# 2º Trabalho de Inteligência Computacional - Data Mining

## Autor: Gustavo Rezende Silva, Orientador:Gina Maira Barbosa de Oliveira

<sup>1</sup>Faculdade de Computação Universidade Federal do Uberlândia (UFU) Uberlândia – MG – Brasil

gustavorezendesilva@hotmail.com, gina@ufu.br

**Resumo.** Este trabalho tem como intuito desenvolver um algoritmo genético para realizar uma mineração de dados em uma base dermatológica, de maneira a gerar regras de classificação para as diversas classes da doença erythemato-squamous presentes neste banço de dados.

**Palavras-Chave.** algoritmos genéticos, data mining, mineração de dados, inteligência computacional

### 1. Introdução

Atualmente a quantidade de banco de dados armazenados por supermecados, hospitais, empresas de marketing e muitos outros tipos de empresas vem aumentado significativamente. Com isso, técnicas para extrair e analisar informações de grandes bases têm sido desenvolvidas e este conjunto de ferramentas e técnicas é definida como mineração de dados (data mining).

Na área médica existe um grande interesse em desenvolver ferramentas para extrair informações de bancos de dados de pacientes com o intuito de auxiliar os médicos a diagnosticar os mesmos, devido ao grande potêncial do data mining este vem sido utilizado. Neste trabalho aplicou-se algumas técnicas de mineração de dados para extrair informações de um banco de dados dermatológico disponível no repositório *Machine Learning repository* da UCI.

Esta base contém informações sobre pacientes diagnosticados com alguma das variações da doença chamada *erythemato-squamous*. O diagnóstico dos diferentes tipos desta enfermidade é complicada, pois a variação entre elas é pequena.

O objetivo final da análise deste banco de dados é conseguir extrair uma regra para cada uma das variações da doença *erythemato-squamous* presentes na base. De formar que um paciente informando apenas os seus sintomas possa ser classificado com alguma das enfermidades, e este resultado pode ser utilizado pelo médico como referência para o diagnóstico.

#### 2. Desenvolvimento

Com o intuito de gerar regras para a base dermatológica citada na sec. 1 foi desenvolvido um algoritmo genético seguindo o artigo de Fidelis *et al.*(2000) e Miranda *et al.*(2003). Nestes trabalhos o indivíduo é uma abstração da regra de classificação de uma classe específica da doença *erythemato-squamous*, ou seja, o que define se um paciente é pré-diagnosticado com uma das enfermidades classificadas.

No banco de dados utilizado cada paciente é representado por um número e 34 atributos que representam os sintomas das doenças, por isso um indivíduo é composto por 34 atributos, 34 operadores matemáticos(=, ! =, < ou >=) e 34 pesos (fig. 1). De forma que um atributo relacionado com um operador define uma condição para um paciente ser classficado naquela

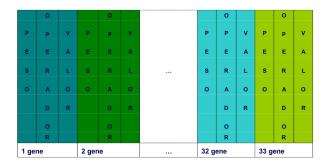


Figura 1. Representação de um indivíduo

regra (Ex.: atributo 2 >= 1), e ainda o peso é um número real entre 0 e 1 que representa a probabilidade de uma condição aparecer na regra.

Como parâmetros do algoritmo genético adotou-se uma população inicial de 50 individuos, taxa de cross over de 100%, taxa de mutação de 30%, número de gerações igual a 50. A atualização da população é feita ordenando os pais e filhos em ordem crescente de avaliação e pegando os 50 melhores, e configurou-se como valor máximo de peso para o qual uma condição aparece sendo 0,3.

Para treinar o algoritmo e gerar as regras foi utilizado 2/3 da base, em seguida após as regras serem geradas as mesmas foram testadas utilizando o 1/3 restante dos dados.

Outro aspecto importante é o método de avaliação adotado, neste trabalho para atribuir uma nota ao indivíduo foi utilizada a equação 3.

$$S_e = \frac{T_p}{T_p + F_n} \tag{1}$$

$$S_e = \frac{T_n}{T_n + F_n} \tag{2}$$

$$Av = S_e * S_n \tag{3}$$

Onde:

- $T_p$  Verdadeiros positivos
- $\bullet$   $T_n$  Verdadeiros negativos
- $F_p$  Falsos positivos
- $F_n$  Falsos negativos

#### 3. Análise dos Resultados

O algoritmo genético desenvolvido para mineração de dados foi executado 10 vezes para cada classe da doença, em seguida o melhor resultado de cada classe foi escolhido e estes estão representados na tabela 1.

Classe	Regras	Fitness	Fitness
das Doenças		Treinamento	Teste
1	thinning of the suprapapillary epidermis >= 1	0.929820	0.996290
	focal hypergranulosis < 3		
2	follicular papules = 0	0.829799	0.765162
	fibrosis of the papillary dermis < 1		
	exocytosis != 0		
	clubbing of the rete ridges < 1		
	vacuolisation and damage of basal layer != 3		
	band-like infiltrate < 1		
3	elongation of the rete ridges = 0	0.992163	0.994803
	band-like infiltrate >= 2		
4	family history != 1	0.811603	0.814022
	eosinophils in the infiltrate < 1		
	elongation of the rete ridges == 0		
	saw-tooth appearance of retes < 2		
	follicular horn plug != 2		
5	koebner phenomenon == 0	0.996662	0.991431
	oral mucosal involvement != 1		
	knee and elbow involvement < 3		
	PNL infiltrate != 2		
	fibrosis of the papillary dermis != 0		
	band-like infiltrate!= 3		
6	follicular papules >= 1	0.991299	0.986765
	perifollicular parakeratosis != 0		

Tabela 1. Resultados do data mining

Para escolher o melhor resultado foi levado em consideração a avaliação de treinamento e teste, além do número de condições resultantes. Uma vez que uma regra com muitas condições não representa o problema de forma genérica, apenas aquele conjunto de dados específicos.

Com a intenção de se ter uma ideia geral do resultado das regras obtidas, as médias dos melhores resultados de cada execução para a mesma regra foram calculadas e se encontram na tabela 2.

Classe	Média Fitness	Média Fitness
	Treinamento	Teste
1	0.932275	0.981523
2	0.825245	0.850473
3	0.976214	0.969933
4	0.850391	0.725871
5	0.896647	0.831133
6	0.939587	0.809525

Tabela 2. Média das avaliações de cada classe

Ao comparar os resultados obtidos neste trabalho com os de Fidelis *et al.*(2000) e Miranda *et al.*(2003), percebe-se que as avaliações das regras encontradas são bem próximas. Entretanto, as condições associadas a cada regra variam em quase todas elas.

Em seguida, alguns parâmetros do algoritmo genético foram alterados com o intuito de obter melhores resultados.

Alterando apenas o peso máximo permito o resultados encontrados foram diferentes (tabela 3). Para algumas classes como a 1 e 2 as avaliações foram melhores, já na classe 5 o fitness foi igual porém a quantidade de condições foi menor, a classe 3 permaneceu com as mesmas avaliações e a 4 e 6 pioraram.

Classe	Regras	Fitness	Fitness
das Doenças		Treinamento	Teste
1	clubbing of the rete ridges >= 1	0.955212	0.996290
	perifollicular parakeratosis < 1		
2	koebner phenomenon < 1	0.910755	0.923471
	polygonal papules == 0		
	fibrosis of the papillary dermis < 1		
	elongation of the rete ridges != 3		
	thinning of the suprapapillary epidermis $== 0$		
	perifollicular parakeratosis == 0		
3	PNL infiltrate == 0	0.992163	0.994803
	band-like infiltrate >= 2		
4	koebner phenomenon >= 1	0.892055	0.727187
	scalp involvement != 2		
	melanin incontinence < 1		
	exocytosis >= 1		
	clubbing of the rete ridges < 1		
	thinning of the suprapapillary epidermis != 2		
	polygonal papules < 2	0.996662	0.991431
	fibrosis of the papillary dermis != 0		
5	hyperkeratosis < 3		
	focal hypergranulosis < 1		
6	koebner phenomenon != 2	0.986932	0.986765
	perifollicular parakeratosis >= 1		

Tabela 3. Peso máximo = 0.2

As maiores dificuldades na reprodução dos modelos foi conseguir analisar os resultados e compara-los, e ainda entender que um número grande de condições não satisfaz o propósito da mineração de dados.

#### Referências

- Clay R. S. Miranda; OLIVEIRA, G. M. B.; SANTOS, J. B. (2003). Algoritmos genéticos aplicados em data mining para obtenção de regras simples e precisas. *SBAI2003: Simpósio Brasileiro de Automação Inteligente*.
- Fidelis M. V., L. H. S. and Freitas, A. A. (2000). Discovering comprehensible classification rules with a genetic algorithm. *Proc. Congress on Evolutionary Computation*, pages 805–810.
- OLIVEIRA, G. M. B. (2017). Apresentação em sala de aula algoritmos genéticos aplicados em data mining para ontenção de regras de classificação.