



Ministério da  
Ciência e Tecnologia



sid.inpe.br/mtc-m19/2010/10.26.11.36-RPQ

## **CLASSIFICAÇÃO DA COBERTURA DO SOLO URBANO UTILIZANDO IMAGENS IKONOS II E DADOS LiDAR**

Leonardo Rodrigues de Deus

Relatório final da disciplina Princípios e Aplicações de Mineração de Dados (CAP-359) do Programa de Pós-Graduação em Computação Aplicada, ministrada pelo professor Rafael Santos.

URL do documento original:

<<http://urlib.net/8JMKD3MGP7W/38FQ4QL>>

INPE  
São José dos Campos  
2010

**PUBLICADO POR:**

Instituto Nacional de Pesquisas Espaciais - INPE

Gabinete do Diretor (GB)

Serviço de Informação e Documentação (SID)

Caixa Postal 515 - CEP 12.245-970

São José dos Campos - SP - Brasil

Tel.:(012) 3208-6923/6921

Fax: (012) 3208-6919

E-mail: pubtc@sid.inpe.br

**CONSELHO DE EDITORAÇÃO E PRESERVAÇÃO DA PRODUÇÃO INTELLECTUAL DO INPE (RE/DIR-204):****Presidente:**

Dr. Gerald Jean Francis Banon - Coordenação Observação da Terra (OBT)

**Membros:**

Dr<sup>a</sup> Inez Staciarini Batista - Coordenação Ciências Espaciais e Atmosféricas (CEA)

Dr<sup>a</sup> Maria do Carmo de Andrade Nono - Conselho de Pós-Graduação

Dr<sup>a</sup> Regina Célia dos Santos Alvalá - Centro de Ciência do Sistema Terrestre (CST)

Marciana Leite Ribeiro - Serviço de Informação e Documentação (SID)

Dr. Ralf Gielow - Centro de Previsão de Tempo e Estudos Climáticos (CPT)

Dr. Wilson Yamaguti - Coordenação Engenharia e Tecnologia Espacial (ETE)

Dr. Horácio Hideki Yanasse - Centro de Tecnologias Especiais (CTE)

**BIBLIOTECA DIGITAL:**

Dr. Gerald Jean Francis Banon - Coordenação de Observação da Terra (OBT)

Marciana Leite Ribeiro - Serviço de Informação e Documentação (SID)

Deicy Farabello - Centro de Previsão de Tempo e Estudos Climáticos (CPT)

**REVISÃO E NORMALIZAÇÃO DOCUMENTÁRIA:**

Marciana Leite Ribeiro - Serviço de Informação e Documentação (SID)

Yolanda Ribeiro da Silva Souza - Serviço de Informação e Documentação (SID)

**EDITORAÇÃO ELETRÔNICA:**

Vivéca Sant´Ana Lemos - Serviço de Informação e Documentação (SID)



Ministério da  
Ciência e Tecnologia



sid.inpe.br/mtc-m19/2010/10.26.11.36-RPQ

## CLASSIFICAÇÃO DA COBERTURA DO SOLO URBANO UTILIZANDO IMAGENS IKONOS II E DADOS LiDAR

Leonardo Rodrigues de Deus

Relatório final da disciplina Princípios e Aplicações de Mineração de Dados (CAP-359) do Programa de Pós-Graduação em Computação Aplicada, ministrada pelo professor Rafael Santos.

URL do documento original:

<<http://urlib.net/8JMKD3MGP7W/38FQ4QL>>

INPE  
São José dos Campos  
2010



## RESUMO

A quantidade de dados que são produzidos nas mais diversas áreas de conhecimento e o conseqüente armazenamento destes dados em repositórios cada vez com mais espaço de disco, leva à necessidade de novas teorias e ferramentas computacionais que auxiliem a mão de obra humana no processo de descobrir conhecimento ao se analisar estes dados. Neste sentido a mineração de dados tem sido empregada nos mais diversos campos de pesquisa para se extrair conhecimento a partir dos grandes volumes de dados disponíveis. O sensoriamento remoto é um campo que tem se beneficiado com a utilização de técnicas de mineração de dados, uma vez que estas podem aumentar o potencial de análise e aplicação dos dados deste tipo. Visando utilizar técnicas de descoberta de conhecimento em banco de dados, este trabalho tem por objetivo utilizar dados de sensoriamento remoto, compostos por imagens do sensor *IKONOS II* e dados *LiDAR*, para classificar a cobertura do solo urbano, analisando se as informações acrescentadas pelos dados *LiDAR* possibilitam a obtenção de uma classificação mais precisa, a partir do uso de um software de mineração de dados para predição da classificação. Foram empregadas três classificadores diferentes implementados no software Weka, o IBK, o J48 e MLP. Os resultados mostraram que para todos os classificadores, a utilização de imagens *IKONOS II* em conjunto com dados *LiDAR* possibilitou que a cobertura do solo urbano fosse classificada com maior precisão. A contribuição dos dados *LiDAR* para classificar a cobertura do solo urbano esta na informação de altura dos alvos que sistema fornece, e a utilização destes dados propiciou uma melhora na precisão das instâncias classificadas corretamente da ordem de 5% para a área de estudo.

Palavras-chaves: mineração de dados, sensoriamento remoto, cobertura do solo urbano, classificação do solo urbano, imagens IKONOS II, dados LiDAR.

## ABSTRACT

The amount of data that are produced in several areas of knowledge and the consequent storage of these data in repositories with increasing disk space, leads to the need for new theories and computational tools to assist the labor human in the process of discovery knowledge when analyzing these data. In this sense, data mining has been used in various fields of research in order to extract knowledge from large volumes of data. Remote sensing is a field that has benefited from the use of data mining techniques, since these can increase the potential for analysis and application of such data. Aiming to use techniques of knowledge discovery in databases, this study aims to use remote sensing data, comprising images of the IKONOS sensor II and LiDAR data to classify land cover urban, evaluating whether the information added by LiDAR data help to obtain a more precise classification, from the use of a data mining software for predicting classification. We used three different classifiers implemented in Weka software, IBK, the J48 and MLP. The results showed that for all classifiers, the use of IKONOS II images in association with LiDAR data enabled the urban land cover be classified more accurately. The contribution of LiDAR data to classify the coverage of urban land is information of height the targets that the system provides and the use of these data led to an improvement in the accuracy of instances correctly classified in the order of 5% for the study area.

Keywords: data mining, remote sensing, urban land cover, classification of urban land, IKONOS II images, LiDAR data.

## LISTA DE FIGURAS

Figura 1 - Visão geral das etapas que compõem o processo de KDD .....	2
Figura 2 - Delimitação da área de Estudo .....	6
Figura 3: Localização geográfica de Uberlândia .....	7
Figura 4: Imagem <i>IKONOS II</i> da área de estudo .....	8
Figura 5: Imagem MDA da área de estudo .....	9
Figura 6: Imagem Intensidade da área de estudo .....	10
Figura 7a: Amostras selecionadas para os 12 alvos trabalhados .....	11
Figura 7b: Amostras selecionadas para os 12 alvos trabalhados .....	12
Figura 8: Amostra do arquivo ARFF utilizado no software WEKA .....	13
Figura 9: Árvore de decisão gerada utilizando o terceiro grupo de atributos .....	28

## LISTA DE TABELAS

	<b>Pág.</b>
Tabela 1 - Resultados Classificador IBK .....	18
Tabela 2: Matriz de confusão da classe Telhado Azul – Classificador IBK .....	20
Tabela 3: Matriz de confusão da classe Asfalto – Classificador IBK .....	21
Tabela 4: Matriz de confusão da classe Solo Exposto – Classificador IBK .....	21
Tabela 5: Matriz de confusão da classe Edifício – Classificador IBK .....	22
Tabela 6: Matriz de confusão da classe Telhado Branco – Classificador IBK .....	22
Tabela 7: Matriz de confusão da classe Telhado Escuro – Classificador IBK .....	23
Tabela 8: Matriz de confusão da classe Telhado Marrom – Classificador IBK .....	23
Tabela 9: Matriz de confusão da classe Telhado Marrom Escuro – Classificador IBK .....	24
Tabela 10: Matriz de confusão da classe Vegetação – Classificador IBK .....	24
Tabela 11: Resultados Classificador J48 .....	25
Tabela 12: Resultados Classificador MLP .....	29
Tabela 13: Precisão instâncias classificadas corretamente para o terceiro grupo de Atributos .....	31



## LISTA DE QUADROS

	<b>Pág.</b>
Quadro 1: Matriz de confusão – Segundo grupo de atributos – Classificador J48 .....	26
Quadro 2: Matriz de confusão – Primeiro grupo de atributos – Classificador J48 .....	27
Quadro 3: Matriz de confusão – Terceiro grupo de atributos – Classificador J48 .....	27
Quadro 4: Matriz de confusão – Classificador MLP – 1 camada escondida .....	30
Quadro 5: Matriz de confusão – Classificador MLP – 50 camadas escondidas .....	31

## SUMÁRIO

	<b>Pág.</b>
RESUMO .....	iii
ABSTRACT .....	iv
LISTA DE FIGURAS .....	v
LISTA DE TABELAS .....	vi
LISTA DE QUADROS .....	vii
1. INTRODUÇÃO .....	1
2. OBJETIVOS .....	4
3. JUSTIFICATIVA .....	5
4. MATERIAIS E MÉTODOS .....	6
4.1 Área de Estudo .....	6
4.2 Dados Utilizados .....	7
4.3 Procedimentos Metodológicos .....	11
5. RESULTADOS E DISCUSSÕES .....	18
5.1 Resultados Classificador IBK .....	18
5.2 Resultados Classificador J48 .....	25
5.3 Resultados Classificador MultiLayerPerceptron (MLP) .....	29
6. CONSIDERAÇÕES FINAIS .....	33
7. REFERÊNCIAS BIBLIOGRÁFICAS .....	35

## **1. Introdução**

O estágio atual de desenvolvimento da sociedade é tratado por muitos como a Era da Informação, onde a demanda por informação e conhecimento é muito grande, promovendo a possibilidade de se acompanhar os acontecimentos em tempo real.

Isto leva à produção de uma quantidade enorme de dados, sobre os mais diversos assuntos, e que são armazenados em sistemas de bancos de dados variados. Fayyad et al. (1996) ressalta que, por meio de uma variedade de campos, os dados têm sido coletados e acumulados em um ritmo muito acelerado.

Esta grande quantidade de dados digitais disponíveis para o desenvolvimento de uma série de pesquisas, leva à necessidade urgente de uma nova geração de teorias e ferramentas computacionais para auxiliar a mão de obra humana na extração de informações úteis, ou seja, na extração de conhecimento, a partir dos volumosos repositórios de dados (FAYYAD et al., 1996).

A demanda por novas formas de se trabalhar com estas grandes bases de dados levou ao surgimento de uma área de conhecimento, no campo da tecnologia da informação, chamada de Descoberta de Conhecimento em Banco de Dados, termo amplamente conhecido na literatura como KDD (Knowledge Discovery in Databases).

A descoberta de conhecimento em banco de dados é um conceito abrangente que, segundo Fayyad et al. (1996), refere-se a todo o processo de descobrir conhecimento útil, e que não está explícito nos dados a serem trabalhados.

Han e Kamber (2006) colocam que esta descoberta de conhecimento consiste de uma sequência iterativa que envolve as seguintes etapas:

1. Limpeza dos dados, para remover ruídos e dados inconsistentes;
2. Integração de dados, permitindo que dados de múltiplas fontes sejam combinados;

3. Seleção dos dados, onde os relevantes para a tarefa de análises são recuperados dos sistemas de armazenamento;
4. Transformação dos dados, onde estes são transformados ou consolidados em formas apropriadas para realizar operações de mineração ou agregação, por exemplo;
5. Mineração de dados, etapa a qual é um processo essencial, onde métodos inteligentes são aplicados com o objetivo de extrair padrões dos dados;
6. Avaliação dos padrões, para identificar os padrões verdadeiramente interessantes para determinada análise, representando a base do conhecimento em algumas medidas;
7. Apresentação do conhecimento, onde técnicas de visualização e representação do conhecimento são utilizadas para apresentar o novo conhecimento obtido pelo usuário.

A Figura 1 mostra as etapas de evolução do processo de descoberta do conhecimento em banco de dados segundo a visão de Fayyad et al. (1996).

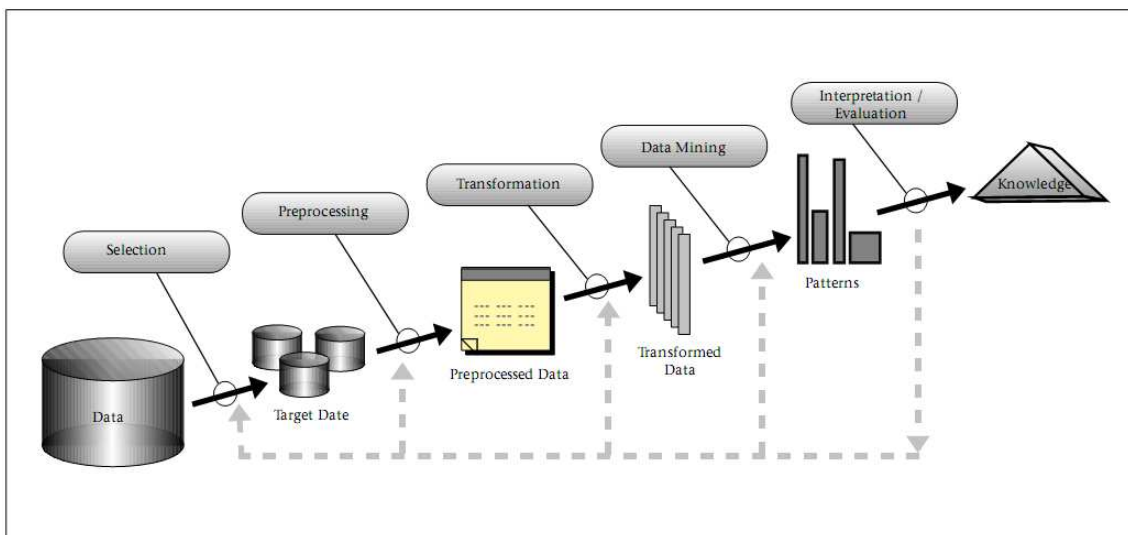


Figura 1: Visão geral das etapas que compõem o processo de KDD.  
Fonte: Fayyad et al. (1996)

Dentre estas etapas, a mineração de dados é a aplicação de algoritmos específicos para extração de padrões a partir dos dados (FAYYAD et al., 1996). A mineração de dados se refere à extração, ou mineração, de conhecimento a partir de grandes quantidades de dados (HAN e KAMBER, 2006).

Han e Kamber (2006) ressaltam que a mineração de dados tem atraído recentemente, uma grande atenção na indústria da informação e na sociedade em geral, devido a disponibilidade de enormes quantidades de dados e a necessidade de transformar estes dados em informações úteis e em conhecimento. E colocam que estas informações e o ganho de conhecimento podem ser usados para aplicações que vão desde análise de marketing, detecção de fraudes e conservação de clientes, até controle de produção e pesquisa científica.

Uma área de pesquisa que tem empregado ferramentas de mineração de dados e obtido bons resultados é área de sensoriamento remoto.

A quantidade de satélites imageadores que estão disponíveis hoje registram um volume enorme de imagens da superfície terrestre, com diversas resoluções espaciais, espectrais, radiométricas e temporais, possibilitando a geração de imensos repositórios de dados de imagens.

Com a grande quantidade de imagens diversos trabalhos podem ser realizados e a utilização de ferramentas de mineração de dados pode aumentar o potencial de análise e aplicações de dados de sensoriamento remoto (Korting et al., 2009).

Os estudos de sensoriamento remoto em áreas urbanas podem se beneficiar bastante da aplicação de ferramentas de mineração de dados, uma vez que nestas áreas encontra-se grande variedade de alvos sendo difícil a distinção entre vários deles, exigindo a utilização de técnicas mais apuradas para extração de informações.

## 2. Objetivos

O problema que este estudo aborda é o de classificar a cobertura do solo urbano a partir da utilização de dados de sensoriamento remoto.

Assim, o objetivo principal deste estudo consiste em utilizar dados de sensoriamento remoto, compostos por imagens do sensor *IKONOS II* e dados *LiDAR*, para classificar a cobertura do solo urbano, analisando se as informações acrescentadas pelos dados *LiDAR* possibilitam uma classificação mais precisa, a partir do uso de um software de mineração de dados para predição da classificação.

Além deste objetivo principal, espera-se também verificar a eficiência de diferentes tipos de classificadores implementados no software utilizado, para predição das classes escolhidas.

### 3. Justificativa

As imagens de sensoriamento remoto são uma das fontes de informações mais utilizadas para produção de mapas detalhados do espaço urbano, e estas informações são extraídas por meio de processos de interpretação visual das imagens ou por métodos automáticos de classificação.

Para grandes áreas, a interpretação visual não é o procedimento adequado, pois é um processo lento e caro, uma vez que demanda muita mão-de-obra qualificada para função. E para a classificação automática, as imagens de alta resolução implicam em algumas limitações importantes.

As imagens de alta resolução possuem baixa resolução espectral, trabalhando apenas na faixa de comprimento de onda do visível e infravermelho próximo, o que segundo Pinho et al. (2007), dificulta a distinção de uma série de alvos urbanos com comportamento espectral semelhante para estes comprimentos de onda, como é o caso de ruas pavimentadas com asfalto e edificações com cobertura de amianto escuro.

Dessa forma se faz necessário que novos métodos sejam desenvolvidos para auxiliar nos processos de classificação automática das imagens.

Nesse sentido o uso combinado de imagens *IKONOS II*, por exemplo, com dados *LiDAR*, pode propiciar que a classificação da cobertura do solo urbano seja feita com mais precisão, uma vez que os dados *LiDAR* fornecem informação de altura dos objetos, o que pode favorecer a discriminação entre alvos que apresentam respostas espectrais semelhantes, mas tem diferença de altura, como alguns tipos de telhados e o solo exposto.

## 4. Materiais e Métodos

### 4.1 Área de Estudo

A área de estudo deste trabalho consiste de uma porção central da cidade de Uberlândia-MG, composta pelos bairros Centro e Fundinho, e por parte dos bairros Martins, Osvaldo Rezende e Tabajaras (Figura 2). Esta área de estudo coincide com a região da cidade que apresenta a maior aglomeração de edificações verticalizadas (TOMÁS, 2010).

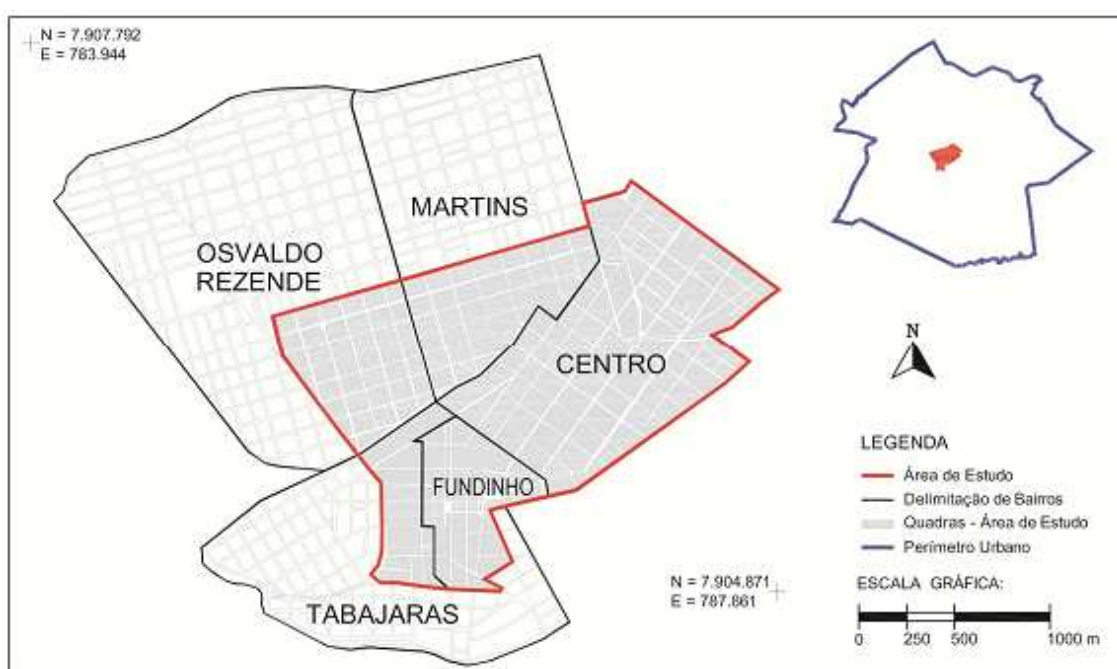


Figura 2: Delimitação da área de Estudo.  
Fonte: Tomás, 2010.

O município de Uberlândia localiza-se na região do Triângulo Mineiro, no estado de Minas Gerais e se destaca por ser um pólo regional, que exerce forte influência tanto nas cidades que estão ao seu entorno, como em outras regiões do estado e também sobre o centro-oeste brasileiro, principalmente sobre a região sul do estado de Goiás, estando estrategicamente localizada às margens da BR-050, principal ligação rodoviária entre Brasília e São Paulo (Figura 3).



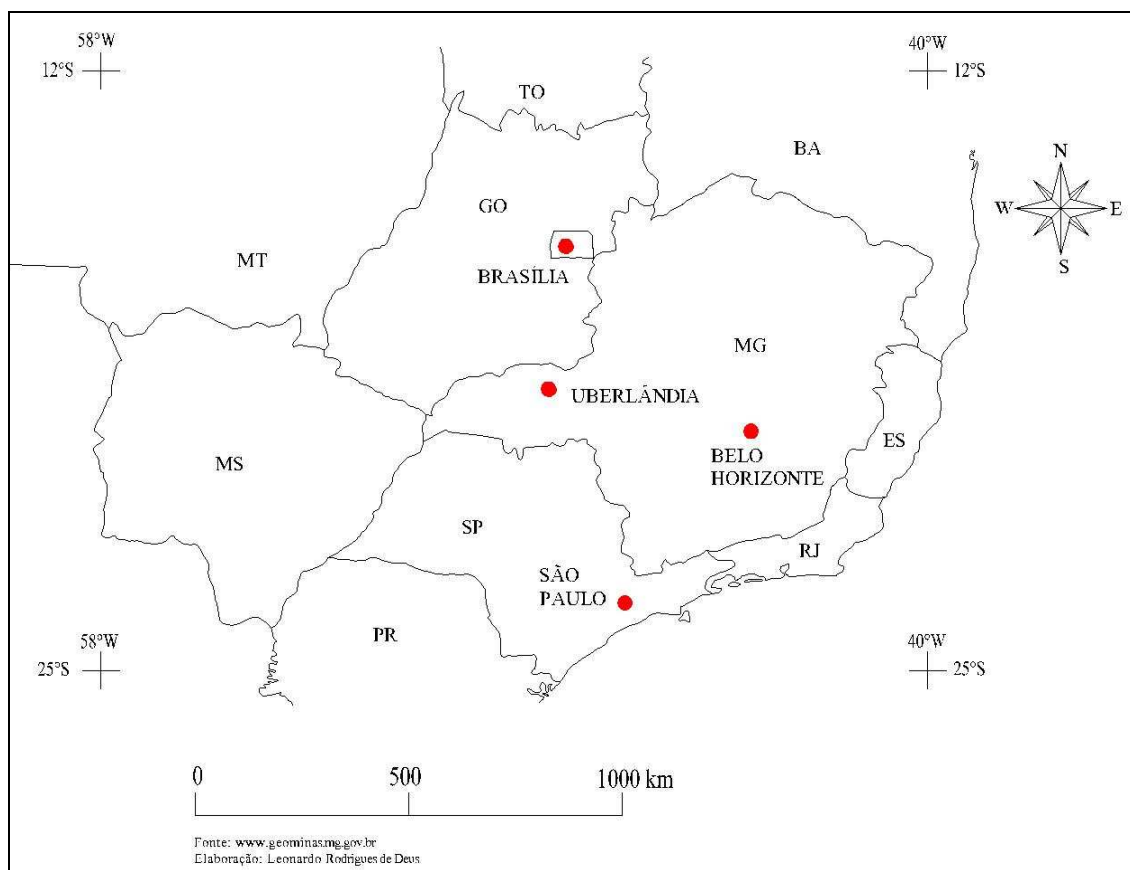


Figura 3: Localização geográfica de Uberlândia.

#### 4.2 Dados utilizados

Para a realização deste trabalho foram utilizados dois grupos de dados: imagens do sensor *IKONOS II* e dados *LiDAR*.

A imagem *IKONOS II* é uma imagem óptica, multiespectral, composta por 4 bandas onde cada uma corresponde a determinado comprimento de onda. A Banda 1 é a banda do azul, a Banda 2 é a banda do verde, a Banda 3 do vermelho, e a Banda 4 é a banda do infravermelho próximo.

Estas 4 bandas possuem resolução espacial de 1 metro, pois são resultados de uma fusão, realizada por Tomás (2010), entre as 4 bandas multiespectrais originais do sensor, que possuem 4 metros de resolução espacial, com a banda pancromática que possui 1 metro.

A Figura 4 mostra a imagem *IKONOS II* utilizada, em composição colorida 3R2G1B, conhecida como composição real, pois apresenta as cores reais dos objetos.



Figura 4: Imagem *IKONOS II* da área de estudo.

Quanto aos dados *LiDAR* (Light Detection And Ranging), Tomás (2010) coloca que este é um sistema topográfico que utiliza a porção infravermelha do espectro eletromagnético e se baseia na emissão e registro do retorno do sinal. O instrumento emite milhares de pulsos laser por segundo, de luz infravermelha, e mede as distâncias, a intensidade da energia refletida pelos objetos e os parâmetros de atitude do feixe (azimute e elevação), com o objetivo de determinar as elevações da superfície.

Foram utilizados então dois arquivos raster também gerados por Tomás (2010), criados a partir dos dados *LiDAR*. O primeiro corresponde ao Modelo Digital de Altura (MDA) da área de estudo, o qual é resultado da subtração entre o Modelo Digital de Superfície (MDS) e o Modelo Digital de Terreno (MDT), obtidos a partir da nuvem de pontos

*LiDAR*. O segundo arquivo raster é uma imagem Intensidade, que corresponde à intensidade da energia refletida pelos objetos atingidos pela luz infravermelha emitida.

Cada pixel do raster MDA apresenta um valor de altitude referente aos alvos que foram atingidos pelos pulsos do sensor *LiDAR*, e cada pixel da imagem Intensidade apresenta a intensidade de retorno da luz infravermelha para o sensor.

A Figura 5 mostra a imagem MDA utilizada no estudo. As cores mais claras na imagem indicam as áreas mais alta em relação ao solo.

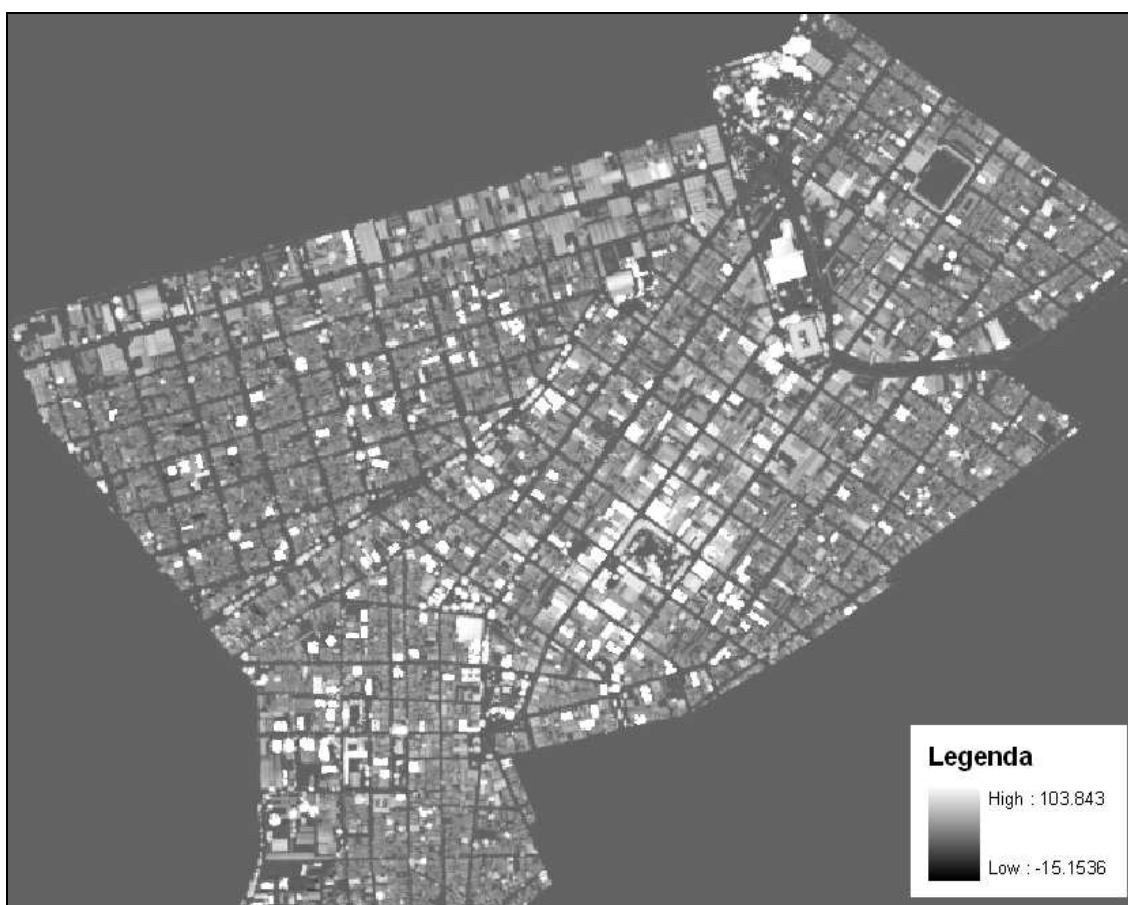


Figura 5: Imagem MDA da área de estudo.

Já a Figura 6 mostra a imagem Intensidade obtida para a área de estudo. Os tons mais claros na imagem indicam as áreas em a luz infravermelha teve índices mais altos de retorno.

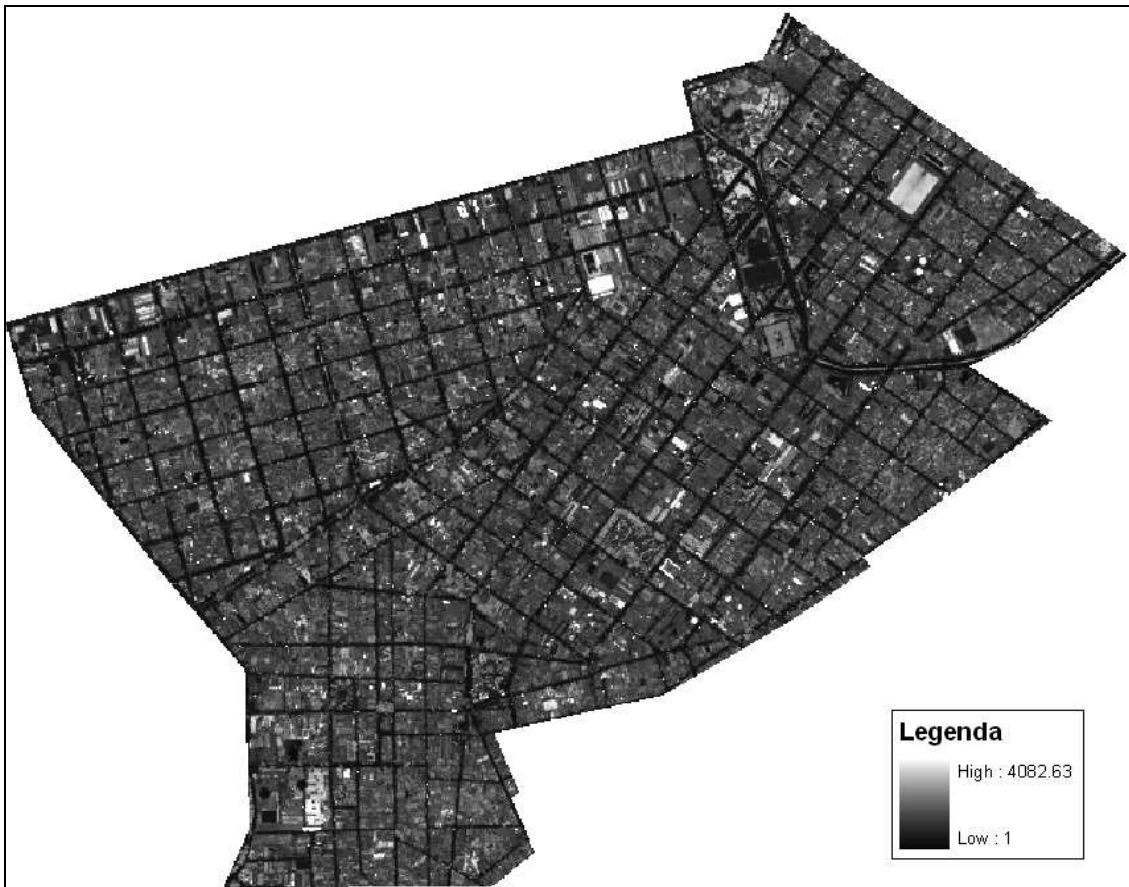


Figura 6: Imagem Intensidade da área de estudo.

O conjunto total de atributos disponíveis para este estudo foi composto então por 4 cenas do sensor *IKONOS II*, referentes às bandas multiespectrais do azul, verde, vermelho e infravermelho, e 2 imagens oriundas dos dados *LiDAR*, que são as imagens de intensidade e Modelo Digital de Altura.

Para avaliar a eficiência dos classificadores escolhidos e avaliar os benefícios destes atributos para prever a classificação da cobertura do solo urbano, os testes foram realizados com 3 grupos de atributos diferentes.

No primeiro grupo foram utilizadas apenas as 4 bandas do sensor *IKONOS II* (bandas azul, verde, vermelho e infravermelho) para predição da classificação da cobertura do solo urbano na área de estudo. O segundo grupo contou apenas com as informações de intensidade e de altitude fornecidas pelo sistema *LiDAR*. E o terceiro grupo de atributos foi composto por todos os dados do sensor *IKONOS II* e do sistema *LiDAR*.



### 4.3 Procedimentos Metodológicos

De posse das 6 imagens da área de estudo, 4 cenas *IKONOS II*, uma imagem MDA, e uma imagem Intensidade, o primeiro passo para se obter a classificação da cobertura do solo urbano foi a obtenção de amostras dos alvos a serem discriminados.

Devido ao grande número de alvos urbanos foram obtidas nesta etapa, amostras de 12 alvos diferentes: asfalto, edifício, piscina, solo exposto, sombra, telhado azul, telhado branco, telhado escuro, telhado marrom, telhado marrom escuro, vegetação e vegetação rasteira.

A Figura 7 (a;b) mostra uma das amostras selecionadas para um dos 12 alvos.

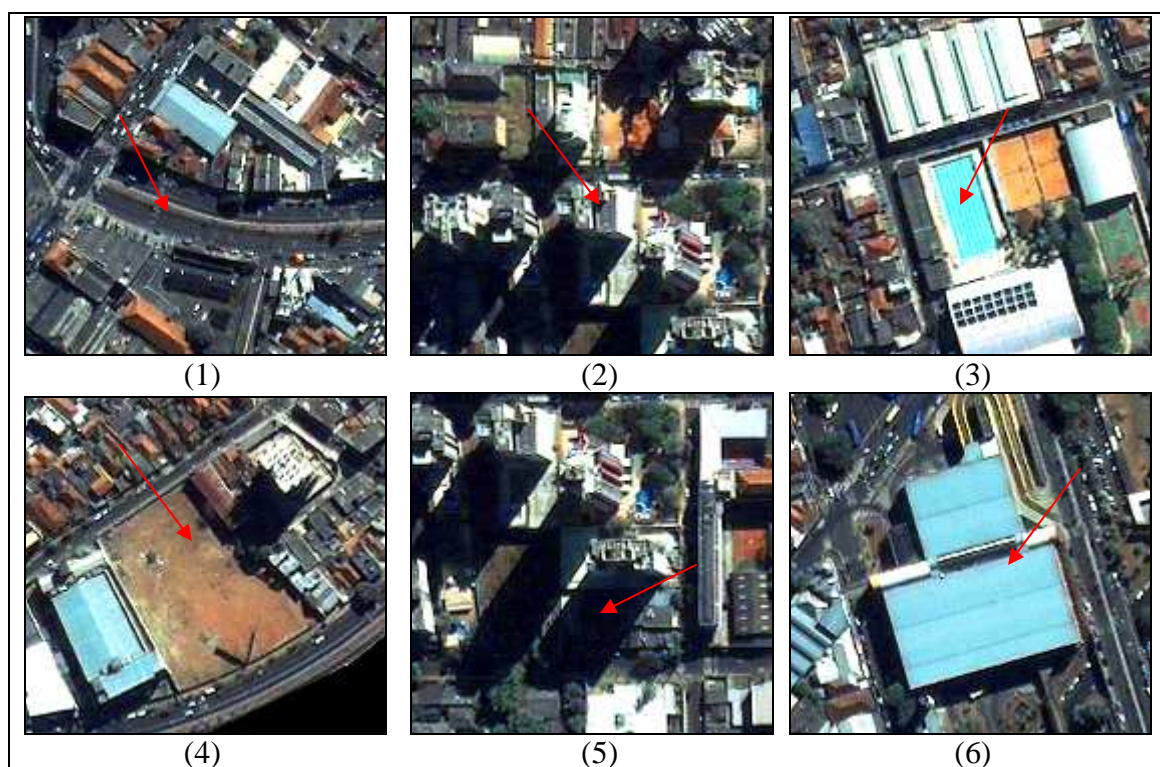


Figura 7a: 1-Asfalto; 2-Edifício; 3-Piscina; 4-Solo Exposto; 5-Sombra; 6-Telhado Azul.

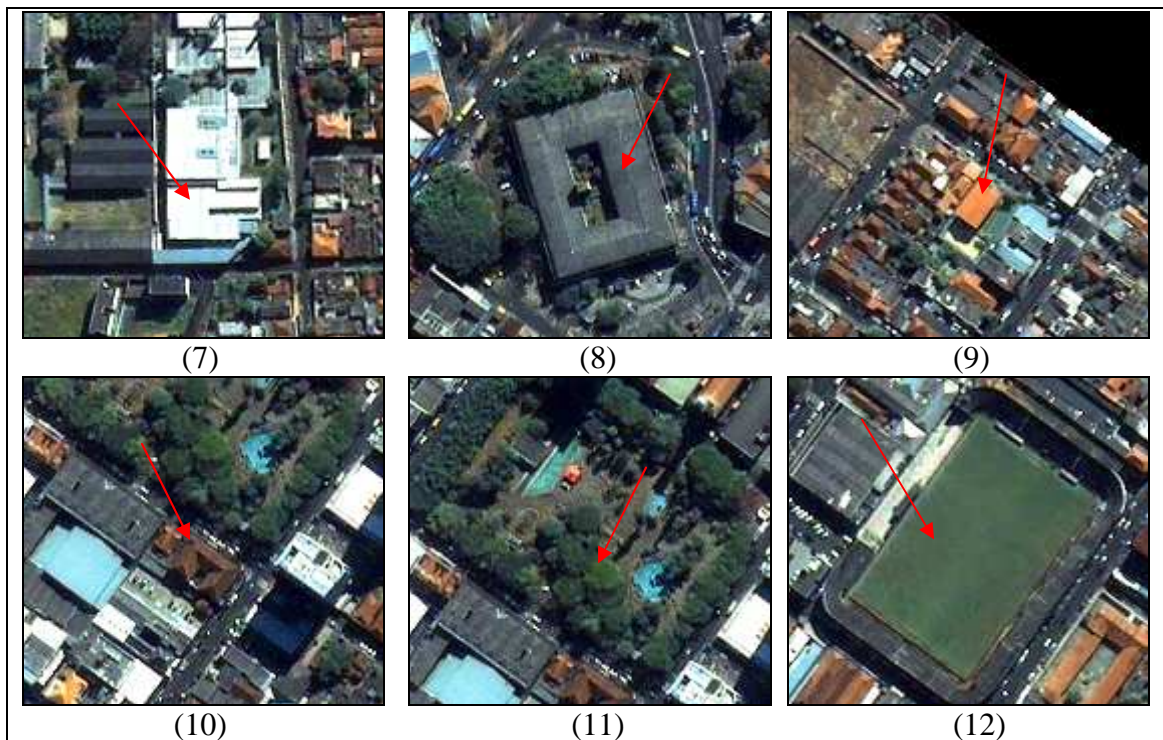


Figura 7b: 7-Telhado Branco; 8-Telhado Escuro; 9-Telhado Marrom; 10-Telhado Marrom Escuro; 11-Vegetação; 12-Vegetação rasteira.

Selecionadas as amostras dos alvos foi então preparado um arquivo tipo .ARFF com o conjunto de amostras, para que este fosse utilizado no processo de predição da classificação da cobertura do solo.

O arquivo ARFF é uma tabela de dados, onde cada pixel das amostras coletadas passa a ser uma instância da tabela e o valor associado a cada pixel, em cada uma das imagens, são os atributos da tabela, bem como a classe de uso do solo determinada para os alvos.

A Figura 8 mostra parte do arquivo ARFF, contendo as instâncias e os atributos utilizados. O arquivo completo contém 32.726 instâncias e 7 atributos, além dos valores de linha e colunas de cada pixel. Importante lembrar que todas as imagens utilizadas possuem exatamente os mesmos números de linhas e colunas, tendo o pixel o tamanho de 1 x 1 metro.

Relation: images-samples									
No.	linha Numeric	coluna Numeric	ikonos_b1 Numeric	ikonos_b2 Numeric	ikonos_b3 Numeric	ikonos_b4 Numeric	mda Numeric	intensidade Numeric	classe Nominal
29010	1421.0	1520.0	279.0	304.0	282.0	315.0	9.627136	37.362198	TELHADO_MARROM_ESCURO
24191	1030.0	735.0	280.0	240.0	187.0	148.0	53.067871	37.975357	EDIFICIO
30336	1858.0	732.0	280.0	185.0	113.0	93.0	0.712036	27.296883	SOMBRA
30530	1866.0	725.0	280.0	170.0	107.0	92.0	-0.027649	22.03653	SOMBRA
2466	301.0	1873.0	281.0	178.0	117.0	95.0	0.140503	15.737288	SOMBRA
26495	1296.0	1442.0	281.0	260.0	171.0	486.0	9.207214	66.859085	VEGETACAO
29825	1838.0	748.0	281.0	206.0	125.0	117.0	1.862915	39.061935	SOMBRA
30486	1864.0	723.0	281.0	174.0	109.0	105.0	-0.247437	21.90346	SOMBRA
2577	308.0	1858.0	282.0	185.0	126.0	98.0	4.477478	37.381714	SOMBRA
24918	1080.0	1747.0	282.0	171.0	106.0	83.0	10.928223	26.740057	SOMBRA
25140	1091.0	1739.0	282.0	180.0	120.0	85.0	5.507385	24.780926	SOMBRA
29152	1452.0	1357.0	282.0	262.0	206.0	185.0	59.037781	41.892498	EDIFICIO
29589	1753.0	699.0	282.0	243.0	164.0	140.0	73.322754	25.513657	EDIFICIO
29893	1841.0	738.0	282.0	153.0	94.0	87.0	-0.240906	58.136295	SOMBRA
31000	1916.0	869.0	282.0	280.0	297.0	341.0	8.661499	34.678867	TELHADO_MARROM_ESCURO
26673	1303.0	1569.0	283.0	281.0	167.0	411.0	12.030884	22.460241	VEGETACAO
29712	1833.0	738.0	283.0	198.0	112.0	84.0	-0.328064	43.80658	SOMBRA
30985	1915.0	890.0	283.0	289.0	253.0	336.0	9.12262	33.14716	TELHADO_MARROM_ESCURO
1490	131.0	1849.0	284.0	230.0	165.0	508.0	24.360168	30.791203	VEGETACAO
24674	1068.0	1756.0	284.0	175.0	103.0	88.0	3.756348	36.063671	SOMBRA
28225	1405.0	1516.0	284.0	284.0	297.0	290.0	11.228943	38.492081	TELHADO_MARROM_ESCURO
24667	1068.0	1749.0	285.0	173.0	117.0	80.0	4.81189	29.45829	SOMBRA
25148	1092.0	1730.0	285.0	185.0	113.0	80.0	3.671814	29.404619	SOMBRA
26266	1285.0	1454.0	285.0	301.0	172.0	530.0	3.383423	60.971657	VEGETACAO
29104	1450.0	1355.0	285.0	278.0	228.0	231.0	58.39679	44.98304	EDIFICIO
29972	1844.0	739.0	285.0	202.0	117.0	91.0	1.862366	34.517288	SOMBRA
3007	348.0	1789.0	286.0	242.0	152.0	306.0	16.998718	55.303787	VEGETACAO
3008	348.0	1790.0	286.0	242.0	151.0	347.0	17.32196	54.327282	VEGETACAO
24077	1027.0	728.0	286.0	251.0	161.0	122.0	46.329895	36.666904	EDIFICIO
24227	1031.0	734.0	286.0	229.0	178.0	143.0	54.197998	33.179497	EDIFICIO
26080	1271.0	1461.0	286.0	261.0	157.0	494.0	8.266907	50.625706	VEGETACAO
29848	1839.0	745.0	286.0	176.0	101.0	93.0	3.096436	69.720673	SOMBRA
30309	1857.0	732.0	286.0	196.0	118.0	94.0	3.489136	27.296883	SOMBRA
30365	1859.0	735.0	286.0	177.0	100.0	98.0	2.2146	28.173883	SOMBRA
2312	240.0	1981.0	287.0	280.0	294.0	314.0	6.178589	30.949707	TELHADO_MARROM_ESCURO
2535	305.0	1870.0	287.0	191.0	125.0	87.0	0.641907	62.919506	SOMBRA
24800	1074.0	1755.0	287.0	239.0	138.0	97.0	8.896667	30.063276	SOMBRA
25817	1257.0	722.0	287.0	278.0	276.0	334.0	3.651306	37.701973	TELHADO_MARROM_ESCURO
29715	1833.0	741.0	287.0	202.0	113.0	96.0	0.12207	37.003895	SOMBRA

Figura 8: Amostra do arquivo ARFF utilizado no software WEKA.

Para se trabalhar com o arquivo ARFF gerado e realizar a predição da classificação da cobertura do solo urbano, foi o utilizado o software WEKA (Waikato Environment for Knowledge Analysis), o qual é uma coleção de algoritmos de aprendizagem de máquina para realizar tarefas de mineração de dados, que contem várias ferramentas para pré-processamento dos dados e visualização, além de algoritmos para regressão, classificação, agrupamentos, mineração de regras de associação e seleção de atributos (HALL et al., 2009).

Foram então utilizados três classificadores diferentes que estão implementados no software WEKA: classificador IBK, classificador J48, e classificador MultiLayerPerceptron (MLP).

O classificador IBK utilizado no software WEKA é um algoritmo de aprendizagem baseado em instâncias (Instance-Based Learning – IBL), o que prevê uma classificação usando somente instâncias específicas e não mantém um grupo de abstrações derivadas das instâncias (AHA et al., 1991).

O algoritmo de aprendizagem baseado em instâncias é derivado do classificador de vizinho mais próximo, sendo bastante similar a algoritmos de vizinhos mais próximos editados, os quais também armazenam e utilizam apenas instâncias selecionadas para gerar previsões de classificações (AHA et al., 1991).

De acordo com a Aha et al. (1991), diferentemente dos algoritmos de vizinhos mais próximos editados, que não são incrementais e tem como objetivo principal manter uma consistência perfeita com os dados iniciais de treinamento, os algoritmos de aprendizagem baseados em instâncias são incrementais e seus objetivos incluem a maximização da acurácia da classificação sobre as instâncias apresentadas na seqüência do processo.

O parâmetro K existente no algoritmo determina o número de vizinhos mais próximos a serem usados para predição de uma classificação. O classificador IBK é um aprendiz baseado em instâncias simples que usa a classe do vizinho mais próximo das instâncias de treinamento para classificar as instâncias de testes.

A classificação por K vizinhos mais próximo (KNN – K nearest neighbor), encontra um grupo de K objetos em um grupo de treinamento que estão próximos a um objeto de teste, e baseia a atribuição de um rótulo sobre a predominância de uma classe particular nesta vizinhança (STEINBACH e TAN, 2009).

O algoritmo KNN funciona da seguinte forma: dado um grupo de treinamento D e um objeto de teste z, o qual é um vetor de valores de atributos e tem uma classe



desconhecida, o algoritmo computa a distância (ou similaridade) entre  $z$  e todos os objetos de treinamento para determinar sua lista de vizinhos mais próximos e em seguida atribuir uma classe para  $z$ , tendo em conta a classe que é maioria entre os objetos da vizinhança (STEINBACH e TAN, 2009).

Quando o valor de  $K$  é muito pequeno, o resultado pode ser sensível a pontos ruidosos existentes no conjunto de dados, e se o valor de  $K$  for muito grande, então a vizinhança pode incluir muitos pontos de outras classes (STEINBACH e TAN, 2009), podendo levar a predição de uma classe que não é a real para determinada instância.

Como classificador IBK utiliza um número de  $K$  vizinhos mais próximos para predizer a classificação dos dados analisados, para este estudo, o classificador foi testado utilizando a variável  $K$  no valor de 1, 3, 5, 8, 11 e 15 vizinhos mais próximos, para os três grupos de atributos considerados no trabalho.

Já o classificador J48 é baseado no algoritmo C4.5 e gera um classificador na forma de uma árvore de decisão, com uma estrutura onde uma folha pode indicar uma classe, ou um nó de decisão, que especifica algum teste a ser realizado sobre um único valor de atributo, com um ramo ou sub-árvore para cada saída possível do teste (QUINLAN, 1993).

C4.5 é um conjunto de algoritmos para problemas de classificação em aprendizagem de máquina e mineração de dados, sendo orientado para aprendizagem supervisionada. Dado um grupo de dados com valores de atributos, onde as instâncias são descritas por coleções de atributos e pertencem a uma classe, de um grupo de classes mutuamente exclusivas, o C4.5 aprende um mapeamento dos valores dos atributos para as classes que podem ser aplicados para classificar novas instâncias (RAMAKRISHNAN, 2009).

Quando se trabalha com o algoritmo C4.5, a entrada consiste em uma coleção de casos de treinamento, cada um tendo uma tupla de valores para um grupo fixo de atributos e um atributo de classe (KOHAVI e QUINLAN, 1999).

Uma árvore de decisão é uma série de questões sistematicamente organizadas, de modo que cada questão consulta um atributo e ramificações baseadas no valor do atributo (RAMAKRISHNAN, 2009).

A respeito dos sistemas de árvores de decisão, RAMAKRISHNAN (2009) explica que todos os métodos de indução de árvores começam com um nó raiz que representa o todo, determinado grupo de dados e os dados recursivamente divididos em subgrupos menores por meio de testes para um determinado atributo em cada nó. As sub-árvores denotam as divisões do grupo de dados original que satisfazem valores testados de um atributo específico. Este processo continua normalmente até um subgrupo ser “puro”, isto é, todas as instâncias em um subgrupo pertencerem à mesma classe, momento este em que o crescimento da árvore é encerrado.

O classificador MultilayerPerceptron (MLP), por sua vez, é um tipo de rede neural artificial que utiliza um sistema de aprendizado supervisionado por correção de erros (backpropagation) para classificar as instâncias.

As redes neurais artificiais são um paradigma de programação que procuram simular a microestrutura de um cérebro, e são bastante utilizados em problemas de inteligência artificial, desde simples tarefas de reconhecimento de padrões até manipulação simbólica avançada (NORIEGA, 2005).

As redes neurais podem ser consideradas como processadores de sinais que estão arquitetados para “passar valores”. Cada neurônio pode receber um número de entradas a partir da camada de entrada ou outras camadas, e então produz uma única saída que pode ser uma saída para outros neurônios ou uma saída para uma rede global (SHIH, 2010).

De acordo com Shih (2010), MultilayerPerceptron são redes de alimentação com uma ou mais camadas, chamadas camadas escondidas, de células neurais, chamadas neurônios escondidos, entre as camadas de entrada e saída. O processo de aprendizagem das MLPs são conduzidas com o algoritmo de aprendizagem de propagação e retorno do erro.

O algoritmo de propagação e retorno do erro (backpropagation) funciona da seguinte forma: primeiro é apresentado o padrão de entrada, e então a resposta de uma unidade é propagada como entrada para as unidades na camada seguinte, até a camada de saída, onde é obtida a resposta da rede e o erro é calculado; em seguida é realizado o processo de volta, desde a camada de saída até a camada de entrada, sendo feitas alterações nos pesos calculados.

Karpagavalli et al. (2009) expõe que MLPs são ferramentas valiosas em problemas quando se tem pouco ou nenhum conhecimento sobre a forma das relações entre os vetores de entrada e suas saídas correspondentes.

Para avaliar a eficiência do classificador MLP, o mesmo foi executado utilizando 1, 5, 10, 25 e 50 camadas escondidas, para cada um dos três grupos de atributos utilizados.

## 5. Resultados e Discussões

### 5.1 Resultados Classificador IBK

A Tabela 1 sintetiza os resultados da precisão com que os dados amostrais foram classificados, utilizando o classificador IBK com os diferentes números de vizinhos mais próximos.

Tabela 1: Resultados Classificador IBK

Valor K	Total Instâncias Classificadas Corretamente			Precisão Instâncias Classificadas Corretamente (%)		
	Primeiro Grupo de Atributos	Segundo Grupo de Atributos	Terceiro Grupo de Atributos	Primeiro Grupo de Atributos	Segundo Grupo de Atributos	Terceiro Grupo de Atributos
1	32726	32726	32726	100	100	100
3	31765	28360	32623	97,0635	86,6589	99,6853
5	31624	27553	32578	96,6326	84,1930	99,5478
8	31447	26960	32518	96,0918	82,3810	99,3644
11	31345	26718	32484	95,7801	81,6415	99,2605
15	31268	26498	32453	95,5448	80,9693	99,1658

Analisando a Tabela 1, observa-se que, para os três grupos de atributos trabalhados, a aplicação do classificador IBK utilizando K igual a 1, foi a configuração que apresentou os melhores resultados. Nesta situação, a precisão das instâncias classificadas corretamente foi de 100% para os três grupos de atributos.

À medida que o valor de K foi aumentando, a Tabela 1 mostra que a precisão das instâncias classificadas corretamente foi diminuindo, e quanto menos atributos foram utilizados na predição da classificação, maior foi a perda de precisão. Desta forma, o grupo de atributos que utilizou apenas os dados *LiDAR*, teve a maior perda de precisão com o aumento de K, caindo de 100% com K igual a 1, para 80,9% com K igual a 15, enquanto que o grupo que utilizou todos os atributos (*IKONOS* e *LiDAR*), teve a menor queda de precisão ao predizer a classificação das amostras utilizadas, não chegando a perda total de 1%.

Observa-se ainda que quando se utiliza apenas as 4 bandas do sensor *IKONOS II* (primeiro grupo de dados) para prever a classificação das amostras, a precisão das instâncias classificadas corretamente cai menos do que quando se utiliza apenas os dados do sistema *LiDAR* (segundo grupo de dados), ficando sempre acima de 95% de precisão. Isto mostra que as bandas do sensor *IKONOS II* têm maior capacidade para distinguir entre os alvos que compõem a cobertura do solo, indicando que os dados *LiDAR* devem ser utilizados como uma fonte de dados complementar ao uso de imagens ópticas para classificar a cobertura do solo urbano.

Os resultados observados com base na Tabela 1 indicam que o conjunto de amostras apresenta uma grande mistura de informações, visto que o classificador com o valor de  $K$  igual a 1 apresentou os melhores resultados, e independente de quanto se aumente este valor de  $K$ , os resultados são sempre piores, indicando que os dados amostrais não se caracterizam pela presença de outliers, que são pontos ruidosos no conjunto de dados.

Analisando as matrizes de confusão obtidas pela aplicação do classificador IBK, utilizando o terceiro grupo de atributos, o qual apresentou os melhores resultados, é possível verificar quais são as classes que apresentam maior grau de confusão entre as instâncias classificadas.

Utilizando apenas um vizinho mais próximo ( $K = 1$ ), não houve nenhuma confusão entre as classes, todas tiveram 100% de instâncias classificadas corretamente. Mas à medida que se aumenta o número de vizinhos mais próximos no processo, algumas confusões podem ser observadas entre determinados alvos.

Para facilitar a visualização dos resultados, as matrizes de confusão resultantes foram trabalhadas para gerar matrizes de confusão que ilustrem os resultados para cada um dos alvos analisados, com os diferentes valores de vizinhos mais próximos utilizados. Para leitura das tabelas a seguir deve-se considerar:

- a = asfalto
- b = piscina
- c = solo exposto
- d = sombra

- e = edifício
- f = telhado azul
- g = telhado branco
- h = telhado escuro
- i = telhado marrom
- j = telhado marrom escuro
- k = vegetação
- l = vegetação rasteira

As classes piscina, sombra e vegetação rasteira tiveram 100% de acerto independente do número de vizinhos mais próximos empregado, mostrando que os atributos utilizados na predição da classificação são excelentes para distinguir estes 3 tipos de alvos.

Outra classe que também foi muito bem discriminada foi a classe Telhado Azul, que só não teve 100% de acerto para todos os valores de vizinhos mais próximos, pois para K igual a 3, uma instância foi classificada como Telhado Branco, conforme pode ser visto na Tabela 2. Isto pode ter acontecido pelo fato de alguma das amostras da classe Telhado Azul não ter sido obtida com o devido cuidado, ou então por contar a presença de um outlier.

Tabela 2: Matriz de confusão da classe Telhado Azul – Classificador IBK

Valor de K	a	b	c	d	e	f	g	h	i	j	k	l	<-- Classificado como
1	0	0	0	0	0	5638	0	0	0	0	0	0	<b>f = telhado azul</b>
3	0	0	0	0	0	5637	1	0	0	0	0	0	
5	0	0	0	0	0	5638	0	0	0	0	0	0	
8	0	0	0	0	0	5638	0	0	0	0	0	0	
11	0	0	0	0	0	5638	0	0	0	0	0	0	
15	0	0	0	0	0	5638	0	0	0	0	0	0	

As demais classes apresentaram maiores confusões entre as instâncias classificadas, mas todas com valores bem baixos em relação ao número total de instâncias de cada classe.

Para a classe asfalto, a Tabela 3 mostra que à medida que o valor de K foi aumentado, aumentou-se o número de instâncias desta classe classificadas como telhado escuro. Este resultado é aceitável, pois estas duas classes possuem uma resposta espectral semelhante para as bandas do sensor *IKONOS II*. Contudo, para este trabalho, não era

esperado que isto acontecesse, visto que o Modelo Digital de Altitude (MDA) gerado a partir do sistema *LiDAR* foi utilizado no processo de classificação, e esperava-se que dois alvos que estão em alturas diferentes fossem bem discriminados.

Tabela 3: Matriz de confusão da classe Asfalto – Classificador IBK

Valor de K	a	b	c	d	e	f	g	h	i	j	k	l	<-- Classificado como
1	291	0	0	0	0	0	0	0	0	0	0	0	<b>a = asfalto</b>
3	284	0	0	0	0	0	0	7	0	0	0	0	
5	277	0	0	0	0	0	0	14	0	0	0	0	
8	276	0	0	0	0	0	0	15	0	0	0	0	
11	271	0	0	0	0	0	0	20	0	0	0	0	
15	265	0	0	0	0	0	0	26	0	0	0	0	

A classe Solo Exposto, que possui um universo de 2.566 amostras, teve sempre menos de 10 delas classificadas como Telhado Marrom Escuro, como mostra a Tabela 4. Este resultado também é aceitável devido ao comportamento espectral dos alvos serem similares, mas que também não deveria ocorrer devido ao uso do MDA. Uma única instância da classe solo exposto foi classificada como Telhado Marrom, quando o classificar utilizou K igual a 3.

Tabela 4: Matriz de confusão da classe Solo Exposto – Classificador IBK

Valor de K	a	b	c	d	e	f	g	h	i	j	k	l	<-- Classificado como
1	0	0	2566	0	0	0	0	0	0	0	0	0	<b>c = solo exposto</b>
3	0	0	2561	0	0	0	0	0	1	4	0	0	
5	0	0	2560	0	0	0	0	0	0	6	0	0	
8	0	0	2559	0	0	0	0	0	0	7	0	0	
11	0	0	2558	0	0	0	0	0	0	8	0	0	
15	0	0	2557	0	0	0	0	0	0	9	0	0	

A classe Edifício foi uma das classes que apresentou maior nível de confusão, conforme se observa na Tabela 5. Para esta classe, à medida que um maior número de vizinhos mais próximos foi utilizado, aumentou-se o número de instâncias classificadas principalmente como Telhado Branco, e em segundo lugar instâncias classificadas como Telhado Escuro. Para o classificador com valor de K igual ou maior que 5, sempre houve uma instância classificada como sombra, e para o classificador com K igual ou

maior que 8, sempre houve também uma instância classificada como vegetação. Além disso, houve ainda 4 instâncias classificadas como Asfalto, uma para K igual a 3, duas para K igual a 5, e uma para K igual a 8.

Tabela 5: Matriz de confusão da classe Edifício – Classificador IBK

Valor de K	a	b	c	d	e	f	g	h	i	j	k	l	<-- Classificado como
1	0	0	0	0	1308	0	0	0	0	0	0	0	<b>e = edifício</b>
3	1	0	0	0	1271	0	24	12	0	0	0	0	
5	2	0	0	1	1258	0	32	15	0	0	0	0	
8	1	0	0	1	1249	0	35	21	0	0	1	0	
11	0	0	0	1	1245	0	39	22	0	0	1	0	
15	0	0	0	1	1239	0	44	23	0	0	1	0	

Para a classe Telhado Branco, o classificador fez confusão principalmente com as classes Telhado Azul, Telhado Escuro e Telhado Marrom Escuro, com destaque para a classe Telhado Escuro, como mostra a Tabela 6. Esta confusão é bastante curiosa, uma vez que a resposta espectral de um telhado branco é bem diferente da resposta de um telhado escuro.

Tabela 6: Matriz de confusão da classe Telhado Branco – Classificador IBK

Valor de K	a	b	c	d	e	f	g	h	i	j	k	l	<-- Classificado como
1	0	0	0	0	0	0	4301	0	0	0	0	0	<b>g = telhado branco</b>
3	0	0	1	0	3	4	4287	4	0	2	0	0	
5	0	0	0	0	2	5	4287	5	0	2	0	0	
8	0	0	1	0	2	7	4277	11	0	3	0	0	
11	0	0	0	0	0	9	4265	19	2	6	0	0	
15	0	0	0	1	0	8	4263	22	2	5	0	0	

A Tabela 6 mostra ainda que para os valores de K igual a 11 e K igual a 15, duas instâncias foram classificadas como Telhado Marrom. Houve também 7 instâncias desta classe classificadas como Edifício, três para K igual a 3, duas para K igual a 5 e duas para K igual a 8, além de 2 instâncias classificadas como solo exposto, 1 para K igual a 3 e 1 para K igual a 8.



A classe Telhado Escuro apresentou as maiores confusões com a classe Telhado Branco, conforme já havia sido observado anteriormente, na análise inversa entre estes dois alvos. Outras confusões para a classe Telhado Escuro ocorreram, conforme se observa na Tabela 7, com instâncias classificadas como Telhado Marrom Escuro, Asfalto, Sombra e Edifício, que são classes que possuem algumas semelhanças nas respostas espectrais destes alvos.

Tabela 7: Matriz de confusão da classe Telhado Escuro – Classificador IBK

Valor de K	a	b	c	d	e	f	g	h	i	j	k	l	<-- Classificado como
1	0	0	0	0	0	0	0	4883	0	0	0	0	h = telhado escuro
3	3	0	0	0	0	0	8	4871	0	1	0	0	
5	2	0	0	0	0	0	18	4861	0	2	0	0	
8	2	0	0	3	0	0	27	4847	0	4	0	0	
11	2	0	0	6	0	0	27	4844	0	4	0	0	
15	2	0	0	9	1	0	29	4839	0	3	0	0	

A classe Telhado Marrom, por sua vez, foi confundida com apenas duas classes com o aumento no número de vizinhos mais próximos no classificador, como mostra a Tabela 8, que foi a classe Solo Exposto, em maior grau, confusão esta bastante recorrente em estudos urbanos, e em menor grau com a classe Telhado Marrom Escuro.

Tabela 8: Matriz de confusão da classe Telhado Marrom – Classificador IBK

Valor de K	a	b	c	d	e	f	g	h	i	j	k	l	<-- Classificado como
1	0	0	0	0	0	0	0	0	882	0	0	0	i = telhado marrom
3	0	0	6	0	0	0	0	0	873	3	0	0	
5	0	0	6	0	0	0	0	0	872	4	0	0	
8	0	0	13	0	0	0	0	0	864	5	0	0	
11	0	0	15	0	0	0	0	0	861	6	0	0	
15	0	0	16	0	0	0	0	0	860	6	0	0	

Já a classe Telhado Marrom Escuro apresentou níveis de confusão praticamente equivalentes para com as classes Solo Exposto e Telhado Escuro, como mostra a Tabela 9. Houve ainda para esta classe, uma instância classificada como Telhado Branco para K igual a 5 e cinco instâncias classificadas como Telhado Marrom, sendo duas para K igual a 8, duas para K igual a 11 e 1 para K igual a 15.

Tabela 9: Matriz de confusão da classe Telhado Marrom Escuro – Classificador IBK

Valor de K	a	b	c	d	e	f	g	h	i	j	k	l	<-- Classificado como
1	0	0	0	0	0	0	0	0	0	1624	0	0	<b>j = telhado marrom escuro</b>
3	0	0	8	0	0	0	0	7	0	1609	0	0	
5	0	0	9	0	0	0	1	15	0	1599	0	0	
8	0	0	16	0	0	0	0	18	2	1588	0	0	
11	0	0	19	0	0	0	0	20	2	1583	0	0	
15	0	0	24	0	0	0	0	21	1	1578	0	0	

Por fim, a classe Vegetação teve com o aumento do número de vizinhos mais próximos, algumas instâncias classificadas principalmente como vegetação rasteira, conforme se observa na Tabela 10. Esta situação novamente não era algo esperado que ocorresse, pois o atributo MDA deveria separar bem estas duas classes pela diferença de altura. Outras instâncias foram classificadas como sombra e para K igual a 3, 5 e 8 uma instância foi classificada como Telhado Escuro.

Tabela 10: Matriz de confusão da classe Vegetação – Classificador IBK

Valor de K	a	b	c	d	e	f	g	h	i	j	k	l	<-- Classificado como
1	0	0	0	0	0	0	0	0	0	0	3000	0	<b>k = vegetacao</b>
3	0	0	0	0	0	0	0	1	0	0	2997	2	
5	0	0	0	1	0	0	0	1	0	0	2993	5	
8	0	0	0	4	0	0	0	1	0	0	2987	8	
11	0	0	0	4	0	0	0	0	0	0	2986	10	
15	0	0	0	7	0	0	0	0	0	0	2981	12	

Os resultados obtidos com a utilização do classificador IBK mostraram que a utilização de imagem óptica em conjunto com dados *LiDAR* favorece para que se obtenha uma classificação da cobertura do solo urbano com alto grau de precisão, visto que a classificação das amostras de testes teve resultados acima de 99% de precisão.

## 5.2 Resultados Classificador J48

O classificador J48 foi aplicado também para os 3 grupos de atributos definidos para o estudo.

A Tabela 11 apresenta os resultados gerais obtidos com o empregado do classificador J48, o qual funciona como um classificador do tipo árvore de decisão.

Tabela 11: Resultados Classificador J48

Atributos	Total Instâncias Classificadas Corretamente	Precisão Instâncias Classificadas Corretamente (%)	Número de Folhas	Tamanho da Árvore
Primeiro Grupo de Atributos	31739	96,984	453	905
Segundo Grupo de Atributos	29090	88,8896	1573	3145
Terceiro Grupo de Atributos	32616	99,6639	230	459

Os resultados sintetizados na Tabela 11 indicam que o uso das imagens do sensor *IKONOS II* juntamente com os dados *LiDAR* são a melhor opção, entre as testadas, para classificar a cobertura do solo urbano. Para o terceiro grupo de atributos, que engloba as imagens *IKONOS II* e os dados *LiDAR*, a precisão de instâncias classificadas corretamente foi da ordem de 99,6%, melhor resultado obtido, e para este grupo de atributos a árvore de decisão gerada foi a menor, indicando ser esta uma árvore de mais fácil compreensão.

A Tabela 11 mostra também que utilização de apenas o conjunto de dados *LiDAR*, com a aplicação deste classificador, apresenta muitas dificuldades para classificar a cobertura do solo urbano. Os resultados para este segundo grupo de atributos apresentam uma precisão de instâncias classificadas corretamente abaixo de 89% e a árvore de decisão gerada é extremamente grande, sendo 7 vezes maior do que a árvore gerada quando se utiliza o terceiro grupo de atributos. Quando se utiliza apenas o conjunto de dados *IKONOS II*, os resultados mostram que a precisão de instâncias classificadas corretamente é alto, chegando praticamente a 97%, mas a árvore de decisão gerada é

ainda 2 vezes maior que a árvore obtida quando se aplica dados *IKONOS II* e dados *LiDAR* concomitantemente.

Ao analisar as matrizes de confusão obtidas pela aplicação do classificador J48 para os diferentes grupos de atributos, observa-se que quanto menos atributos são utilizados para prever a classificação, maior é o grau de confusão entre as classes analisadas.

O Quadro 1 mostra a matriz de confusão gerada pelo classificador J48 quando se utiliza apenas os dados *LiDAR* (segundo grupo de atributos) para predição da classificação da cobertura do solo urbano.

Quadro 1: Matriz de confusão – Segundo grupo de atributos – Classificador J48

a	b	c	d	e	f	g	h	i	j	k	l	<-- classified as
200	33	8	4	0	6	39	1	0	0	0	0	a = ASFALTO
30	493	60	5	0	0	40	0	4	0	1	19	b = PISCINA
30	31	2338	50	0	13	16	12	9	7	4	56	c = SOLO_EXPOSTO
8	29	95	1304	0	40	53	157	16	86	18	24	d = SOMBRA
0	0	5	6	1234	0	14	8	0	6	35	0	e = EDIFICIO
0	3	15	11	0	5549	2	47	0	0	6	5	f = TELHADO_AZUL
30	3	8	33	1	14	3823	165	86	80	56	2	g = TELHADO_BRANCO
1	3	38	91	1	61	158	4362	27	77	64	0	h = TELHADO_ESCURO
0	0	13	12	0	0	282	34	502	6	14	19	i = TELHADO_MARROM
0	0	13	79	1	1	107	188	3	1223	9	0	j = TELHADO_MARROM_ESCURO
2	5	22	43	0	12	168	299	36	26	2377	10	k = VEGETACAO
0	2	62	2	0	0	0	0	0	0	0	5685	l = VEGETACAO_RASTEIRA

Analisando o Quadro 1 observa-se que ocorre um grande número de confusões entre todas as classes analisadas, sem ser possível determinar um padrão entre as classes que estão sendo confundidas. O que se pode notar é que alvos de diferentes alturas, como Asfalto e Telhado Branco, que não deveriam apresentar confusões entre si no processo de classificação devido à utilização do MDA, tiveram algumas instâncias confundidas entre estes dois alvos. Outro fato interessante que pode ser observado, é a questão do atributo intensidade não ser adequado para distinção de alvos que tem comportamento espectral diferentes, como é o caso do alvo Piscina, que teve instâncias classificadas como Asfalto, Solo Exposto e Telhado Branco, ou ainda a classe Sombra, que teve instâncias classificadas como todas as outras classes, exceto como Edifício.

A matriz de confusão gerada a partir do uso dos atributos do sensor *IKONOS II* (primeiro grupo de atributos), Quadro 2, apresenta um nível de confusão bem menor entre as classes do que o que foi observado anteriormente.

Quadro 2: Matriz de confusão – Primeiro grupo de atributos – Classificador J48

a	b	c	d	e	f	g	h	i	j	k	l	<-- classified as
258	0	0	0	5	0	1	26	1	0	0	0	a = ASFALTO
0	652	0	0	0	0	0	0	0	0	0	0	b = PISCINA
0	0	2471	0	1	0	2	1	9	82	0	0	c = SOLO_EXPOSTO
0	0	0	1827	2	0	0	0	0	1	0	0	d = SOMBRA
10	0	8	1	759	10	166	352	0	2	0	0	e = EDIFICIO
0	0	0	0	1	5636	1	0	0	0	0	0	f = TELHADO_AZUL
1	0	3	0	62	7	4223	5	0	0	0	0	g = TELHADO_BRANCO
36	0	1	0	65	0	14	4763	0	3	1	0	h = TELHADO_ESCURO
0	0	13	0	0	0	1	2	856	10	0	0	i = TELHADO_MARROM
0	0	39	0	1	0	0	13	8	1563	0	0	j = TELHADO_MARROM_ESCURO
0	0	0	1	0	0	0	0	0	0	2985	14	k = VEGETACAO
0	0	0	0	0	0	0	0	0	0	5	5746	l = VEGETACAO_RASTEIRA

O Quadro 2 mostra que as confusões observadas entre as classes quando se utiliza as bandas multiespectrais do sensor *IKONOS II*, estão todas relacionadas às respostas espectrais semelhantes entre alguns alvos, como Solo Exposto e Telhado Marrom Escuro, ou Asfalto e Telhado Escuro por exemplo. Isto indica que as bandas do sensor *IKONOS II* sozinhas, tem maior capacidade para distinguir entre os alvos analisados do que quando se usa somente os dados do sistema *LiDAR*.

Quando se utiliza o terceiro grupo de atributos (dados *IKONOS* e *LiDAR*), a matriz de confusão mostra que o nível de confusões entre as classes é ainda menor (Quadro 3).

Quadro 3: Matriz de confusão – Terceiro grupo de atributos – Classificador J48

a	b	c	d	e	f	g	h	i	j	k	l	<-- classified as
285	0	0	0	0	0	0	6	0	0	0	0	a = ASFALTO
0	652	0	0	0	0	0	0	0	0	0	0	b = PISCINA
0	0	2558	0	0	0	0	0	2	6	0	0	c = SOLO_EXPOSTO
0	0	0	1829	0	0	0	0	0	0	1	0	d = SOMBRA
0	0	0	0	1282	0	14	11	1	0	0	0	e = EDIFICIO
0	0	0	0	0	5637	1	0	0	0	0	0	f = TELHADO_AZUL
0	0	0	0	2	3	4290	3	0	3	0	0	g = TELHADO_BRANCO
0	0	3	0	1	4	8	4864	1	1	1	0	h = TELHADO_ESCURO
0	0	4	0	0	0	1	0	875	2	0	0	i = TELHADO_MARROM
0	0	13	0	0	0	0	10	4	1597	0	0	j = TELHADO_MARROM_ESCURO
2	0	0	0	1	0	0	0	0	1	2996	0	k = VEGETACAO
0	0	0	0	0	0	0	0	0	0	0	5751	l = VEGETACAO_RASTEIRA

Com base no Quadro 3, observa-se que o número de instâncias classificadas incorretamente é praticamente irrelevante frente ao número total de instâncias pertencentes a cada classe. Dentre as confusões que mais se destacam nesta situação, tem-se a classe Edifício que foi classificada como Telhado Branco e Telhado Escuro, o que provavelmente se deve ao fato de as amostras da classe Edifício conterem prédios que tem telhados compostos por materiais diferentes, e a classe Telhado Marrom Escuro que teve algumas instâncias classificadas como Telhado Escuro e como Solo Exposto.



Além disso, como era esperado, o atributo MDA aparece no topo da árvore distinguindo algumas instâncias da classe Edifício dos demais alvos, uma vez que os maiores valores para o atributo MDA são referentes a instâncias da classe Edifícios.

Comparando os resultados obtidos pelo classificador J48 frente aos resultados obtidos pelo classificador IBK, observa-se que a diferença de precisão entre as instâncias classificadas corretamente considerando apenas o uso do terceiro grupo de atributos pode ser considerada irrelevante, uma vez todos os resultados estão acima de 99%. Contudo, o classificador J48 leva uma grande vantagem na questão do tempo necessário para processamento. O melhor resultado obtido com o classificar J48 levou um tempo de processamento de apenas 5,23 segundo, enquanto o melhor resultado obtido com o classificador IBK levou um tempo aproximado de 5 minutos. Isto mostra que o custo computacional do classificador J48 é bem menor, sem perca de precisão na predição da classificação dos dados utilizados.

### 5.3 Resultados Classificador MultiLayerPerceptron (MLP)

O último classificador implementado no software Weka testado neste trabalho foi o classificador MultilayerPerceptron (MLP), o qual como foi dito, é um classificador do tipo rede neural que usa um sistema de aprendizado supervisionado por correção de erros (backpropagation) para prever a classificação das instâncias analisadas.

A Tabela 12 mostra os resultados das instâncias classificadas corretamente obtidos por este classificador.

Tabela 12: Resultados Classificador MLP

Número de Camadas Escondidas	Total Instâncias Classificadas Corretamente			Precisão Instâncias Classificadas Corretamente (%)		
	Primeiro Grupo de Atributos	Segundo Grupo de Atributos	Terceiro Grupo de Atributos	Primeiro Grupo de Atributos	Segundo Grupo de Atributos	Terceiro Grupo de Atributos
<b>1</b>	18509	17011	19383	56,5576	51,9801	59,2281
<b>5</b>	29939	23723	31945	91,4838	72,4898	97,6135
<b>10</b>	30370	24332	32399	92,8008	74,3507	99,0008
<b>25</b>	31368	24673	32527	95,8504	75,3927	99,3919
<b>50</b>	31341	24660	32551	95,7679	75,3529	99,4653

Analisando a Tabela 12, observa-se que à medida que o número de camadas escondidas aumenta, o número de instâncias classificadas corretamente também aumenta, e assim como já havia sido observado para os outros dois classificadores avaliados, nota-se que quando se utiliza o terceiro grupo de atributos, a precisão das instâncias classificadas corretamente é maior.

Tomando os resultados obtidos para o terceiro grupo de atributos para analisar as matrizes de confusão, uma vez que com este grupo de atributos os resultados foram os melhores, a matriz de confusão obtida pelo classificador ao se usar apenas uma camada escondida (Quadro 4), mostra que nesta situação o classificador não consegue distinguir as classes consideradas, classificando todas as instâncias analisadas em apenas 5 classes diferentes, obtendo 100% de erro para as instâncias pertencentes a maioria das classes, como por exemplo a classe Asfalto, que teve todas as instâncias classificadas como Telhado Azul ou Telhado Escuro.

Quadro 4: Matriz de confusão – Classificador MLP – 1 camada escondida

a	b	c	d	e	f	g	h	i	j	k	l	<-- classified as
0	0	0	0	0	14	0	277	0	0	0	0	a = ASFALTO
0	0	0	0	0	652	0	0	0	0	0	0	b = PISCINA
0	0	0	0	0	0	1707	773	0	0	86	0	c = SOLO_EXPOSTO
0	0	0	0	0	13	1	1816	0	0	0	0	d = SOMBRA
0	0	0	0	0	70	16	1214	0	0	3	5	e = EDIFICIO
0	0	0	0	0	5628	0	10	0	0	0	0	f = TELHADO_AZUL
0	0	0	0	0	25	1612	2237	0	0	389	38	g = TELHADO_BRANCO
0	0	0	0	0	1	8	4874	0	0	0	0	h = TELHADO_ESCURO
0	0	0	0	0	0	281	62	0	0	409	130	i = TELHADO_MARROM
0	0	0	0	0	0	731	886	0	0	7	0	j = TELHADO_MARROM_ESCURO
0	0	0	0	0	0	606	61	0	0	1524	809	k = VEGETACAO
0	0	0	0	0	0	0	0	0	0	6	5745	l = VEGETACAO_RASTEIRA

Mas à medida que o número de camadas escondidas aumenta a confusão entre as classes diminui. O Quadro 5 mostra a matriz de confusão gerada ao se usar 50 camadas escondidas para prever a classificação utilizando o terceiro grupo de atributos.

O Quadro 5 mostra que as confusões entre as classes são mínimas, e quando acontecem estão relacionadas a alvos que possuem respostas espectrais semelhantes para as bandas do sensor *IKONOS II*, e a instâncias que o MDA não foi possível de distinguir a diferença de altura, provavelmente por estas terem os valores de altura errados.



Quadro 5: Matriz de confusão – Classificador MLP – 50 camadas escondidas

a	b	c	d	e	f	g	h	i	j	k	l	<-- classified as
285	0	0	0	1	0	0	5	0	0	0	0	a = ASFALTO
0	652	0	0	0	0	0	0	0	0	0	0	b = PISCINA
0	0	2556	0	0	0	0	0	1	9	0	0	c = SOLO_EXPOSTO
0	0	0	1830	0	0	0	0	0	0	0	0	d = SOMBRA
0	0	0	0	1246	0	51	11	0	0	0	0	e = EDIFICIO
0	0	0	0	0	5637	1	0	0	0	0	0	f = TELHADO_AZUL
0	0	0	0	2	3	4273	15	2	6	0	0	g = TELHADO_BRANCO
4	0	0	1	1	0	3	4870	1	3	0	0	h = TELHADO_ESCURO
0	0	5	0	0	0	0	0	871	6	0	0	i = TELHADO_MARROM
0	0	16	0	0	0	0	22	2	1584	0	0	j = TELHADO_MARROM_ESCURO
0	0	0	0	0	0	0	1	0	0	2997	2	k = VEGETACAO
0	0	0	0	0	0	0	0	0	0	1	5750	l = VEGETACAO_RASTEIRA

Para saber se o classificador MLP conseguiria classificar todas as instâncias com 100% de precisão, sem que nenhuma instância fosse confundida, foi testado para o terceiro grupo de atributos, a execução do classificador com números maiores de camadas escondidas. Foram então utilizadas 100, 200 e 500 camadas escondidas.

Tabela 13: Precisão instâncias classificadas corretamente para o terceiro grupo de atributos

Número de Camadas Escondidas	Precisão Instâncias Classificadas Corretamente (%)
1	59,23
5	97,61
10	99,01
25	99,40
50	99,46
100	99,41
200	99,50
500	97,47

A Tabela 13 mostra, para o conjunto de dados utilizado, que a precisão de instâncias classificadas corretamente atinge um ponto de estabilização entre 25 e 50 camadas escondidas. Para valores maiores o ganho de precisão é insignificante, e após ultrapassar as 200 camadas escondidas, a tabela indica que ocorre a perda de precisão.

É importante mencionar que à medida que o número de camadas escondidas aumenta o custo computacional para realizar a predição da classificação também aumenta. O tempo de processamento deste classificador também é bem mais alto do que os dois classificadores testados anteriormente, indicando que para o conjunto de dados

utilizado, o classificador MLP demanda mais processamento para obtenção de uma precisão semelhante a do classificador J48, que mostrou ser bem mais rápido para apresentar os resultados.

## 6. Considerações Finais

A quantidade de dados produzidos atualmente é muito maior do que a capacidade humana pode suportar para analisar, sendo necessária a utilização de meios computacionais que auxiliem a extração de informações úteis e de conhecimento no meio do grande volume de dados armazenados.

Neste sentido, a utilização de softwares que realizam processos de mineração de dados é muito importante na obtenção de conhecimentos que não estão explícitos quando analisamos os dados.

Para este trabalho, a utilização do sistema de mineração de dados WEKA para análise dos dados utilizados mostrou que a combinação entre as imagens das bandas do sensor *IKONOS II* e os dados *LiDAR* são bastante eficientes para prever a classificação da cobertura do solo urbano.

Os resultados obtidos com a utilização dos três classificados testados, mostraram que as bandas do sensor *IKONOS II*, quando utilizadas de forma exclusiva, têm maior capacidade para distinguir entre os alvos que compõem a cobertura do solo do que os dados *LiDAR*, quando são utilizados como dados isolados para tal tarefa. Isto indicar que os dados *LiDAR* devem ser utilizados como uma fonte de dados complementar ao uso de imagens ópticas para classificar a cobertura do solo urbano.

A utilização dos dados *LiDAR* para distinção entre os alvos considerados, propiciou uma melhora na precisão das instâncias classificadas corretamente da ordem de 5%. Enquanto as imagens *IKONOS II* contribuíram com praticamente 95% de precisão ao se classificar as instâncias analisadas.

A contribuição dos dados *LiDAR* para classificar a cobertura do solo urbano esta na informação de altura dos alvos que sistema fornece, sendo então uma informação muito importante. Mas os resultados mostraram que o Modelo Digital de Altura (MDA) utilizado continha alguns problemas, pois instâncias de alvos diferentes, com diferença considerável de altura, acabaram sendo confundidas em determinados testes realizados.

Dentre os classificadores testados, o classificador IBK com um vizinho mais próximo, obteve o melhor resultado possível, sendo perfeito para prever a classificação da cobertura do solo urbano, alcançando 100% de precisão na predição da classificação para os dados de testes analisados. Mas o classificador J48, e o classificador MLP com 25 camadas escondidas, aplicados para o terceiro grupo de atributos (bandas *IKONOS II* e dados *LiDAR*), alcançaram mais de 99% de instâncias classificadas corretamente, sendo então classificadores que se mostraram adequados para se trabalhar com o conjunto de dados utilizado.

O classificador J48 se mostrou, empiricamente, ser o classificador com melhor relação custo x benefício, pois o tempo de processamento computacional para apresentar os resultados foi bem inferior aos outros dois classificadores, e a precisão da classificação foi alta como as dos demais.

Para próximos estudos devem-se aplicar os algoritmos testados para classificar o conjunto total de dados e analisar se as precisões da predição da classificação da cobertura do solo urbano serão altas como as obtidas para as amostras de treinamento, além de apresentar os resultados obtidos no formato de imagens, para permitir que a imagem classificada possa ser comparada visualmente com imagens ópticas de alta resolução ou fotos aéreas da área de estudo, possibilitando que a sensibilidade e experiência do intérprete possam ser empregadas para avaliar a qualidade da classificação alcançada.

## 7. Referências Bibliográficas

AHA, D. W.; KIBLER, D.; ALBERT, M. K. Instance-Based Learning Algorithms. **Machine Learning**, v. 6, p. 37-66, 1991.

FAYYAD, U.; PIATETSKY-SHAPIO, G.; SMYTH, P. From Data Mining to Knowledge Discovery in Databases. **American Association for Artificial Intelligence Magazine**. p. 37-54, 1996, Disponível em:  
<<http://www.daedalus.es/fileadmin/daedalus/doc/MineriaDeDatos/fayyad96.pdf>>. Acesso em 16 Set 2010.

HALL, M.; FRANK, E.; HOLMES, G.; PFAHRINGER, B.; REUTEMANN, P.; WITTEN, I. H. The WEKA Data Mining Software: An Update. **SIGKDD Explorations**, v. 11, n. 1, p. 10-18, 2009.

HAN, Jiawei; KAMBER, Micheline. **Data Mining: Concepts and Techniques**. San Francisco: Morgan Kaufmann Publisher, 2006. 743p.

KARPAGAVALLI, S.; JAMUNA, K. S.; VIJAYA, M. S. Machine learning approach for preoperative anaesthetic risk prediction. **International Journal of Recent Trends in Engineering**, v. 1, n. 2, p.19-22, May 2009.

KOHAVI, Ron; QUINLAN, J. Ross. Decision-tree discovery. In: KLOSGEN, Will; ZYTJOW, Jan M. (Ed.). **Handbook of Data Mining and Knowledge Discovery**. Oxford University Press, 2002. cap. 16.1.3, p. 267-276. Disponível em:  
<<http://robotics.stanford.edu/~ronnyk/ronnyk-bib.html>>. Acesso em: 13/09/2010.

KORTING, T. S.; FONSECA, L. M. G.; ESCADA, M. I. S.; CÂMARA, G. GeoDMA: um sistema para mineração de dados de sensoriamento remoto. In: SIMPÓSIO BRASILEIRO DE SENSORIAMENTO REMOTO, 14, 2009, Natal. **Anais...** São José dos Campos: INPE, 2009. p. 7813-7820. Disponível em:  
<<http://marte.dpi.inpe.br/col/dpi.inpe.br/sbsr@80/2008/11.17.21.17/doc/7813-7820.pdf>>. Acesso em: 27 Set 2010.

NORIEGA, Leonardo. **Multilayer Perceptron Tutorial**. pdf. Disponível em:  
<[http://www.cs.sun.ac.za/~kroon/courses/machine\\_learning/lecture5/mlp.pdf](http://www.cs.sun.ac.za/~kroon/courses/machine_learning/lecture5/mlp.pdf)>. Acesso em: 14 setembro 2010.

PINHO, C. M. D.; ALMEIDA, C. M.; KUX, H. J. H.; RENNÓ, C. D.; FONSECA, L. M. G. Classificação de cobertura do solo de ambientes intra-urbanos utilizando imagens de alta resolução espacial e classificação orientada a objetos. In: ALMEIDA, C. M.; CÂMARA, G.; MONTEIRO, A. M. V. (Ed.). **Geoinformação em urbanismo: cidade real x cidade virtual**. São Paulo: Oficina de Textos, 2007. cap. 8, p. 171-192.

QUINLAN, John Ross. **C4.5: Programs for Machine Learning**. San Mateo: Morgan Kaufmann Publishers, 1993. 302 p.

RAMAKRISHNAN, Naren. C4.5. In: WU, Xindong; KUMAR, Vipin (Ed.). **The Top Ten Algorithms in Data Mining**. Boca Raton: Chapman and Hall/CRC, Taylor and Francis Group, 2009. cap. 8, p. 151-161.

SHIH, F. Y. **Image Processing and Pattern Recognition: Fundamental e Techniques**. New Jersey: John Wiley e Sons, Inc., Hoboken, 2010. 537 p.

SOUSA, R. C. A.; KUX, H. J. H. Comportamento espectral e alvos urbanos: simulação com as bandas espectrais do satélite CBERS. SIMPÓSIO BRASILEIRO DE SENSORIAMENTO REMOTO, 12, 2005, Goiânia. **Anais...** São José dos Campos: INPE, 2005. p. 1-8. Disponível em: <[http://www.obt.inpe.br/cbers/cbers\\_XIISBSR/351\\_comp\\_esp\\_urbano.pdf](http://www.obt.inpe.br/cbers/cbers_XIISBSR/351_comp_esp_urbano.pdf)>. Acesso em: 05 Out 2010.

STEINBACH, Michael; TAN, Pang-Ning. kNN: k-Nearest Neighbors. In: WU, Xindong; KUMAR, Vipin (Ed.). **The Top Ten Algorithms in Data Mining**. Boca Raton: Chapman and Hall/CRC, Taylor and Francis Group, 2009. cap. 8, p. 151-161.

TOMÁS, L. V. **Inferência populacional urbana baseada no volume de edificações residenciais usando imagens IKONOS-II e dados LIDAR**. 2010, 120 p. Tese (Doutorado em Sensoriamento Remoto) – Instituto Nacional de Pesquisas Espaciais (INPE), São José dos Campos, 2010.