



Universidade Federal de Pernambuco
Centro de Informática

Graduação em Engenharia da Computação

Análise comparativa de técnicas de seleção de protótipos

Dayvid Victor Rodrigues de Oliveira

Trabalho de Graduação

Recife
22 de novembro de 2011

Universidade Federal de Pernambuco
Centro de Informática

Dayvid Victor Rodrigues de Oliveira

Análise comparativa de técnicas de seleção de protótipos

Trabalho apresentado ao Programa de Graduação em Engenharia da Computação do Centro de Informática da Universidade Federal de Pernambuco como requisito parcial para obtenção do grau de Bacharel em Engenharia da Computação.

Orientador: *Prof. Dr. George Darmiton*

Recife
22 de novembro de 2011

*Eu dedico este trabalho a João Rodrigues de Silva, meu
avô.*

Agradecimentos

Agradeço a Deus.

What can I give back to God, for the blessings You pour out on me?
—BONO (Boston, 2001)

Resumo

RESUMO

Palavras-chave: PORTUGUES

Abstract

ABSTRACT

Keywords: INGLES

Sumário

1	Técnicas de Seleção de Protótipos	1
1.1	Motivação e Contextualização	1
1.1.1	Seleção de Protótipos	1
1.1.2	Bases Desbalanceadas	1
1.2	Objetivo	1
1.3	Estrutura do Trabalho	1
1.3.1	Sessões	1
1.3.2	Metodologia Utilizada	1
1.3.3	Bases de dados	1
2	Técnicas de Seleção de Protótipos	2
2.1	ENN	2
2.2	CNN	3
2.3	Tomek Links	4
2.4	LVQ	5
2.4.1	LVQ 1	5
2.4.2	LVQ 2.1	5
2.4.3	LVQ 3	5
2.5	SGP	5
2.6	SGP 2	5
2.7	CCNN	5

Lista de Figuras

Lista de Tabelas

Técnicas de Seleção de Protótipos

1.1 Motivação e Contextualização

Classificadores são (...)

*"I have a dream for the Web [in which computers] become capable of analyzing all the data on the Web, the content, links, and transactions between people and computers. A **Semantic Web** which should make this possible has yet to emerge, but when it does, the day-to-day mechanisms of trade, bureaucracy and our daily lives will be handled by machines talking to machines. The **intelligent agents** people have touted for ages will finally materialize."*Tim Berners-Lee

Tradução literal: *"Eu tenho um sonho para a Web [em que os computadores] tornam-se capazes de analisar todos os dados na Web, o conteúdo, links, e as transações entre pessoas e computadores. A **Web Semântica** que deve tornar isso possível ainda está para surgir; mas quando isso acontecer, os mecanismos dia-a-dia da burocracia do comércio e nossas vidas diárias serão tratados por máquinas falando com máquinas. Os **agentes inteligentes** que as pessoas têm falado por anos vão finalmente se concretizar."*Tim Berners-Lee

1.1.1 Seleção de Protótipos

1.1.2 Bases Desbalanceadas

1.2 Objetivo

1.3 Estrutura do Trabalho

1.3.1 Sessões

1.3.2 Metodologia Utilizada

1.3.3 Bases de dados

Técnicas de Seleção de Protótipos

Neste capítulo, serão mostradas as técnicas de seleção de protótipos abordadas neste trabalho. Cada uma das sessões abaixo abordará uma técnica, será mostrado o conceito da técnica, assim como o pseudo-código e as características de cada uma destas técnicas.

2.1 ENN

Edited Nearest Neighbor Rule[CPZ11] é uma técnica de seleção de protótipos puramente seletiva proposta por Wilson em 1976. De uma forma geral, esta técnica foi projetada para funcionar como um filtro de ruídos, ela elimina pontos na região de fronteira, região de alta susceptibilidade a erros, e com isso elimina ruídos.

Por atuar apenas na região de fronteira, esta técnica possui uma baixa capacidade de redução, deixando as instâncias que não se encontram na região de fronteira intactas, exceto pelos ruídos extremos.

Uma desvantagem desta técnica é que ela possui uma baixa capacidade de redução de elementos, visto que ela não elimina redundância.

Segue abaixo o algoritmo da execução do ENN e, logo após, alguns comentários sobre este algoritmo.

Algorithm 1 ENN

Require: *list*: uma lista

1. **for all** instância e_i da base de dados original **do**
 2. Aplique o KNN sobre e_i
 3. **if** e_i foi classificado erroneamente **then**
 4. salve e_i em *list*
 5. **end if**
 6. **end for**
 7. Remova da base de dados todos os elementos que estão em *list*
-

O valor de K pode variar de acordo com o tamanho da base de dados, porém, tipicamente, utiliza-se o valor de K=3. Tipicamente, O valor de K é inversamente proporcional a quantidade de instâncias que serão eliminadas, ou seja, para que o filtro elimine todos os possíveis ruídos, deve-se utilizar K=1.

Uma vantagem do ENN é que ele independe da ordem que a base de dados foi apresentada, ou seja, o ENN aplicado a uma base de dados, com o mesmo valor de K, sempre terá o mesmo resultado.

2.2 CNN

Condensed Nearest Neighbor [Har68] é uma técnica de seleção de protótipos puramente seletiva que tem como objetivo eliminar informação redundante. Diferentemente do ENN [CPZ11], o CNN não elimina instâncias nas regiões de fronteira, a técnica mantém estes elementos pois estes que "são importantes" para distinguir entre duas classes.

A ideia geral do CNN é encontrar o menor subconjunto da base de dados original que, utilizando o 1-NN, classifica todos os padrões da base de dados original corretamente. Fazendo isso, o algoritmo elimina os elementos mais afastados da região de indecisão, da fronteira de classificação.

O algoritmo do CNN será mostrado abaixo, e logo após, comentários a respeito do mesmo.

Algorithm 2 CNN

Require: *list*: uma lista

1. Escolha um elemento de cada classe *aleatoriamente* e coloque-os em *list*
 2. **for all** instância e_i da base de dados original **do**
 3. Aplique o KNN sobre e_i utilizando os elementos em *list* para treinamento
 4. **if** e_i foi classificado erroneamente **then**
 5. salve e_i em *list*
 6. **end if**
 7. **end for**
 8. Remova da base original todos os elementos que não estão em *list*
-

Podemos observar que este algoritmo possui uma abordagem diferente do ENN, por exemplo, pois ele começa com um conjunto mínimo de instâncias (uma de cada classe) e depois adiciona instâncias conforme a necessidade de mantê-las para que todos os elementos da base de dados original sejam classificados corretamente.

Uma coisa que pode-se observar no algoritmo, é a palavra *aleatoriamente*, o que significa que o CNN aplicado numa mesma base de dados com um mesmo valor de K para o KNN, nem sempre resulta nos mesmos protótipos. O primeiro fato para que isso ocorra é a seleção aleatória dos protótipos iniciais. Existem algumas adaptações para o CNN, onde os protótipos iniciais são escolhidos utilizando técnicas como o SGP[FHA07] para obter as instâncias mais centrais. Modificações no CNN são muito comuns [Tom76], porém, mesmo com estas modificações, o CNN ainda não é determinístico, pois a ordem em que as instâncias são classificadas afeta o resultado final.

Para o caso de estudo abordado neste trabalho, o CNN pode ser utilizado de forma adaptada. A adaptação consiste em manter todos os elementos da classe minoritária e o mínimo possível da classe majoritária. O próprio CNN se encarrega de remover os elementos redundantes da classe majoritária, assim, basta apenas selecionar todos os elementos da classe minoritária aos protótipos iniciais. Segue abaixo o algoritmo desta adaptação:

Com o algoritmo CNN adaptado para bases desbalanceadas, os elementos redundantes da classe majoritária são removidos, e todos os elementos da classe minoritária são mantidos. Esta adaptação do CNN irá reduzir a base de dados, ocasionando as vantagens de redução, e ainda reduzirá o desbalanceamento da base.

Algorithm 3 CNN para bases desbalanceadas

Require: *list*: uma lista

1. Coloque todos os elementos da classe minoritária em *list*
 2. **for all** instância e_i da base de dados original **do**
 3. Aplique o KNN sobre e_i utilizando os elementos em *list* para treinamento
 4. **if** e_i foi classificado erroneamente **then**
 5. salve e_i em *list*
 6. **end if**
 7. **end for**
 8. Remova da base original todos os elementos que não estão em *list*
-

2.3 Tomek Links

Mantendo a mesma linha do ENN, Tomek Links é uma técnica de seleção de protótipos puramente seletiva que elimina os elementos das regiões de fronteiras e instâncias com probabilidade de ser ruído. Tomek Links podem ser definidos da seguinte forma: Dadas duas instâncias e_i e e_j , o par (e_i, e_j) é chamado de Tomek Link se não existe nenhuma instância e_k , tal que, para todo e_k $dist(e_i, e_j) < dist(e_i, e_k)$ e $dist(e_i, e_j) < dist(e_j, e_k)$. Segue abaixo o algoritmo detalhado:

Algorithm 4 Seleciona Tomek Links

Require: *list*: uma lista

1. **for all** instância e_i da base de dados original **do**
 2. e_j = instância mais próxima de e_i
 3. **if** instância mais próxima de e_j for e_i **then**
 4. **if** classe de e_i for diferente da classe de e_j **then**
 5. salve o par (e_i, e_j) em *list*
 6. **end if**
 7. **end if**
 8. **end for**
 9. Retorne os pares em *list*, os Tomek Links
-

Os Tomek Links representam elementos da região de fronteira e prováveis ruídos, e a técnica de seleção de protótipos consiste em remover os Tomek Links da base de dados original. Segue abaixo o algoritmo da seleção de protótipos:

Enquanto o CNN remove os elementos que estão longe da região de indecisão, o Tomek Links remove os elementos que estão próximos desta região, o que causa uma maior separação entre as classes.

Observa-se facilmente que os Tomek Links pode remove todas as instâncias da fronteira, inclusive as instâncias da classe minoritária, assim sendo, uma possível adaptação dos Tomek Links é eliminar apenas os elementos das classes majoritárias. Nesse caso, ainda ocorreria uma separação entre as classes, mas apenas as instâncias da classe majoritária seriam removidos, diminuindo assim o nível de desbalanceamento. Segue abaixo o algoritmo desta adaptação:

Algorithm 5 Tomek Links

Require: *list*: uma lista

1. *list* = *SelecionaTomekLinks* da base original
 2. **for all** (e_i, e_j) em *list* **do**
 3. remova e_i da base original
 4. remova e_j da base original
 5. **end for**
 6. Retorne a base original filtrada
-

Algorithm 6 Tomek Links

Require: *list*: uma lista

1. *list* = *SelecionaTomekLinks* da base original
 2. **for all** (e_i, e_j) em *list* **do**
 3. **if** e_i for da classe majoritária **then**
 4. remova e_i da base original
 5. **end if**
 6. **if** e_j for da classe majoritária **then**
 7. remova e_j da base original
 8. **end if**
 9. **end for**
 10. Retorne a base original filtrada
-

Com esta adaptação, a classe minoritária é mantida, evitando o aumento do desbalanceamento ou a remoção por alta probabilidade de ruído.

2.4 LVQ**2.4.1 LVQ 1****2.4.2 LVQ 2.1****2.4.3 LVQ 3****2.5 SGP****2.6 SGP 2****2.7 CCNN**

Referências Bibliográficas

- [CPZ11] Ruiqin Chang, Zheng Pei, and Chao Zhang. A modified editing k-nearest neighbor rule. *JCP*, 6(7):1493–1500, 2011.
- [dSPC08] Cristiano de Santana Pereira and George D. C. Cavalcanti. Prototype selection: Combining self-generating prototypes and gaussian mixtures for pattern classification. In *IJCNN*, pages 3505–3510. IEEE, 2008.
- [EJJ04] Andrew Estabrooks, Taeho Jo, and Nathalie Japkowicz. A multiple resampling method for learning from imbalanced data sets. *Computational Intelligence*, 20(1):18–36, 2004.
- [FHA07] Hatem A. Fayed, Sherif Hashem, and Amir F. Atiya. Self-generating prototypes for pattern classification. *Pattern Recognition*, 40(5):1498–1509, 2007.
- [Har68] P. E. Hart. The condensed nearest neighbor rule. *IEEE Transactions on Information Theory*, 14:515–516, 1968.
- [HKN07] Jason Van Hulse, Taghi M. Khoshgoftaar, and Amri Napolitano. Experimental perspectives on learning from imbalanced data. In Zoubin Ghahramani, editor, *ICML*, volume 227 of *ACM International Conference Proceeding Series*, pages 935–942. ACM, 2007.
- [Koh86] Teuvo Kohonen. Learning vector quantization for pattern recognition. Report TKK-F-A601, Laboratory of Computer and Information Science, Department of Technical Physics, Helsinki University of Technology, Helsinki, Finland, 1986.
- [Koh88] Teuvo Kohonen. Learning vector quantization. *Neural Networks*, 1, Supplement 1:3–16, 1988.
- [PI69] E. A. Patrick and F. P. Fischer II. A generalization of the k-nearest neighbor rule. In *IJCAI*, pages 63–64, 1969.
- [Sav] Sergei Savchenko. <http://cgm.cs.mcgill.ca/>.
- [TC67] P.E. Hart T.M. Cover. Nearest neighbor pattern classification. *IEEE Transactions on Information Theory*, IT-13:21 – 27, 1967.
- [Tom76] I. Tomek. Two Modifications of CNN. *IEEE Transactions on Systems, Man, and Cybernetics*, 7(2):679–772, 1976.

- [WP01] G. Weiss and F. Provost. The effect of class distribution on classifier learning, 2001.