

Mineração de dados e adaptação de modelos de classificação de cobertura e uso da terra para imagem Worldview 2

Lídice Cabral Nascimento¹
Carla Bernadete Madureira Cruz¹

Universidade Federal do Rio de Janeiro – UFRJ ¹
Depto. de Geografia – Grupo de Sensoriamento Remoto ESPAÇO
Campus da Ilha do Fundão, prédio CCMN, bl. I, s/ 012 – CEP 21941-590 – Rio de Janeiro, RJ
{lidicecabral, carlamad}@gmail.com

Abstract. The studies about forest fragmentation have increased since the past 3 decades, together with the increase of the discussion about conservation and preservation. In this way, the landscape and land use maps are important tools for this analysis, as well as other remote sensing techniques. The object oriented analysis classifies the image according to patterns as texture, color, shape, and context. To identify which value and attribute is necessary to the classification, data mining technique has been used, because it looks for patterns inside the data. This technique helps to accelerate the process; however, it was necessary to adapt the model. In this way, the land use and land cover classification was made using the values that data mining software has provided. Afterwards, this classification was analyzed, and the verification was made by the confusion matrix, which was generated through the creation of random points shapefile in ArcGis 9.3, though this it was possible to evaluate the classification accuracy (58,89%). From these results, the model has been adapted. However, these adaptations have been done in the classes “water” and “forest. As result, 2 new classes have been created: shadow and forest2. Also, some values were changed and some attributes added, for example, brightness, to identify dark areas. In the end, the accuracy was 89,33%, however this result doesn’t show some errors which are still on the model.

Palavras-chave: data mining, object oriented analysis, landscape ecology, model adaptation, mineração de dados, análise orientada ao objeto, ecologia da paisagem, adaptação de modelos.

1. Introdução

As florestas, que antes da intervenção humana se distribuíam por vastas áreas foram reduzidas a vários arquipélagos de fragmentos florestais muito pequenos, bastante separados entre si. Na década de 1980 a preocupação com a conservação e preservação dos ecossistemas começou a se tornar uma questão tanto da sociedade em geral quanto do meio científico. Assim, deu-se início a discussões a respeito da fragmentação florestal e a melhor estratégia para se conservar a biodiversidade. Neste âmbito surgiram algumas teorias, como por exemplo, o SLOSS (sigla de *single large or several small reserves*), onde o debate estava em torno da conservação de uma grande área ou vários pequenos fragmentos com o objetivo de proteger mais espécies (WILCOX e MURPHY, 1985). A partir de então cada vez mais foram se desenvolvendo estudos acerca do tema “fragmentação florestal”. Dentro desta perspectiva, os mapeamentos de cobertura e uso da terra através de ferramentas de sensoriamento remoto vieram se notabilizando como metodologia para esta análise. Contudo, o uso destas ferramentas não é algo homogêneo, contendo diversos tipos de metodologias neste sentido, classificadas como automáticas, semi-automáticas e manuais.

Todas baseiam-se na resposta espectral dos alvos, contudo a análise orientada a objeto busca “simular técnicas de interpretação visual através da modelagem do conhecimento para a identificação de feições, baseada na descrição de padrões identificadores, tais como, cor, textura, métrica, contexto” (CRUZ et al, 2007). Para realizar esta modelagem é necessário identificar inicialmente os objetos, sendo assim necessário o processo de segmentação. Segundo Pinho et al (2005), “a segmentação multi-resolução parte do pressuposto de que as informações contextuais são importantes e, por isto, a interpretação de uma cena deve considerar não apenas a dimensão espectral, como também a dimensão espacial”. Por isso este processo deve ser moldado de acordo com a resolução da imagem e a escala de interesse (ANTUNES, 2003 apud PINHO et al, 2005). A partir daí realiza-se o processo de modelagem dos dados, onde se definem os descritores para a realização da classificação a partir de modelos booleanos ou fuzzy. Este tem sido um dos desafios

para a sua plena aplicação em mapeamentos de cobertura e uso da terra, pois o conhecimento sobre a melhor forma de caracterização de suas classes ainda é considerado insuficiente para que se alcance bons resultados de forma automática.

Neste sentido, a mineração de dados aparece como uma metodologia auxiliar, que otimiza o processo de escolha de descritores e limiares das classes, pois procura definir padrões, associações, mudanças, anomalias e estruturas significativas entre os dados. Essa técnica extrai informações de uma determinada base de dados, criada por meio da tarefa de classificação e da técnica de árvores de decisão. Ficando a cargo do analista apenas os processos de elaboração da rede hierárquica, segmentação e coleta de amostras (SOUZA, 2012).

Dois conceitos importantes que são utilizados durante o processo de mineração de dados devem ser ressaltados: tarefas e técnicas. As tarefas consistem na especificação do que se quer buscar no conjunto de dados, por exemplo: regras de associação, padrões sequenciais, análise de *outliers*, classificação e predição, análise de clusters (ou agrupamentos). As técnicas consistem na especificação de métodos que garantam como descobrir os padrões de interesse da pesquisa, como por exemplo: técnicas estatísticas, técnicas de aprendizado de máquina e técnicas baseadas em crescimento-poda-validação (SOUZA, 2012).

No presente trabalho, será utilizada a tarefa de classificação, que é um modelo constituído de regras que permitem classificar os objetos do banco de dados dentro de um número de classes pré-determinado. É criado a partir de um banco de dados de treinamento, cujos elementos são chamados de amostras ou exemplos (VINICI, 2005 apud SOUZA, 2012). A árvore de decisão é uma forma de classificador, ou seja, uma técnica, que vem sendo utilizado em diversas pesquisas na área das geociências, sendo considerada no presente estudo.

A mineração de dados seja uma técnica comprovadamente eficiente na classificação de uso e cobertura do solo, entretanto, tem suas limitações. Numa área de estudo tão heterogênea como o município de Silva Jardim, com coberturas naturais e antrópicas, presença de um grande corpo d'água, além de atividades rurais e urbanas, foi avaliado que para realizar uma classificação mais próxima da realidade seria necessário uma modelagem de dados complementar.

Neste trabalho será apresentada a modelagem realizada para as classes “Água” e “Floresta”, cujas respostas espectrais são mais homogêneas do que as outras classes presentes neste mapeamento. Assim, o objetivo do presente trabalho foi realizar a adaptação de modelos de uso e cobertura do solo para o município de Silva Jardim, a partir da mineração de dados, utilizando imagem Worldview2. É importante ressaltar que este mapeamento serve como um dos dados de entrada para uma análise cuja temática principal é a ecologia de paisagens. Isto será o fator determinante para a metodologia empregada neste mapeamento.

2. Metodologia de Trabalho

Para realizar a classificação de uso e cobertura do solo, foi utilizada a imagem do satélite Worldview 2, cujo dado original conta com pixel de 0,5m. Contudo, por conta da dificuldade de processamento, este dado foi reamostrado para 5m de tamanho do pixel, sem perda de qualidade na escala de análise desejada, 1:25.000. O mapa de uso e cobertura do solo contou com o auxílio da ferramenta WEKA de mineração de dados, proporcionando uma classificação mais eficiente e rápida. Esta ferramenta procura padrões nos dados nela inseridos, de forma que para a classificação orientada ao objeto, realizada no programa Definiens Developer 7.0, auxiliou na definição dos descritores das classes. Sendo os procedimentos descritos a seguir:

Inicialmente, definiu-se a legenda de interesse do mapeamento, que inicialmente foram as classes: água, solo exposto, urbano, floresta, pastagem e agricultura. Sabe-se que para este tipo de imagem a legenda pode ser mais detalhada. Contudo, o programa utilizado realiza a classificação de forma hierárquica. Por isso optou-se por uma legenda mais simplificada e assim, com o desenvolver da modelagem ela iria se complexificando, de forma que é importante melhorar a discretização de cada classe para aumentar o detalhamento. No programa Definiens Developer 7.0, realizou-se a

segmentação da imagem Worldview2 do município de Silva Jardim, com o parâmetro de escala 50, padrão de forma 0.1 e compacidade 0.5. Posteriormente, foram recolhidas amostras para cada classe, através da ferramenta de edição manual e exportadas para uma tabela (.csv) com os dados referentes aos descritores de média, brilho, desvio padrão, textura e NDVI. O produto gerado nesta etapa constitui-se numa tabela onde as linhas referem-se às instâncias das classes de cobertura e uso da terra e as colunas referem-se às respostas espectrais de cada amostra por descritor, como no exemplo da tabela 1:

Tabela1. Exemplo hipotético de tabela gerada pelo programa Definiens Developer 7.0

	Descritor 1	Descritor 2	Descritor 3	Descritor 4
Classe 1	Valor	valor	valor	valor
Classe 2	Valor	valor	valor	valor

Para esta tabela ser inserida no programa neozelandês WEKA de mineração de dados, é necessário que se atribuam códigos aos descritores, como na tabela abaixo (Tabela 2), que mostra uma parte da tabela utilizada para o presente mapeamento.

Tabela 2: exemplo do produto gerado pelo programa Definiens 7.0 que foi inserido no minerador de dados WEKA.

CLASSE	A1	A2	A3	A4	A5	A6	A7	A8	A9	A10	A11	A12	A13	A14	A15	A16
Floresta	71.48422	9.185193	7.275025	7.917636	6.827911	8.936582	2.073769	30.9139	50.06019	28.3605	176.6023	0.702058	18.52539	26.26131	17.03093	63.84832
Floresta	76.11444	9.083855	7.365117	8.089352	6.918194	8.445315	2.104781	31.54976	54.01554	29.3441	189.5484	0.714609	18.01941	25.68382	16.56116	57.26349
Floresta	71.82036	9.053092	7.340191	7.804029	6.842851	9.249879	1.936842	34.61076	51.55582	31.00507	170.1098	0.661873	15.94664	22.07055	14.92742	57.78371
Floresta	71.21253	9.017577	7.286193	7.869793	6.752527	8.849435	2.130137	30.65995	48.4985	26.9996	178.6921	0.707097	16.53619	24.13081	15.50477	62.60945
Floresta	70.45023	9.200468	7.396143	7.967702	6.934812	9.262429	1.983764	31.41946	52.20015	29.21233	168.969	0.686414	18.22668	24.79458	16.62867	62.22444
Floresta	75.86005	8.843608	7.231887	7.768308	6.617034	8.339464	2.260219	32.16256	48.62577	25.59579	197.0561	0.719372	15.61202	22.00764	14.47816	52.88169
Floresta	73.91886	8.865286	7.319773	7.798682	6.711584	8.693055	2.123895	33.16178	49.99521	27.76127	184.7572	0.69565	15.4431	21.62074	13.98101	54.58153
Floresta	69.49897	8.970523	7.168381	7.726112	6.696115	9.13931	2.049416	31.94141	47.21942	28.20137	170.6337	0.684646	14.87513	21.48355	14.2155	57.20511

Os descritores apresentados na tabela acima estão representados em códigos, sendo:

A1 – Brilho; A2 - *GLCM Entropy (quick 8/11) (all dir.)*; A3 - *GLCM Entropy (quick 8/11) Layer 1 (all dir.)*; A4 - *GLCM Entropy (quick 8/11) Layer 2 (all dir.)*; A5 - *GLCM Entropy (quick 8/11) Layer 3 (all dir.)*; A6 - *GLCM Entropy (quick 8/11) Layer 4 (all dir.)*; A7 - Maxima diferença ; A8 - Média banda 1; A9 - Média banda 2; A10 – Média banda 3; A11- Média banda 4 ; A12 – NDVI; A13 - Desvio padrão banda 1; A14 - Desvio padrão banda 2 ; A15 - Desvio padrão banda 3 ; A16 - Desvio padrão banda 4.

No WEKA, foi utilizado o classificador do tipo árvore de decisão, com o algoritmo J48. Este algoritmo produz uma sequência de regras partindo da raiz, criando-se sub-árvores até chegar às folhas, o que implica em uma divisão hierárquica em múltiplos subproblemas de decisão, os quais tendem a ser mais simples que o problema original (SOUZA, 2012).

A partir dessa árvore de decisão, foi criada a legenda hierárquica, e os limiares de cada descritor foram inseridos no programa Definiens Developer 7.0, utilizando classificadores booleanos. Esta opção foi feita seguindo a metodologia proposta por Martins e Fonseca (2009), uma vez que na árvore de decisão os limiares são definidos seguindo a lógica booleana. A partir de então foi gerada a classificação de cobertura e uso da terra.

A validação deste mapeamento foi realizada no programa ArcGis9.3 através da geração de 30 pontos aleatórios por classe; a partir daí cada um dos pontos foi verificado através de análise visual na imagem base (Worldview 2). Desta verificação, foi gerada uma matriz de confusão, que tem como resultado final o total percentual de acerto do mapeamento e por classes, podendo-se identificar confusões entre as classes. A partir daí, identificou-se as duas classes mais homogêneas em termos de resposta espectral, que conseqüentemente, tiveram mais acertos na classificação e iniciou-se a adaptação da modelagem realizada pelo minerador de dados, alterando limiares e inserindo descritores, sem edição manual. Esta modelagem foi avaliada utilizando a mesma metodologia empregada na validação do mapeamento utilizando a mineração de dados, contudo, houve uma generalização de classes (solo exposto e urbano – Outros2/ agricultura e pastagem -

outros3). Esta classe de “outros” entra na modelagem utilizando o critério inverso de similaridade, e está inserida na legenda de forma hierárquica para deixar de fora da modelagem classes que não são de interesse no momento. Posteriormente todas as classes serão modeladas.

3.Resultados e Discussão

Foram inseridas no minerador de dados um total de 142 amostras, divididas em 6 classes, a distribuição de amostras por classe está discretizada na Figura 2. Embora se tenha buscado coletar um número de amostras igualitário por cada classe, ainda houve uma disparidade entre as classes, com um desvio padrão de 5,64.

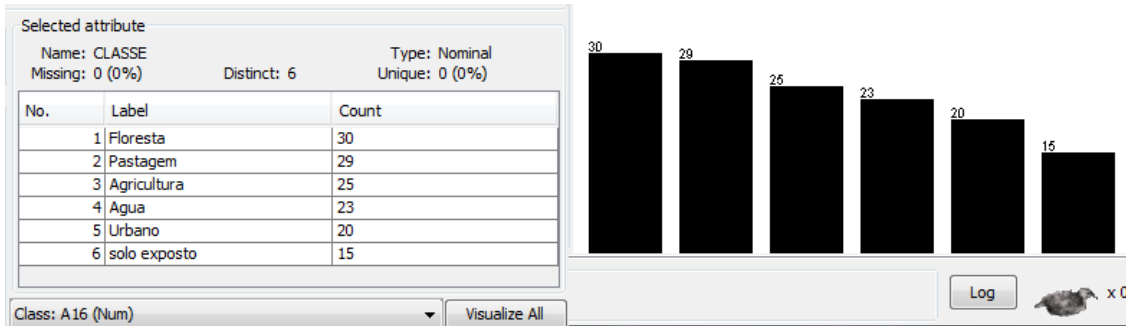


Figura 2. Quantidade de amostras por classe de uso e cobertura do solo inserida no minerador de dados WEKA.

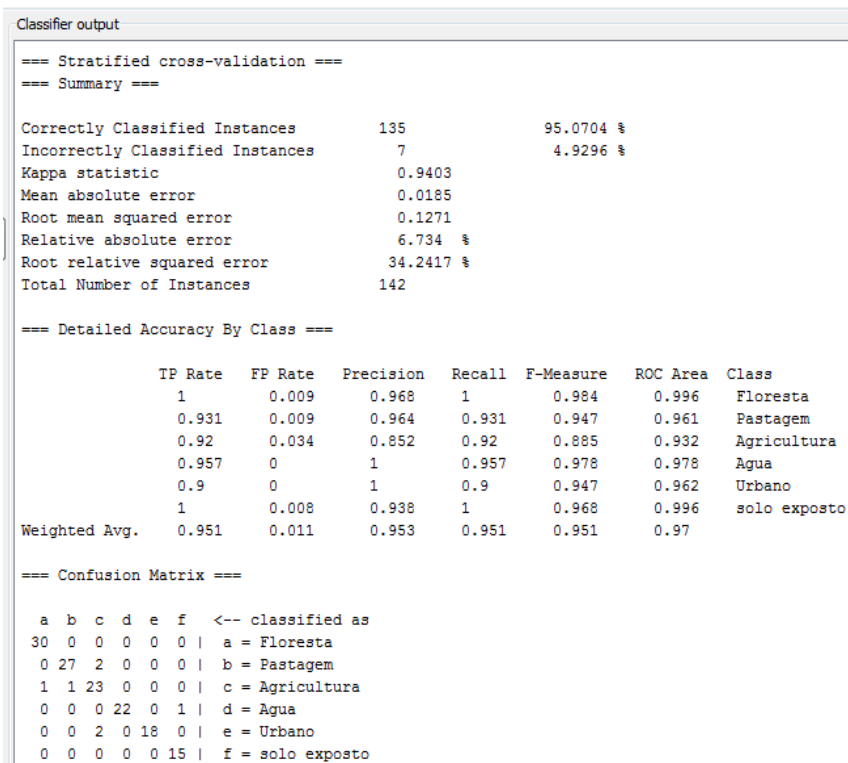


Figura 3. Relatório de desempenho de classificação do minerador de dados

Ao rodar o classificador J48, o programa oferece um relatório de desempenho (Figura 3), onde gera o índice kappa, que no caso em análise foi de 0.9403. Ou seja, uma acurácia muito boa. Na matriz de confusão, as três classes que apresentaram maiores confusões foram: pastagem, agricultura e urbano. Todas com apenas duas amostras erroneamente classificadas. É interessante notar que ambas as classes de pastagem e urbano tiveram suas duas amostras erroneamente classificadas como Agricultura.

Esse processo de classificação teve como produto final a árvore de decisão apresentada abaixo (Figura 4), onde o primeiro nó considerou o NDVI como o descritor diferenciador entre as classes de vegetação e as outras. A Figura 5 apresenta o mapa de cobertura e uso da terra gerado a partir do resultado da classificação dos parâmetros do minerador de dados.

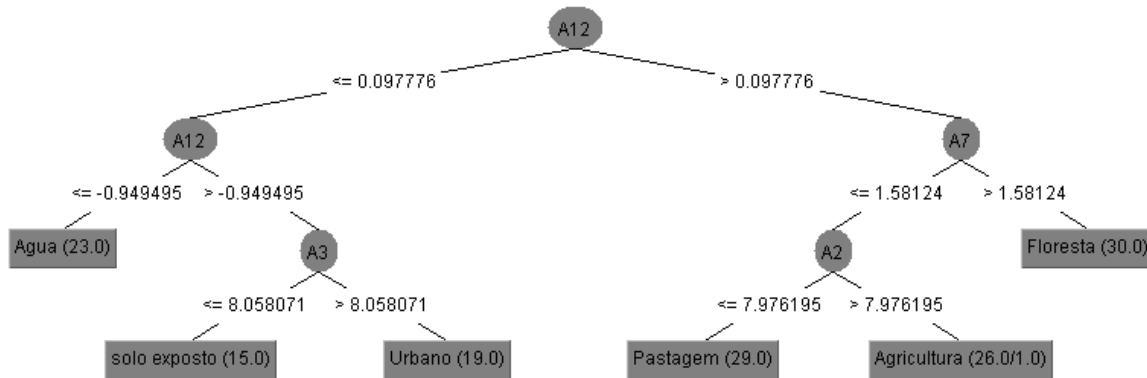


Figura 4: Árvore de decisão gerada pelo minerador de dados, onde os nós (círculos) representam os códigos dos descritores, e as linhas contêm os limiares. A12 (NDVI) / A7 (Máxima diferença) / A3 (Entropia GLCM, banda 1) / A2 (Entropia GLCM, todas direções).

A partir da árvore de decisão, que forneceu os descritores e os limites, foi gerada a classificação de cobertura e uso da terra do município de Silva Jardim, RJ. Cujo mapa final encontra-se na Figura 5. Embora o classificador do minerador de dados tenha encontrado um índice kappa que indica ótima acurácia (Figura 3), ao analisarmos a matriz de confusão (Tabela 3), podemos perceber que houve uma taxa de acerto de 58,89%, principalmente nas classes de água e floresta, que possuem resposta espectral mais homogênea, se comparada às outras classes. Esta diferença de acurácia deve-se à metodologia empregada para o cálculo do índice kappa. Enquanto a validação do minerador só leva em consideração as amostras, a validação do mapeamento considera toda a área do município. A classe classificada como agricultura, teve 100% de confusão com a classe de pastagem, segundo a matriz de confusão do mapeamento (tabela 3), de forma que este é um resultado completamente insatisfatório. Outra classe que apresentou alta taxa de confusão foi a de urbano, onde 16 amostras foram consideradas Pastagem, isto se deu pela própria natureza da configuração espacial desta classe, uma vez que não se apresenta como áreas urbanas de alta intensidade, possuindo muita vegetação e algumas áreas não construídas. É importante ressaltar que o município de Silva Jardim é essencialmente rural.

Tabela 3. Matriz de confusão da classificação gerada através da mineração de dados

Dados da Classificação (Mineração de Dados)	Classes	Dados da amostragem						Total das linhas
		Agricultura	Água	Floresta	Pastagem	Solo Exposto	Urbano	
Agricultura		0	0	0	30	0	0	30
Água		0	30	0	0	0	0	30
Floresta		0	0	27	3	0	0	30
Pastagem		1	0	6	23	0	0	30
Solo Exposto		0	3	0	7	20	0	30
Urbano		0	2	0	16	6	6	30
Soma das colunas		1	35	33	79	26	6	180
Taxa de acerto total:								58.89%

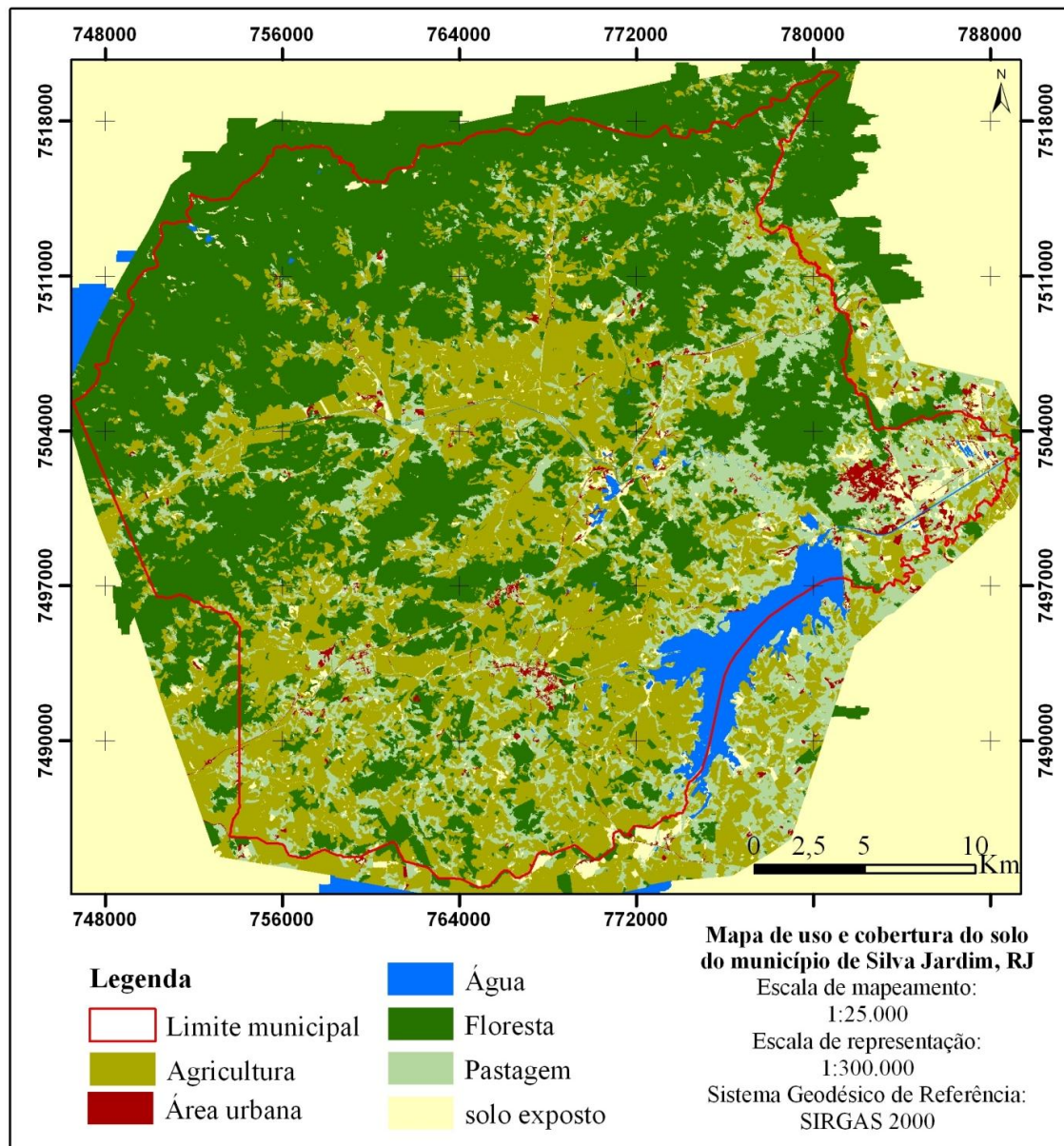


Figura 5. Mapa de cobertura e uso da terra gerado com o auxílio de ferramentas de mineração de dados

Assim, observaram-se inicialmente quais foram as confusões apresentadas pelo primeiro mapeamento, referentes à classe Água. Algumas áreas de sombra foram confundidas com a classe Água, este é uma confusão comum nas modelagens, afinal muitas vezes tem a resposta espectral semelhante em algumas bandas. Assim, optou-se por criar a classe "sombra", subordinada a classe Água.

Desta maneira, foi realizada a modelagem dos dados alterando os limiares dos descritores, assim, consideraram-se todos os valores menores que -0,68 (NDVI) como pertencentes a classe água. Dentro desta classe, separou-se água de sombra, classificando a classe sombra a partir do descritor de brilho, onde se considerou todos os valores menores que 9. A classe de água foi estabelecida a partir do critério inverso de similaridade.

Dentro da categoria Vegetação inicialmente havia 3 classes: Floresta, Agricultura e Pastagem. No primeiro nível de classificação a classe floresta foi separada das duas outras. Contudo, o resultado encontrado subestimava as áreas de floresta, uma vez que esta classe, embora considerada

homogênea, apresente grandes heterogeneidades quanto a sua resposta espectral. De forma que existem áreas florestadas mais escuras, sob influência de sombras, e mais claras. Assim, optou-se por, inicialmente, alterar os limites do descritor Máxima diferença, considerando todos os segmentos com valores superiores a 1.4. Ainda assim, o resultado não foi satisfatório, de forma que foi introduzido mais um descritor de textura (Entropia GLCM todas as direções), cujos valores superiores a 7.89 foram considerados inseridos na classe Floresta. Ainda assim, muitos segmentos mais escuros desta classe não foram inseridos, de forma que se optou por mais um nível de classificação de florestas (classe “floresta2). Assim, utilizou-se o descritor de brilho, considerando todos os segmentos com valores inferiores a 51 como pertencentes à classe de floresta. A legenda final obedeceu à hierarquia apresentada abaixo (Figura 6):

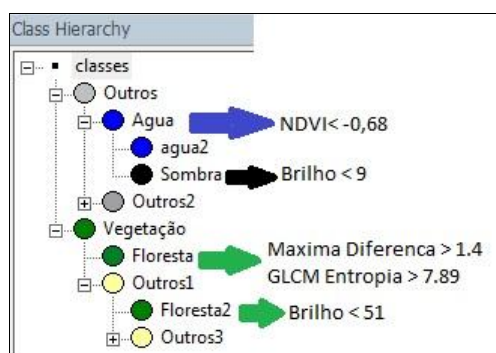


Figura 6. Hierarquia da legenda e os limites e descritores utilizados na classificação.

As classes Agricultura, Pastagem, Solo exposto e Urbano, não foram modeladas, de forma que neste mapeamento as duas primeiras encontram-se contempladas na classe “outros 3” e as duas últimas na classe “outros 2”. Embora tenhamos realizado uma modelagem da classe floresta, a fim de abarcar segmentos com respostas espectrais mais escuras, alguns segmentos mais claros, em áreas mais abertas não foram contemplados. Também se apresentam segmentos que não se inserem na classe floresta, nem de pastagem, pois são áreas de regeneração florestal. Desta maneira, fica indicado que seria necessário inserir mais algumas classes, como por exemplo, refúgios vegetacionais e vegetação secundária.

A tabela 4 apresenta a matriz de confusão, com 89,33% de acerto. Esta é uma taxa alta, contudo, é possível identificar a grande confusão entre as classes outros2 e outros 3, ainda não modeladas. Apesar da alta taxa de acerto nas classes modeladas, este resultado esconde alguns erros presentes no mapeamento. Ainda é possível identificar áreas de pastagem classificadas como floresta, e vice-versa, constituindo erros de omissão e comissão. As classes de água e sombra, entretanto, apresentam grande índice de acertos, por conta da sua alta homogeneidade das respostas espectrais.

Tabela 4. Matriz de confusão da classificação realizada a partir da modelagem das classes de água, sombra e floresta.

Dados da Classificação	Dados da amostragem						Soma das linhas
	Classes	Água	Floresta	Outros2	Outros3	Sombra	
Água		30	0	0	0	0	30
Floresta		0	28	0	2	0	30
Outros2		1	0	16	12	1	30
Outros3		0	0	0	30	0	30
Sombra		0	0	0	0	30	30
Soma das colunas		31	28	16	44	31	150
Total de acertos							89,33%

4. Conclusões

A mineração de dados mostrou-se ser uma técnica eficiente e útil para a classificação de áreas heterogêneas, com grande diversidade de alvos e respostas espectrais. Contudo, para que se obtenha um modelo mais próximo da realidade é necessário realizar modelagem adicional sobre os resultados obtidos. No atual mapeamento optou-se por iniciar o processo com uma legenda pouco complexa, com muitas generalizações, contudo, esta decisão mostrou-se equivocada, pois dificultou a modelagem estatística realizada pelo minerador de dados. Desta maneira, conclui-se que o ideal seria iniciar o processo de classificação com uma legenda mais complexa, e se necessário, realizar a generalização posterior.

É preciso também realizar uma avaliação da segmentação gerada, uma vez que este é o procedimento mais importante da análise baseada em objeto, pois se os alvos a serem classificados forem muito heterogêneos, dificulta a classificação, aumentando a possibilidade de erros.

Referências Bibliográficas

Cruz, C. B. M, et al. Classificação orientada a objetos no mapeamento dos remanescentes da cobertura vegetal do bioma Mata Atlântica, na Escala 1:250.000. In: Simpósio Brasileiro de Sensoriamento Remoto (SBSR), 13, 2007, Florianópolis, **Anais...**, São Jose dos Campos: INPE, 2007.

Jenness, J. Repeating shapes for ArcGIS. Jenness Enterprises, 2012. Available at: http://www.jennessent.com/arcgis/repeat_shapes.htm.

Martins, V.A., Fonseca, L.M.G. Classificacao de uso do solobaseada na analise orientada a objeto e mineração de dados utilizando imagem SPOT/HRG-5. In: Simpósio Brasileiro de Sensoriamento Remoto (SBSR), 14, 2009, Natal, **Anais...**, São Jose dos Campos: INPE, 2009.

Pinho,C. M.D., Feitosa,F.F., Kux,H. Classificação automática de cobertura do solo urbano em imagem IKONOS: comparação entre a abordagem pixel-a-pixel e orientada a objetos. In: Simpósio Brasileiro de Sensoriamento Remoto (SBSR), 12, 2005, Goiânia, **Anais...**, São Jose dos Campos: INPE, 2005.

Souza,E.M.F.R. Diferenças nas respostas espectrais de floresta em encosta por meio de imagem hiperspectral. Tese de doutorado Universidade Federal Fluminense, Niterói, 2012.

Wilcox, B. A.; Murphy, D. D. Conservation Strategy: The Effects of Fragmentation on Extinction. **The American Naturalist**, v. 125, n. 6, p. 879-887, 1985.