes with a subgroup of the originally available training set. Numerical experiments have been done with simulated and real data. The results were very interesting and clearly indicated the potential of the method.

## Keywords

Pattern Classification; Data Selection; LVQ; Learning Vector Quantization; Active Learning; Exhaustive Learning; Neural Networks.

# Sumário

|   |    |
|---|----|
| 1 Introdução  | 13 |
| 2 Uma Proposta para Seleção de Dados em Modelos LVQ                             | 15 |
| 2.1. Seleção de Dados em LVQ  | 15 |
| 2.2. Uma Proposta de Seleção de Dados: Algoritmo de Alexander – Abu-<br>Mostafa | 20 |
| 2.2.1. Aprendizado Exaustivo  | 20 |
| 2.2.2. Generalização sob Aprendizado Exaustivo                                  | 20 |
| 2.2.3. Correlações entre Erro de Treinamento e Erro de Teste                    | 21 |
| 2.2.4. Análise de $\rho$  | 22 |
| 2.2.4.1. Estimação de Ruído   | 22 |
| 2.2.4.2. Fronteira de Decisão e Zona de Risco                                   | 23 |
| 2.3. Uma Proposta de Seleção de Dados para Modelos LVQ                          | 23 |
| 2.3.1. Região de Sorteio de Protótipos  | 24 |
| 2.3.2. Heurística para Seleção de Dados   | 25 |
| 2.3.3. O Algoritmo Principal  | 27 |
| 3 Resultados Numéricos  | 28 |
| 4 Conclusão   | 55 |
| 5 Referências   | 57 |
| 6 Apêndice A – Seleção de Dados   | 59 |
| Active Learning, Query-Based Learning e Sequential Design                       | 60 |
| Uma Comparação Crítica entre os Três Métodos                                    | 62 |
| 7 Apêndice B – LVQ  | 63 |
| LVQ 1   | 64 |
| LVQ 2.1   | 65 |



|  |    |
|--|----|
| LVQ 3                                    | 66 |
| Análise de Margem                        | 67 |
| Uma Avaliação Crítica Sobre o Modelo LVQ | 67 |

## Lista de figuras

|  |    |
|--|----|
| Figura 1 - Distribuições 1 e 2 com protótipos nos centróides | 16 |
| Figura 2 - Distribuições 1 e 2 com protótipos treinados      | 17 |
| Figura 3 - Distribuições 1 e 2 com ponto classificado errado | 18 |
| Figura 4 - Exemplo 1: Distribuições Originais                | 31 |
| Figura 5 - Exemplo 1: Distribuições após corte em zero       | 32 |
| Figura 6 - Exemplo 1: Histograma de Rho                      | 32 |
| Figura 7 - Exemplo 1: Função Estimada pela Janela de Parzen  | 33 |
| Figura 8 - Exemplo 1: Distribuições após corte no mínimo     | 33 |
| Figura 9 - Exemplo 2: Distribuições Originais                | 36 |
| Figura 10 - Exemplo 2: Distribuições após corte em zero      | 37 |
| Figura 11 - Exemplo 2: Histograma de Rho                     | 37 |
| Figura 12 - Exemplo 2: Função Estimada pela Janela de Parzen | 38 |
| Figura 13 - Exemplo 3: Distribuições Originais               | 40 |
| Figura 14 - Exemplo 3: Distribuições após corte em zero      | 41 |
| Figura 15 - Exemplo 3: Histograma de Rho                     | 41 |
| Figura 16 - Exemplo 3: Função Estimada pela Janela de Parzen | 42 |
| Figura 17 - Exemplo 3: Distribuições após corte no mínimo    | 42 |
| Figura 18 - Exemplo 4: Histograma de Rho                     | 45 |
| Figura 19 - Exemplo 4: Função Estimada pela Janela de Parzen | 45 |
| Figura 20 - Exemplo 5: Histograma de Rho                     | 48 |
| Figura 21 - Exemplo 5: Função Estimada pela Janela de Parzen | 48 |
| Figura 22 - Exemplo 6: Histograma de Rho                     | 51 |
| Figura 23 - Exemplo 6: Função Estimada pela Janela de Parzen | 52 |
| Figura 24 - Exemplo 7: Histograma de Rho                     | 54 |
| Figura 25 - Exemplo 7: Função Estimada pela Janela de Parzen | 54 |

## Lista de tabelas

|  |    |
|--|----|
| Tabela 1 - Classificação antes e depois da Atualização               | 18 |
| Tabela 2 - Exemplo 1 – Quantidade de Ruído Detectado pelo Método     | 30 |
| Tabela 3 - Exemplo 1 – Quantidade de Pontos Relevantes               |    |
| Equivocadamente Eliminados   | 30 |
| Tabela 4 - Exemplo 1 - Classificação                                 | 31 |
| Tabela 5 - Exemplo 2 – Quantidade de Ruído Detectado pelo Método     | 35 |
| Tabela 6 - Exemplo 2 – Quantidade de Pontos Relevantes               |    |
| Equivocadamente Eliminados   | 35 |
| Tabela 7 - Exemplo 2 - Classificação                                 | 36 |
| Tabela 8 - Exemplo 3 – Quantidade de Ruído Detectado pelo Método     | 39 |
| Tabela 9 - Exemplo 3 – Quantidade de Pontos Relevantes               |    |
| Equivocadamente Eliminados   | 39 |
| Tabela 10 - Exemplo 3 - Classificação                                | 40 |
| Tabela 11 - Exemplo 4 – Quantidade de Ruído Detectado pelo Método    | 43 |
| Tabela 12 - Exemplo 4 – Quantidade de Pontos Relevantes              |    |
| Equivocadamente Eliminados   | 44 |
| Tabela 13 - Exemplo 4 - Classificação                                | 44 |
| Tabela 14 - Exemplo 5 – Quantidade de Ruído Detectado pelo Método    | 46 |
| Tabela 15 - Exemplo 5 – Quantidade de Pontos Relevantes              |    |
| Equivocadamente Eliminados   | 47 |
| Tabela 16 - Exemplo 5 - Classificação                                | 47 |
| Tabela 17 - Exemplo 6 - Classificação do banco Iris (leave-one-out)  | 50 |
| Tabela 18 - Exemplo 7 - Classificação do banco Glass (leave-one-out) | 53 |

A sabedoria é suprema; portanto, adquira a sabedoria.  
Sim, com tudo o que possuiis adquira o entendimento.  
(Provérbios 4:7)