

**UNIVERSIDADE FEDERAL DE SANTA MARIA
CENTRO DE CIÊNCIAS RURAIS
PROGRAMA DE PÓS-GRADUAÇÃO EM CIÊNCIA DO SOLO**

**MAPEAMENTO DIGITAL DE SOLOS:
METODOLOGIAS PARA ATENDER A DEMANDA
POR INFORMAÇÃO ESPACIAL EM SOLOS**

TESE DE DOUTORADO

Alexandre ten Caten

Santa Maria, RS, Brasil

2011

**MAPEAMENTO DIGITAL DE SOLOS: METODOLOGIAS
PARA ATENDER A DEMANDA POR INFORMAÇÃO
ESPACIAL EM SOLOS**

Alexandre ten Caten

Tese apresentada ao Curso de Doutorado do Programa de Pós-Graduação em
Ciência do Solo, Área de Concentração em Processos Físicos e Morfogenéticos
do Solo, da Universidade Federal de Santa Maria (UFSM, RS), como requisito
parcial para obtenção do grau de
Doutor em Ciência do Solo.

Orientador: Prof. Dr. Ricardo Simão Diniz Dalmolin

Santa Maria, RS, Brasil

2011

T289m Ten Caten, Alexandre
Mapeamento digital de solos : metodologias para atender a demanda por
informação espacial em solos / por Alexandre Ten Caten. – 2011.
106 f. ; il. ; 30 cm

Orientador: Ricardo Simão Diniz Dalmolin
Tese (doutorado) – Universidade Federal de Santa Maria, Centro de Ciências
Rurais, Programa de Pós-Graduação em Ciência do Solo, RS, 2011

1. Mapeamento digital de solos 2. Pedometria 3. Levantamento de solos
4. Wavelet 5. Mapa cloroplético 6. Árvore de decisão I. Dalmolin, Ricardo
Simão Diniz II. Título.

CDU 631.4:528.7/9

Ficha catalográfica elaborada por Cláudia Terezinha Branco Gallotti – CRB 10/1109
Biblioteca Central UFSM

© 2011

Todos os direitos autorais reservados a Alexandre ten Caten. A reprodução de partes ou do todo deste trabalho só poderá ser feita mediante a citação da fonte.

Endereço: Rua Arsênio Machado Soares, 90 ap 204. CEP: 97110-110

Fone (0xx) 55 9945 3935; End. Eletr: acaten@yahoo.com.br

**Universidade Federal de Santa Maria
Centro de Ciências Rurais
Programa de Pós-Graduação em Ciência do Solo**

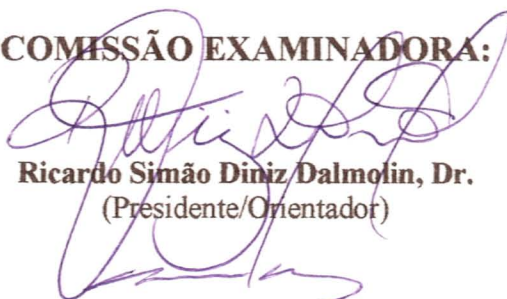
A Comissão Examinadora, abaixo assinada,
aprova a Tese de Doutorado

**MAPEAMENTO DIGITAL DE SOLOS: METODOLOGIAS PARA
ATENDER A DEMANDA POR INFORMAÇÃO ESPACIAL EM SOLOS**

elaborada por
Alexandre ten Caten

como requisito parcial para obtenção do grau de
Doutor em Ciência do Solo

COMISSÃO EXAMINADORA:



Ricardo Simão Diniz Dalmolin, Dr.
(Presidente/Orientador)

Fabricio de Araújo Pedron, Dr. (UFSM)



Jean Paolo Gomes Minella, Dr. (UFSM)



Ivan Luiz Zilli Bacic, PhD. (EPAGRI)



Elvio Giasson, Dr. (UFRGS)

Santa Maria, 07 de novembro de 2011.

DEDICATÓRIA

Aos meus pais pela pergunta: '- Já fez os tema guri?'.

AGRADECIMENTOS

Às Forças Positivas do Universo que conspiraram para que eu tivesse saúde para concluir essa etapa de minha vida.

À Universidade Federal de Santa Maria por possibilitar o acesso ao conhecimento e as oportunidades de crescimento pessoal em todas as etapas de minha vida acadêmica na instituição.

À Coordenação de Aperfeiçoamento de Pessoal de Nível Superior (CAPES) e ao Conselho Nacional de Desenvolvimento Científico e Tecnológico (CNPq) pelo financiamento à pesquisa e pela disponibilização do portal de periódicos.

Ao Programa de Pós-Graduação em Ciência do Solo da Universidade Federal de Santa Maria pela oportunidade de cursar o doutorado.

Ao professor Ricardo Simão Diniz Dalmolin pela oportunidade e exemplo de ser humano e profissional.

Aos membros da banca examinadora da defesa de tese de doutorado Ivan Luiz Zilli Bacic, Elvio Giasson, Jean Paolo Gomes Minella e Fabrício de Araújo Pedron pela disponibilidade e contribuições para a melhoria do trabalho.

Aos professores do PPGCS pelos ensinamentos e ao funcionário Heverton Heinz pela prestatividade.

Ao Instituto Federal Farroupilha por viabilizar que eu continuasse meus estudos até a conclusão do doutorado.

Aos colegas e estudantes do Instituto Federal Farroupilha campus Júlio de Castilhos pela compreensão para com minha ausência.

Ao estudante do curso Tecnólogo em Geoprocessamento Luis Fernando Chimelo Ruiz pela valorosa ajuda com a produção dos dados.

Aos colegas do setor de Pedologia e do PPGCS pelas inúmeras trocas.

À Leo por ter estado ao meu lado durante todos esses anos de minha formação acadêmica.

Obrigado

EPÍGRAFE

But there is still *terra incognita* ahead of soil science ...
(Sabine Grundwald)

RESUMO

Tese de Doutorado
Programa de Pós-Graduação em Ciência do Solo
Universidade Federal de Santa Maria

MAPEAMENTO DIGITAL DE SOLOS: METODOLOGIAS PARA ATENDER A DEMANDA POR INFORMAÇÃO ESPACIAL EM SOLOS

AUTOR: ALEXANDRE TEN CATEN

ORIENTADOR: RICARDO SIMÃO DINIZ DALMOLIN

Data e Local da Defesa: Santa Maria, 07 de novembro de 2011.

O solo é cada vez mais reconhecido como tendo um importante papel nos ecossistemas, assim como para a produção de alimentos e regulação do clima global. Por esse motivo, a demanda por informações relevantes e atualizadas em solos está em uma crescente. O Mapeamento Digital de Solos (MDS) possibilita gerar essas informações demandadas em diferentes resoluções espaciais e com indicadores de qualidade associados. O objetivo deste estudo foi analisar as principais abordagens metodológicas utilizadas nos mapeamentos digitais de classes de solos através de uma revisão de literatura dos trabalhos nacionais, assim como propor procedimentos para a análise dos dados a serem utilizados em projetos de mapeamento digital de classes de solos. O emprego de técnicas de MDS para o mapeamento de classes de solos é recente no país, a primeira publicação nesse sentido ocorreu apenas em 2006. Entre as funções preditivas utilizadas predomina o emprego da técnica de regressões logísticas. Quanto à avaliação da qualidade dos modelos preditivos o emprego da matriz de erros e do índice kappa têm sido os procedimentos mais usuais. O emprego da transformada wavelet mostrou-se como uma metodologia de grande potencial para a análise da resolução espacial de máxima variabilidade de atributos de terreno a serem usados em projetos de MDS. A metodologia proposta de exclusão dos dados oriundos de covariáveis ambientais localizadas na bordas dos polígonos de solos possibilitou a geração de modelos por Árvore de Decisão (AD) menos complexos e mais precisos. Assim como o volume de dados necessários para o treinamento de modelos preditivos por AD está entre cinco e 15% do conjunto total de dados como mostrou este estudo. Observações coletadas a campo indicaram uma acurácia dos mapas preditos próxima a 70% para os modelos oriundos dessas densidades de amostragem.

Palavras-chave: ondaleta, árvore de decisão, pedometria, levantamento de solos.

ABSTRACT

Doctoral Thesis
Graduation Program in Soil Science
Federal University of Santa Maria

DIGITAL SOIL MAPPING: METHODS TO MEET THE DEMAND FOR SOIL SPATIAL INFORMATION

AUTHOR: ALEXANDRE TEN CATEN

ADVISER: RICARDO SIMÃO DINIZ DALMOLIN

Defense Place and Date: Santa Maria, November 07th, 2011.

Soil has increasingly being recognized as having an important role in ecosystems as well as for food production and global climate regulation. For this reason, the demand for relevant and updated information on soil is increasing. Digital Soil Mapping (DSM) provides this information at different spatial resolution with associated quality indicators. The aim of this study was to analyze the main methodological approaches used for DSM of soil classes through a literature review of national researches and to propose procedures for data analysis in DSM projects of soil classes. The use of DSM techniques for mapping soil classes in Brazil is recent, the first publication on this subject occurred only in 2006. Among the predictive functions, logistic regressions is the predominantly used technique. Quality evaluation of the predictive models employed error matrix and kappa index in most cases. The use of wavelet transform proved to be a methodology of great potential for analyzing the spatial resolution of terrain attributes maximum variability. The proposed methodology of data exclusion for environmental covariates located too near at the border of soil classes polygons has enabled the generation of less complex and more accurate Decision Tree (DT) models. It was also shown that the amount of data required for DT model training is between five and 15% of the total data set. Collected field observations indicated a predicted accuracy close to 70% for DT models produced by those sampling densities.

Key words: wavelet, decision tree, pedometric, soil survey

SUMÁRIO

INTRODUÇÃO	11
ARTIGO 1- Mapeamento digital de classes de solos: características da abordagem brasileira	14
RESUMO	14
ABSTRACT	15
INTRODUÇÃO	15
Demandas para a ciência do solo	15
Mapeamento Digital de Solos (MDS)	16
DESENVOLVIMENTO	18
Metodologia adotada	18
O MDS no Brasil é recente	18
Modelos preditivos empregados	20
Fatores de formação empregados	22
Distribuição e dimensão dos estudos	23
Procedimentos metodológicos	25
Qualidade e disponibilidade da informação	27
Expectativas para o MDS no Brasil	30
CONSIDERAÇÕES FINAIS	32
REFERÊNCIAS	33
ARTIGO 2 - Resolução espacial de um modelo digital de elevação definida pela função wavelet	41
Resumo	41
Abstract	41
INTRODUÇÃO	42
MATERIAL E MÉTODOS	45
Área de estudo	45
Atributos de terreno	46
Análise wavelet	46
RESULTADOS E DISCUSSÃO	49
Atributos de terreno	49
Transformada wavelet	50
CONCLUSÕES	59

AGRADECIMENTOS	59
REFERÊNCIAS	59
ARTIGO 3 - Mapeamento digital de solos: estratégia de pré-processamento de dados.	62
RESUMO	62
SUMMARY	62
INTRODUÇÃO	62
MATERIAL E MÉTODOS	65
Mapa de solos para treinamento	65
Covariáveis preditoras	66
Árvore de Decisão (AD)	67
Mapa de solos	68
Qualidade dos modelos e mapas	68
RESULTADOS E DISCUSSÃO	68
CONCLUSÕES	75
AGRADECIMENTOS	75
LITERATURA CITADA	75
ARTIGO 4 - Mapeamento de classe de solos por árvore de decisão: impacto do volume de dados	78
Resumo	78
INTRODUÇÃO	78
MATERIAL E MÉTODOS	80
Área de estudo	80
Covariáveis ambientais	81
Árvores de decisão	84
Mapa de solos	84
Qualidade dos modelos e mapas	85
RESULTADOS E DISCUSSÃO	85
Estudo da multicolinearidade entre os preditores	85
Complexidade do modelo por árvore de decisão	89
Mapas de solos	92
Qualidade dos mapas de solos	94
CONCLUSÃO	96
AGRADECIMENTOS	97
REFERÊNCIAS	97

DISCUSSÃO	100
REFERÊNCIAS BIBLIOGRÁFICAS	105

INTRODUÇÃO

O Mapeamento Digital de Solos (MDS), em nível nacional, está ultrapassando a posição de nova possibilidade para pesquisa e caminha para sua consolidação como técnica de auxílio ao mapeamento convencional de solo. Esse último há muito tido e dito como insuficiente para fornecer as informações de forma quantitativa e qualitativa de acordo com as demandas.

Algumas iniciativas em termos nacionais têm demonstrado a maior atenção dos cientistas brasileiros a esta linha de pesquisa. Entre elas podem ser citadas a criação da Comissão Especializada em Pedometria dentro da divisão Solo no Espaço e no Tempo pela Sociedade Brasileira de Ciência do Solo (SBCS); a formação da RedeMDS, sob os auspícios da EMBRAPA Solos, visando à reunião de pesquisadores interessados em desenvolver essa área do conhecimento no país; além do aumento das pesquisas em MDS realizadas no Brasil, que tem se refletido em um maior número de publicações científicas relacionadas a esse tema.

Este estudo parte do pressuposto de que o MDS já é também uma realidade nacional e, como já foram realizadas as primeiras demonstrações em dissertações e teses, faz-se necessário agora instrumentalizar os pesquisadores de procedimentos metodológicos para a aplicação dessa técnica de mapeamento nas situações em que a informação espacial em solos seja necessária. Esta tese de doutorado está dividida em quatro capítulos, os quais representam indagações metodológicas dos autores em projetos de MDS.

O objetivo geral deste estudo foi avaliar procedimentos metodológicos a serem empregados em projetos de mapeamento digital de classes de solos. Esse objetivo geral foi apoiado em uma revisão de literatura dos trabalhos nacionais relacionados ao mapeamento digital de classes de solos. A revisão de trabalhos nacionais teve como intuito identificar métodos os quais necessitem de uma maior contribuição por parte da pesquisa em MDS.

Os objetivos específicos do presente trabalho foram: (i) avaliar o potencial da transformada wavelet no estudo da variabilidade espacial e na definição da resolução espacial de atributos de terreno potencialmente empregáveis ao MDS; (ii) testar um método de exclusão dos dados de covariáveis preditoras localizados nos locais de maior variabilidade na transição entre os polígonos de classes de solos; (iii) verificar a influência do volume de dados da área de treinamento utilizado para gerar os modelos preditivos por Árvore de Decisão (AD) aplicados ao MDS.

O primeiro capítulo traz a revisão de literatura dos resultados de pesquisa realizados por pesquisadores brasileiros no mapeamento digital de classes de solos. A importância que esse tema vem ganhando entre os pesquisadores se reflete no número de trabalhos produzidos. São cada vez mais frequentes as apresentações na forma de pôsteres e orais nos eventos nacionais relacionados à ciência do solo, assim como no número de publicações em periódicos científicos. Até onde se pode constatar, uma revisão de literatura com esse enfoque é inédita no país, o que demonstra a importância de se reunir e discutir as principais características dos trabalhos já realizados. Essa revisão irá possibilitar uma perspectiva mais ampla dos caminhos que o mapeamento automatizado de solos vem seguindo no país, além de nortear trabalhos e demandas futuras.

Na sequência da tese, são apresentados os resultados de um estudo no qual se buscou aplicar a análise Wavelet (português Ondaleta) para explorar a variabilidade espacial de atributos de terreno. Entre as covariáveis preditoras mais aplicadas no MDS estão os atributos de terreno ligados ao fator relevo de formação do solo. A geração de alguns desses atributos de terreno é sensível à resolução espacial em que se está trabalhando. Esse é o caso, por exemplo, do índice de posição topográfica que é utilizado para classificar a paisagem em classes morfológicas. Neste capítulo da tese, foi aplicada a análise Wavelet nos atributos de terreno elevação, declividade, curvatura em perfil e índice de umidade topográfica, visando definir a resolução espacial mais apropriada para representar a variabilidade de cada um dos quatro atributos na área do estudo, além de demonstrar a aplicação dessa forma de análise espacial de dados. A aplicação dessa técnica para o estudo de dados em ciência do solo no Brasil ainda é incipiente como se pode constatar na fase de revisão bibliográfica do segundo capítulo.

No terceiro capítulo, são apresentados os resultados de uma estratégia de pré-processamento dos dados a serem utilizados no mapeamento digital de classes de solos. Os mapas de classes de solos existentes são comumente utilizados para treinar os modelos preditivos. Nos mapas de solos, o delineamento das unidades de mapeamento é feito através da interpretação visual de pares estereoscópicos. Intrínseco a esse método, está presente a subjetividade da posição mais adequada para o limite entre os polígonos de solos. A utilização desses polígonos como referência para o treinamento de modelos preditivos implicará a adição de informações desviadas. O estudo demonstra o efeito da utilização de diferentes conjuntos de dados na qualidade de mapas de classes de solos, com e sem pré-tratamento. O pré-tratamento proposto é a utilização da função *buffer* para direção interna aos polígonos

originais. Com isso, informações próximas às bordas dos polígonos não são utilizadas para treinar os modelos preditivos por AD.

No quarto e último capítulo, foi avaliado o efeito do volume de dados utilizado para gerar os modelos preditivos. A disponibilização de modelos digitais de terreno e imagens com resoluções espaciais cada vez maiores faz com que um grande volume de dados esteja disponível e tenha de ser manipulado no MDS. O intuito deste estudo foi avaliar, dentro de um universo de informações disponíveis, qual o volume de dados mais adequado para a geração de modelos preditivos por AD a serem aplicados no mapeamento de classes de solos, visando evitar a manipulação desnecessária de dados, assim como evitar que um conjunto muito pequeno de dados implique a impossibilidade dos modelos preditivos em capturar a complexidade inerente à distribuição espacial de classes de solos.

ARTIGO 1- Mapeamento digital de classes de solos: características da abordagem brasileira

Digital soil mapping: characteristics of the brazilian approach

-REVISÃO BIBLIOGRÁFICA-**

RESUMO

O solo é cada vez mais reconhecido como tendo um importante papel nos ecossistemas, assim como para a produção de alimentos e regulação do clima global. Por esse motivo, a demanda por informações relevantes e atualizadas em solos está em uma crescente. Pesquisadores em ciência do solo estão sendo demandados a gerar informações em diferentes resoluções espaciais e com qualidade associada dentro do que está sendo chamado de Mapeamento Digital de Solos (MDS). Devido ao crescente número de trabalhos relacionados ao MDS, faz-se necessário reunir e discutir as principais características dos estudos relacionados ao mapeamento automatizado de classes de solos no Brasil, o que irá possibilitar uma perspectiva mais ampla dos caminhos, além de nortear trabalhos e demandas futuras. O mapeamento de classes de solos empregando técnicas de MDS é recente no país, a primeira publicação nesse sentido ocorreu apenas em 2006. Entre as funções preditivas utilizadas, predomina o emprego da técnica de regressões logísticas. O fator de formação relevo foi empregado na totalidade dos estudos revisados. Quanto à avaliação da qualidade dos modelos preditivos, constatou-se que a matriz de erros e o índice kappa têm sido os procedimentos mais usuais. A consolidação dessa abordagem automatizada como ferramenta auxiliar ao mapeamento convencional passa pelo treinamento dos jovens pedólogos para a utilização de tecnologias da geoinformação e de ferramentas quantitativas dos aspectos de variabilidade do solo.

Palavras-chave: pedometria, classes de solos, levantamento de solos.

25 **ABSTRACT**

26 Soil is increasingly being recognized as having an important role in ecosystems, as
27 well as for food production and global climate regulation. For this reason, the demand for
28 relevant and updated soil information is increasing. Soil science researchers are being
29 demanded to produce information in different spatial resolutions with associated quality in
30 what is being called Digital Soil Mapping (DSM). Due to an increasing number of papers
31 related to the DSM in Brazil, it is necessary to discuss the main characteristics of those
32 studies related to the automated mapping of soil classes, which will enable a broader
33 perspective of the subject and guide future works and demands. The mapping of soil classes
34 using DSM techniques is recent in the country, the first publication in this topic occurred just
35 in 2006. Among the predictive functions the predominant is logistic regression. The soil
36 formation factor relief was used in all studies reviewed. Quality of predictive models was
37 evaluated employing error matrix and kappa which were the most common procedures. The
38 consolidation of this automated approach as an auxiliary tool to the conventional soil mapping
39 will demand training of young soil scientists to use geoinformation technologies and
40 quantitative tools to handle aspects of soil variability.

41 **Key words:** pedometric, soil classes, soil survey.

42

43 **INTRODUÇÃO**

44 **Demandas para a ciência do solo**

45 Para alimentar a população mundial, em 2050, mais de um bilhão de hectares nativos
46 terão de ser convertidos em áreas agricultáveis, caso os modos de produção continuem como
47 os atuais. Visando tornar a agricultura sustentável, é preciso monitorar os efeitos de sua
48 prática em todo o planeta. Historicamente, estratégias agrícolas têm focado em lucratividade e
49 produção. É necessário que esta estratégia seja convertida em sustentabilidade ambiental,
50 segurança alimentar, saúde humana e no bem estar econômico e social. Para que isso seja

51 alcançado, um conjunto de dados, em diferentes escalas, precisa ser coletado e disponibilizado
52 aos agricultores, pesquisadores e políticos (Sachs et al., 2010).

53 O solo é cada vez mais reconhecido como tendo um importante papel não só para os
54 ecossistemas, como também para a produção de alimentos e regulação do clima global. Por
55 esse motivo, a busca por informações relevantes e atualizadas em solos está em uma
56 crescente. Contudo, a comunicação da informação acerca do solo é um desafio devido à
57 divergência de termos, desatualização, generalização e imprecisão dos métodos (Sanchez et
58 al., 2009). A pesquisa em solo tem respondido a suas próprias dúvidas e não tem
59 disponibilizado respostas quantitativas a antigas perguntas, de maneira a ser diretamente
60 utilizada por usuários do solo e tomadores de decisão (Hartemink & McBratney, 2008).

61 Para Hartemink & McBratney (2008), o conhecimento e a pesquisa em solos vêm
62 sendo mais valorizados na medida em que será necessário alimentar 8 bilhões de habitantes,
63 produzir energia a partir da agricultura e aumentar a produção animal embora, mundialmente,
64 a informação sobre o solo não seja acurada ou digitalmente disponível nem atualizada. Entre
65 as oportunidades para a ciência do solo, estão as várias técnicas e métodos já disponíveis aos
66 cientistas de solos, além dos softwares e hardwares que necessitam serem explorados e
67 melhorados. Uma das oportunidades é o Mapeamento Digital de Solos (MDS). No MDS,
68 amostras de solos e covariáveis ambientais são coletadas por metodologias tradicionais ao
69 nível do solo, em aviões ou por satélites. Esses dados são empregados em modelos preditivos
70 que possibilitam gerar novas informações com estatísticas de qualidade associadas (Minasny
71 et al., 2008).

72 **Mapeamento Digital de Solos (MDS)**

73 O termo Pedometria foi criado pelo pesquisador Alex McBratney da Universidade de
74 Sidney para descrever o estudo quantitativo da variação do solo (Burrough et al., 1994). A
75 predição e o Mapeamento Digital de Solos (MDS) – *Digital Soil Mapping* – tiveram suas

76 bases estabelecidas por McBratney et al. (2003) e definidas por Lagacherie & McBratney
77 (2007) como “a criação e a população de sistemas de informação espacial de solos por meio
78 de modelos numéricos visando inferir as variações espaciais e temporais de classes e
79 propriedades do solo, a partir de observações, conhecimento e dados de covariáveis
80 ambientais relacionados”.

81 Como principal aplicação dessa abordagem está a predição, por meio de modelos
82 matemáticos, das classes e propriedades de solos e o mapeamento digital dos resultados de
83 forma contínua e espacial, criando a possibilidade de organizar um amplo conjunto de dados
84 para análise e interpretações em qualquer época, não sendo o mapa o único produto
85 (McBratney et al., 2003). Essa abordagem iniciou nos anos 70 e teve grande desenvolvimento
86 nos anos 80 devido aos avanços tecnológicos nas áreas de tecnologia da informação,
87 sensoriamento remoto, estatística, modelagem, posicionamento global, sistemas de medida e,
88 mais recentemente, acesso instantâneo à informação através da rede mundial de
89 computadores.

90 Para Sanchez et al. (2009), entre as principais etapas do MDS estão: (i) reunião e
91 unificação das informações disponíveis, representando os fatores de formação como modelos
92 digitais de elevação, mapas de vegetação, clima e geológicos, além de informações
93 disponíveis em relatórios e mapas de solos sobre sua distribuição espacial; (ii) mapeamento
94 de classes e propriedades do solo através da aplicação de relações matemáticas entre essas e
95 os fatores de formação do solo; (iii) aplicação das informações mapeadas com vistas a gerar
96 novas informações, as quais são mais difíceis de serem medidas do que estimadas; (iv) a
97 informação gerada na etapa anterior é confrontada com as atividades antrópicas, gerando
98 mapas de conflitos, riscos e possibilidades e, por fim, (v) a geração de um conjunto de
99 medidas e recomendações baseado no conhecimento da distribuição espacial do solo.

100 A presente revisão bibliográfica justifica-se na medida em que é crescente o número
101 de trabalhos publicados a respeito da aplicação da técnica de mapeamento digital de classes
102 de solos no Brasil. Faz-se necessário reunir e discutir as principais características dos
103 trabalhos já realizados, o que irá possibilitar uma perspectiva mais ampla dos caminhos que o
104 mapeamento automatizado de solos vem seguindo no país, além de nortear trabalhos e
105 demandas futuras.

106

107 **DESENVOLVIMENTO**

108 **Metodologia adotada**

109 Artigos relacionados ao mapeamento digital de classes de solos realizados no Brasil
110 foram reunidos neste estudo. Para que o artigo publicado fosse selecionado para esta revisão,
111 deveria ter aplicado o MDS em um sentido *stricto sensu*. Interpreta-se o *stricto sensu* como o
112 emprego de covariáveis preditoras de uma ou mais funções matemáticas e de um conjunto de
113 dados de treinamento. Abordagens *lato sensu* do MDS, em que não ocorre uma classificação
114 numérica de solos, mas apenas uma delimitação de classes de solos apoiada por sistemas
115 informatizados, como realizado por Sarmiento et al. (2008), não foram incluídas neste estudo.

116 Para este trabalho, efetuou-se uma busca nas plataformas de dados *Scielo*, *Scopus* e
117 *Web of Science*, utilizando palavras-chave relacionadas ao tema ‘mapeamento digital de
118 classes de solos’. A busca não fez distinção entre revistas científicas nacionais ou
119 estrangeiras.

120 **O MDS no Brasil é recente**

121 Seguindo os critérios apresentados, dez estudos foram reunidos e as principais
122 informações foram tabuladas para fomentar a discussão nesta revisão. O mapeamento de
123 classes de solos empregando técnicas de MDS é recente no país (Tabela 1). A primeira
124 publicação nesse sentido ocorreu apenas em 2006 (Giasson et al., 2006), mesmo ano em que
125 foi sediado no Rio de Janeiro o 2º Workshop Global em Mapeamento Digital de Solos

126 (Hartemink et al., 2008). Esse evento pioneiro visou divulgar entre os cientistas brasileiros
127 que estudam solo os métodos e técnicas empregados no MDS.

128 **Tabela 1** – Informações descritivas de dez estudos em mapeamento digital de classes de solos
129 realizados no Brasil.

Estudos	Modelos	Fatores	Atributos usados	Localização	Escala	Área (km ²)
Giasson et al.(2011)	AD	r, s	dec, pla, per, umi, cur, dir, acu, com, pod	Bento Gonçalves (RS)	1:10.000	6,14
ten Caten et al. (2011a)	REG	r, c	ele, dec, pla, per, dis, con, umi, sed, rad	São Pedro do Sul (RS)	1:50.000	874
ten Caten et al. (2011b)	REG	r, c	ele, dec, pla, per, dis, con, umi, sed, rad	São Pedro do Sul (RS)	1:50.000	874
ten Caten et al. (2011c)	REG	r, c	ele, dec, pla, per, dis, con, umi, sed, rad	São Pedro do Sul (RS)	1:50.000	874
Chagas et al.(2010)	NEU	r, s, p, o	ele, dec, pla, umi, asp, geo, arg, oxi, ndv	Bacia hidrográfica do Rio São Domingos (RJ)	1:100.000 a 1:50.000	9,9
Coelho & Giasson (2010)	REG, AD, BAY, CAR, LMT	r, s	dec, pla, per, umi, cur, dir, acu	Ijuí, Bozano e Coronel Barros (RS)	1:50.000	1.018
Crivelenti et al. (2009)	AD	r, p	dec, pla, per, dis, con, geo	Carta topográfica Dois Corregos (SP)	1:100.000	772
Carvalho et al.(2009)	NEB	r, p, o	ele, dec, geo, veg	Mucugê (Ba)	1:50.000	382,2
Figueiredo & Giasson (2008)	REG	r, s	dec, pla, per, dis, umi, dir, acu, com	Ibirubá e Quinze de Novembro (RS)	1:80.000	720
Giasson et al. (2006)	REG	r, s	dec, pla, per, umi, cur, dir, acu, com, pod	Sentinela do Sul (RS)	1:50.000	253

130

131 REG: regressões logísticas múltiplas, AD: árvore de decisão, BAY: classificação Bayes, CAR: classificação
132 hierárquica CART, LMT: classificação hierárquica LMT, NEU: redes neurais, NEB: lógica nebulosa, r: relevo,
133 c: clima, s: solo, p: material de origem, o: organismos, ele: elevação, dec: declividade, pla: curvatura planar, per:
134 curvatura de perfil, dis: distância à drenagem, con: área de contribuição, umi: índice de umidade topográfica,
135 sed: capacidade de transporte de sedimento, cur: curvaturas combinadas, dir: direção de fluxo, acu: fluxo
136 acumulado, com: comprimento de fluxo, pod: índice de poder de córrego, asp: aspecto, geo: geologia, arg: índice
137 de argila, oxi: óxido de ferro, ndv: NDVI, veg: vegetação, rad: radiação relativa disponível.
138

139 Esse despertar mais recente dos cientistas brasileiros para o MDS pode estar
140 relacionado (i) à disponibilização bem mais tardia de software e hardware aqui no país; (ii) ao
141 conservadorismo de muitos pedólogos que relutam em utilizar sistemas automatizados
142 capazes de contribuir para o mapeamento de solos; (iii) à carência de pessoal qualificado para
143 o emprego da tecnologia da informação na ciência do solo; (iv) à popularização mais recente
144 no Brasil de tecnologias, como sistema de posicionamento global e sensoriamento remoto,
145 que possibilitam gerar dados sobre o ambiente visando a geração dos modelos.

146 Por outro lado, o mapeamento de classes de solos empregando técnicas automatizadas
147 já vem sendo praticado há muitos anos em outros países. Webster & Burrough (1972)
148 demonstraram, há quase quatro décadas, o emprego de estatística multivariada aplicada em
149 sistemas informatizados para mapear classes de solos no sul da Inglaterra.

150 Estudos relacionados ao MDS estão se tornando mais frequentes. Em 2011, foram
151 quatro trabalhos publicados (Giasson et al., 2011; ten Caten et al., 2011a; ten Caten et al.,
152 2011b; ten Caten et al., 2011c). Com a criação da Comissão Especializada em Pedometria
153 dentro da divisão Solo no Espaço e no Tempo pela Sociedade Brasileira de Ciência do Solo
154 (SBCS), haverá uma maior divulgação do MDS entre cientistas do solo no Brasil.

155 **Modelos preditivos empregados**

156 Entre as funções matemáticas utilizadas, predomina o emprego da técnica de
157 regressões logísticas (Tabela 1). Através de uma relação linearizada entre covariáveis
158 predictoras e classes de solos, as regressões logísticas geram um valor de pertinência para cada
159 classe de solos a ser mapeada sobre a paisagem (Giasson et al., 2006). Esse valor de
160 probabilidade de ocorrência é posteriormente avaliado entre todas as classes de solos para um
161 mesmo ponto da paisagem, sendo que entre elas a de maior valor é atribuída à coordenada em
162 análise, gerando-se um mapa de classe de solos. As regressões logísticas apresentam como
163 desvantagem uma sensibilidade à proporção relativa entre classes de solos no conjunto de

164 treinamento dos modelos (ten Caten et al., 2011a) e uma relativa complexidade para a
165 interpretação dos parâmetros estatísticos dos modelos logísticos gerados (Kempen et al.,
166 2009).

167 Metodologias que busquem relacionar fatores de formação e classes de solos de forma
168 mais similar ao raciocínio do pedólogo têm sido aplicadas pelos estudos que empregam a
169 técnica de mineração de dados por Árvore de Decisão (AD) (Crivelenti et al., 2009; Giasson
170 et al., 2011). O emprego de AD é fruto de pesquisas que buscam aplicar técnicas robustas
171 para a extração de padrões em grandes conjuntos de dados (Witten & Frank, 2005). A
172 abordagem por AD tem sido empregada por apresentar a vantagem de possibilitar a expressão
173 das relações solo-paisagem de maneira explícita (Kheir et al., 2010a). Além de permitir o
174 agrupamento e a busca por padrões, a AD possibilita o entendimento de como esses dados são
175 inter-relacionados (Kheir et al., 2010b). Nos estudos até aqui realizados no país, não se tem
176 dado uma ênfase à análise das regras de decisão geradas durante a modelagem. Isso permite
177 afirmar que a AD vem sendo empregada mais pela sua robustez como técnica preditiva, com
178 mais ênfase ao mapa final gerado, do que ao seu potencial em explicitar e esclarecer as
179 relações entre fatores de formação e classes de solos.

180 Entre as estratégias referidas como ‘caixa preta’ (Ballabio, 2009), a técnica de redes
181 neurais artificiais (RNA) foi empregada por Chagas et al. (2010). O uso de redes neurais
182 possibilita uma grande acurácia na predição, contudo, os problemas aparecem quando a rede
183 precisa ser aplicada como conhecimento formalizado devido à complexidade da rede e dos
184 pesos empregados (Qi & Zhu, 2003). Kheir et al. (2010b) concordam com essa dificuldade
185 para a implementação das RNA. Para esses autores, as redes neurais não possibilitam um
186 modelo de fácil compreensão que permita os pesquisadores ter um entendimento completo da
187 natureza dos dados que estão sendo analisados. Também entre as estratégias referidas como

188 modelo ‘caixa preta’ está o emprego de Máquinas de Vetores de Suporte (Ballabio, 2009), o
189 qual ainda não foi empregado em nenhum estudo no Brasil.

190 **Fatores de formação empregados**

191 O fator de formação relevo foi empregado na totalidade dos estudos revisados (Tabela
192 1). Possivelmente isso se deva a sua ampla disponibilidade através de Modelo Digital de
193 Elevação (MDE) oriundo do SRTM (*Shuttle Radar Topography Mission*) e de cartas
194 topográficas. Também contribui para sua aplicação a possibilidade de derivar a partir desse
195 um maior número de covariáveis preditoras; a clara relação existente entre o relevo e o padrão
196 de distribuição espacial das classes de solos, além da resolução espacial dos arquivos raster
197 associados a este fator.

198 Uma das principais fontes de informação para o MDS é o MDE. MDE derivados de
199 cartas topográficas e das plataformas SRTM e ASTER necessitam ser avaliados, bem como
200 novas fontes de dados de maior resolução espacial, a exemplo da plataforma TerraSAR-X, os
201 quais carecem de investigação sobre sua aplicabilidade. O posicionamento por satélite, que
202 tem sido realizado a partir da constelação americana (GPS) e russa (GLONASS), apresentará
203 ganhos de qualidade no posicionamento com a disponibilização dos sistemas europeu
204 (Galileo) e chinês (Compass) em um futuro breve.

205 Os fatores material de origem e clima foram utilizados com menor frequência (Tabela
206 1). Isso decorre possivelmente da generalização dos mapas disponíveis dessas informações.
207 Os mapas geológicos amplamente disponíveis para todo o território nacional encontram-se na
208 escala 1:1.000.000 (CPRM, 2011). Já a informação acerca do clima, além da questão da falta
209 de resolução nos mapas disponíveis, existe ainda o fato de que os solos tenham se formado
210 em tempos pretéritos, podendo haver uma baixa correlação entre condições climáticas atuais e
211 as classes de solos.

212 O Brasil apresenta desafios bem particulares ao mapeamento digital de classes de
213 solos. Entre os desafios está o uso de técnicas de sensoriamento remoto para gerar
214 informações acerca dos fatores de formação do solo. Mapas, como de uso da terra e índices de
215 vegetação (NDVI), apresentam grandes desafios à investigação devido à ação antrópica que
216 modifica a cobertura natural e à presença de nuvens frequentes em várias regiões desse país
217 predominantemente tropical. Isso ocorre também no uso do sinal refletido pelo solo para os
218 sensores orbitais, que acaba sendo comprometido pela cobertura do solo em áreas agrícolas
219 onde se pratica o plantio direto e/ou pela cobertura vegetal natural em muitas regiões do país.

220 **Distribuição e dimensão dos estudos**

221 Entre os dez estudos analisados, sete foram realizados no Estado do Rio Grande do Sul
222 (Tabela 1). Isso pode ser atribuído à presença de dois grupos de pesquisa que desenvolvem
223 atividades nessa linha, naquele Estado. É importante que essa técnica seja divulgada e
224 avaliada em outras regiões do país. Nesse aspecto, a Comissão Especializada em Pedometria
225 da divisão Solo no Espaço e no Tempo (SBCS) terá um papel muito importante para a
226 consolidação das pesquisas em MDS, em todas as regiões do país. Essa informação contrasta
227 com os vazios de levantamentos de solos reportados por Mendonça-Santos & Santos (2007).
228 Esses autores descrevem que as regiões norte e noroeste, sobretudo grandes áreas de floresta
229 Amazônica, possuem informação apenas em pequena escala sobre os solos da região. Ainda
230 conforme esses autores, isso decorre da dificuldade de acesso, políticas de preservação da
231 floresta e a baixa fertilidade do solo, que têm prejudicado os levantamentos de solos na
232 região.

233 Como a aplicação do MDS implica em uma fase de treinamento dos modelos
234 preditivos, o efetivo auxílio dessa técnica para o aumento do conhecimento a respeito da
235 distribuição espacial dos solos na região norte e noroeste só irá ocorrer a partir da
236 disponibilização dessas informações. Esse é um desafio a ser superado. Coletas a campo e por

237 técnicas não proximais deverão gerar um grande volume de dados, caso se deseje um
238 mapeamento em escalas mais detalhadas.

239 O MDS aplicado à predição de classes de solos no Brasil tem ocorrido
240 predominantemente em nível semidetalhado (1:50.000) (Tabela 1), possivelmente devido ao
241 emprego de Cartas Topográficas para a geração de MDE a fim de produzir atributos de
242 terreno utilizados para treinar os modelos preditivos. Embora uma das prerrogativas à
243 aplicação de técnicas de MDS seja possibilitar o mapeamento a partir de demandas
244 específicas pelo conhecimento da distribuição do solo (Carré et al., 2007), isso não tem sido
245 reportado nas publicações nacionais. Dessa forma, o volume de informações e,
246 conseqüentemente, a escala do mapeamento têm sido definidos por outros critérios que não a
247 demanda pela informação.

248 Mapas de solos, disponíveis para o limite de alguns municípios, têm sido o conjunto
249 de dados mais utilizado para treinamento dos modelos preditivos (Tabela 1). Estudos que
250 utilizam os limites administrativos dos municípios podem ter sua aplicação limitada pelo fato
251 de uma Bacia Hidrográfica (BH) estar distribuída em um conjunto de municípios. Esse
252 mesmo problema ocorre caso seja utilizado dados oriundos de cartas topográficas. Esse fator
253 tem condicionado a extensão territorial da maioria dos trabalhos. Apenas o estudo
254 desenvolvido por Chagas et al.(2010) foi realizado a partir das dimensões de uma BH. As BH
255 têm sido indicadas como unidade básica para questões importantes. No Brasil, a Lei Federal
256 nº9.433/97 estabelece a bacia hidrográfica como unidade territorial para aplicação da Política
257 Nacional de Recursos Hídricos. Possivelmente seria a unidade mais adequada para novos
258 mapeamentos pela técnica de MDS, uma vez que é crescente a demanda por informação nessa
259 unidade, disponibilizando informação para estudos futuros.

260 **Procedimentos metodológicos**

261 Os estudos até aqui realizados utilizaram de diferentes densidades de amostragem
262 (Tabela 2). Como os estudos são realizados a partir de dados no formato matriz (raster), a
263 unidade mínima desse modelo de representação dos dados, o pixel, é utilizado como amostra
264 para ajuste dos modelos preditivos. O conjunto de dados utilizados para treinar os modelos
265 representou desde menos de um por cento (Giasson et al., 2011) até a totalidade das amostras
266 (Carvalho et al., 2009), indicando que existe uma carência pela padronização do número de
267 amostras a ser utilizada na fase de treinamento dos modelos. Estratégias apropriadas para a
268 seleção das observações conduzem a melhores resultados do que quando os modelos são
269 ajustados ao total de observações disponíveis (Schmidt et al.; 2008).

270 A totalidade dos estudos utilizou de MDE, como já mencionado, possivelmente devido
271 a questões ligadas ao fator de formação relevo. Entre as fontes de MDE utilizadas, destaca-se
272 o SRTM (*Shuttle Radar Topography Mission*). Esse dado tem sido utilizado em sua resolução
273 original de 90 m ou no formato Topodata (Valeriano & Rossetti, 2010), que, por sua vez,
274 consiste em uma interpolação dos dados SRTM para 30 m. As curvas de nível presentes nas
275 cartas topográficas também têm sido utilizadas como fontes de dados altimétricos. A
276 utilização de tecnologias como o LIDAR, Scanners 3D e GNSS RTK para a obtenção da
277 informação sobre as nuances do terreno ainda não foi testada no MDS nacional. Essas
278 tecnologias têm um grande potencial para obtenção de MDE de elevada resolução e exatidão,
279 visando ao levantamento preditivo de solos em pequenas áreas e/ou grandes escalas.

280

281

282

283

284

285

286 **Tabela 2** – Aspectos metodológicos do mapeamento digital de classes de solo no Brasil

Estudos	Total de amostras	Amostras /km²	Total amostrado (%)	Fonte do MDE	Resolução (m)	Software
Giasson et al. (2011)	1.333	217	0,54	foto	5	arcv, weka
ten Caten et al. (2011a)	70.000	80	20,00	srtm	50	arcg
ten Caten et al. (2011b)	58.440	67	16,72	srtm	50	arcg
ten Caten et al. (2011c)	70.000	80	20,00	srtm	50	arcg
Chagas et al.(2010)	3.000	303	27,27	cart	30	arcg, erda, java
Coelho & Giasson (2010)	11.000	11	8,75	srtm	90	arcg, weka
Crivelenti et al. (2009)	794.273	1.029	90,00	cart	30	arcg, ilwi, idri, weka
Carvalho et al.(2009)	424.643	1.111	100,00	srtm	30	arcg
Figueiredo & Giasson (2008)	7.200	10	8,10	srtm	90	arcv
Giasson et al. (2006)	7.500	30	25,00	srtm	92	arcv

287 foto: pares de fotografias estéreo, cart: cartas topográficas, srtm: shuttle radar topographic mission, arcg: ArcGIS
 288 9, arcv: ArcView 3.2, ilwi: Ilwis, idri: Idrisi Andes, erda: Erdas, java: Java Neural Network Simulator, weka:
 289 Weka
 290

291 Predomina entre os estudos de MDS realizados no país a aplicação de softwares
 292 proprietários. Os programas de código fonte fechados podem limitar a comunicação dos
 293 procedimentos executados para a obtenção de determinado resultado (Wood, 2009). Entre as
 294 alegações estaria o fato de que esses programas não possibilitam conhecer completamente as

295 rotinas e fórmulas empregadas nos processamento dos dados. Isso impossibilitaria a
296 padronização de procedimentos de geração da informação. Existe uma variedade de
297 programas com possíveis aplicações ao MDS que possuem seu código fonte aberto e uma
298 grande rede de colaboradores para o desenvolvimento, a saber: SAGA, GRASS, ILWIS e o
299 pacote estatístico R.

300 **Qualidade e disponibilidade da informação**

301 Quanto à avaliação da qualidade dos modelos preditivos, a matriz de erros e o índice
302 kappa têm sido os procedimentos mais usuais (Tabela 3). A matriz de erros é uma boa
303 alternativa para determinar a natureza e a frequência dos erros envolvidos. Contudo, caso a
304 matriz seja gerada de maneira inapropriada, sem ser representativa dos dados, então a análise
305 será sem significado (Congalton, 1991). O autor ressalta que os seguintes fatores precisam ser
306 considerados: amostragem no campo, modo de obtenção do mapa a ser testado, correlação
307 espacial, tamanho das amostras e esquema de amostragem.

308 Entre os estudos, em apenas dois casos foi avaliada a acurácia do mapeamento com
309 dados de campo (Chagas et al., 2010; ten Caten et al., 2011b). Predomina entre os estudos
310 apenas a avaliação do potencial de reprodução do mapa original a partir dos modelos
311 preditivos. Essa metodologia compara o mapa predito com um mapa produzido por
312 metodologias clássicas de levantamento de solos. Entretanto, pode não ser uma abordagem
313 eficaz para indicar a qualidade do mapa de solos gerado, principalmente, devido às diferenças
314 entre as metodologias empregadas.

315

316

317

318

319 **Tabela 3** – Características dos estudos revisados quanto a aspectos de qualidade e
 320 disponibilidade da informação gerada.

Estudos	Qualidade	Índice kappa*	NC	Legenda simplificada	Idioma
Giasson et al. (2011)	mat / kap / não	0,52	2°	realizada	Ing
ten Caten et al. (2011a)	mat / kap / não	0,63**	ass, 1°, 2°	não realizada	Port
ten Caten et al. (2011b)	mat / kap / sim	0,46	ass, 2°	não realizada	Port
ten Caten et al. (2011c)	mat / kap / não	0,37	ass, 2°	não realizada	Port
Chagas et al.(2010)	mat / kap / sim	0,73	ass, 4°	realizada	Port
Coelho & Giasson (2010)	mat / kap / não	0,38	2°, 3°, 4°	realizada	Port
Crivelenti et al. (2009)	mat / kap / não	0,43	ass, 3°, 4°, textura	realizada	Port
Carvalho et al.(2009)	visual / não	-	ass, 2°, 3°, 4°	não realizada	Port
Figueiredo & Giasson (2008)	mat / kap / não	0,38***	ass, 3°, relevo, textura	realizada	Port
Giasson et al. (2006)	mat / kap / não	0,36***	ass, 1°, 2°	realizada	Ing

321 NC: nível categórico do SiBCS, mat: matriz de erros, kap: índice kappa, visual: comparação visual entre mapas, /
 322 sim: verificado no campo, / não: não verificado no campo, ass: associações, *: reprodutibilidade do mapa, **:
 323 dado não publicado, ***: sem simplificação da legenda, Ing: Inglês, Port: Português
 324

325 O valor médio do índice kappa entre os estudos realizados no país é de 0,47 (Tabela
 326 3). Esse valor é similar aos valores reportados na literatura internacional, como em Hengl &
 327 Rossiter (2003) de 0,58 em locais montanhosos e de 0,39 para áreas planas, e por Scull et al.

328 (2005) que observaram valores de 0,44 e 0,52 em áreas montanhosas e planas,
329 respectivamente. A abordagem adotada em ambos os estudos, pela estratificação da área a ser
330 mapeada de acordo com características locais predominantes, tem sido apresentada como uma
331 possibilidade de melhorar o desempenho dos modelos preditivos, uma vez que a identificação
332 e separação das áreas de acordo com critérios pré-estabelecidos é facilitada por sistemas de
333 informação geográfica. Assim, possibilita-se a geração de modelos mais regionalizados, com
334 melhor desempenho em cada uma das regiões para os quais foram desenvolvidos, permitindo-
335 se capturar as peculiaridades de formação dos solos. A abordagem estratificada foi testada por
336 Giasson et al. (2006), com um valor de kappa de 0,31 para a abordagem estratificada e de
337 kappa de 0,36 quando não estratificada. Esses autores relataram que a possível causa do pior
338 desempenho da abordagem estratificada tenha sido o modelo digital de elevação de 90 m
339 utilizado no estudo.

340 A predição das classes de solos foi realizada nos quatro níveis categóricos do Sistema
341 Brasileiro de Classificação de Solos (SiBCS). A quase totalidade dos estudos também buscou
342 mapear classes de solos na forma de associações. Alguns estudos buscaram avaliar o potencial
343 das metodologias utilizadas em mapear níveis categóricos mais baixos, como o 3° e o 4°
344 níveis do SiBCS. No entanto, nesses níveis os solos são classificados de acordo com
345 características de mais difícil associação com os fatores de formação do solo, a exemplo de
346 características físicas e químicas no 3° nível e da presença de variações e características
347 extraordinárias para a classificação no 4° nível. Acertos na distinção e na espacialização
348 dessas classes na paisagem podem estar mais sendo fruto do acaso do que do verdadeiro
349 poder preditivo dos modelos.

350 Alguns estudos têm se utilizado da estratégia de simplificar a legenda em relação a um
351 mapa de referência, visando melhorar capacidade preditiva dos modelos. Essa abordagem
352 consiste em mapear classes na forma de associações de solos ou no mapeamento de classes

353 apenas em níveis mais elevados do SiBCS. Tal estratégia melhora os indicadores de acerto em
354 relação a um mapa de referência, contudo, produz mapas generalizados com um grau cada vez
355 maior de variações dentro da unidade de mapeamento, quanto maior a simplificação adotada.

356 Os resultados das pesquisas em MDS realizadas no país são, na sua maioria,
357 publicados em artigos escritos em língua portuguesa (Tabela 3). Como o Brasil é um país com
358 uma diversidade de biomas e com grandes lacunas na cartografia de seus solos, os resultados
359 aqui alcançados, sem dúvida, são de interesse da comunidade internacional. Entre as
360 experiências já relatadas do MDS, predominam as condições de solos temperados e de regiões
361 áridas do globo (McBratney et al., 2003; Scull et al., 2003; Grunwald, 2009). Faz-se
362 fundamental que os estudos dessa linha de pesquisa realizados em condições de clima tropical
363 sejam divulgados em outros países. A internacionalização das pesquisas em MDS realizadas
364 no Brasil passa pela publicação dos resultados no idioma inglês.

365 **Expectativas para o MDS no Brasil**

366 Os custos do MDS não foram mencionados em nenhum dos estudos revisados. Para
367 Grunwald (2009), a pesquisa relacionada aos investimentos necessários para tornar o MDS
368 efetivo na produção de conhecimento espacial acerca de classes e propriedades de solos
369 necessita mais atenção dos pesquisadores. Custos envolvidos no acesso a levantamentos de
370 solos já existentes e coletas de novas informações, assim como dados oriundos de técnicas de
371 sensoriamento proximal ou remoto deverão ser considerados. Além disso, conhecimento de
372 especialistas em solos, estatística, informática e geomática precisa ser avaliado no custo total
373 de projetos em MDS. O tempo investido em um projeto de MDS também deverá ser relatado
374 em futuros estudos, para que novos projetos possam ser adequadamente planejados.

375 Uma padronização da abordagem para a avaliação da qualidade dos mapas gerados
376 será fundamental. Na medida em que dados existentes (mapas de solos) conjuntamente com
377 dados de apoio de fácil aquisição (SRTM), que possuem erros inerentes, são cada vez mais

378 utilizados para o MDS, é preciso que uma metodologia padronizada seja proposta para avaliar
379 e possibilitar comparações entre resultados (Krol, 2008). Espera-se dos sistemas modernos
380 para o mapeamento de solos mais qualidade do que dos mapeamentos convencionais, pois no
381 MDS são empregados dados mais bem definidos e documentados em termos de
382 posicionamento, atualização e integração de dados (Finke, 2007). Para esse autor, é necessário
383 estabelecer padrões de qualidade em termos da conformidade da informação com a realidade
384 do campo (acurácia) e com a incerteza associada à metodologia automatizada que se utiliza
385 para gerar a informação (precisão).

386 A avaliação da qualidade dos mapas de classes de solos a partir dos resultados da
387 matriz de erros e índice kappa em relação ao mapa original de solos pode estar subestimando
388 a verdadeira qualidade dos mapas gerados pelo MDS. Como os mapas convencionais são a
389 representação cartográfica do conhecimento pedológico empregado na sua produção, e o
390 conhecimento tácito do pedólogo está em constante aprimoramento, haverá sempre distorções
391 entre o mapa gerado pelo pedólogo e a verdadeira classe de solo presente na paisagem. De tal
392 forma que a comparação do mapa gerado pela metodologia automatizada com aquele obtido
393 pela metodologia convencional apenas indica a capacidade de reprodução desse último pelo
394 modelo preditivo. A checagem do mapa predito com informações de campo deve ser o
395 verdadeiro indicador de qualidade a ser empregado.

396 Como uma das prerrogativas para a execução de um projeto de MDS é a sua
397 implementação em um sistema de informação geográfica, a disponibilização da informação
398 gerada em formatos vetorial ou matricial é facilitada. Contudo, não é o que tem sido relatado
399 nos estudos realizados no país. Em nenhum dos artigos revisados, os autores relatam a
400 disponibilização dos mapas gerados. Além de a forma digital possibilitar a rápida atualização
401 dos mapas e sua disponibilização para acesso pela internet, essa abordagem permite que
402 metodologias inovadoras de acesso, como visualização 3D (Grunwald et al., 2007) e

403 tecnologias móveis (iPhone) (Beaudette & O'Geen, 2010), sejam utilizadas. Essas estratégias
404 permitirão maior divulgação das informações geradas. Um passo que precisa ser dado pelos
405 pesquisadores brasileiros em MDS.

406

407 **CONSIDERAÇÕES FINAIS**

408 Os pedólogos têm se questionado sobre as motivações para a redução da valorização
409 de sua atividade profissional. Entre os culpados, são elencadas a valorização da pesquisa
410 aplicada em detrimento da básica, os critérios taxonômicos utilizados pela pedologia, os
411 planejadores sem interesse pela atividade do pedólogo, levantamentos pedológicos
412 insuficientes para aplicação imediata em problemas da sociedade (Ker & Novais, 2003). Em
413 contraste com a argumentada recessão verificada em relação aos levantamentos de solos,
414 verifica-se uma tendência de renovação, principalmente, quanto às técnicas de coleta e
415 processamento de dados, formação e desenvolvimento de modelos e ao MDS (Ramos, 2003).

416 Face ao grande desafio de produzir informação sobre a distribuição espacial de classes
417 de solos para todo o território brasileiro, o mapeamento digital de solos será um grande
418 aliado. A abordagem tem se consolidado com um número cada vez maior de artigos
419 publicados em revistas científicas especializadas, bem como com a participação de
420 pesquisadores em publicações internacionais.

421 Há a necessidade de que sejam somados esforços para o treinamento de novos
422 pedólogos aptos a implementar projetos de MDS, bem como em realizar levantamentos de
423 solos por meio das metodologias convencionais. Ambas as abordagens são complementares e
424 necessitam serem desenvolvidas simultaneamente. O levantamento de solos convencional
425 produz as informações que são utilizadas para treinar os modelos e predizer classes e
426 propriedades de solos em áreas não mapeadas. A abordagem automatizada permite que
427 extensas regiões sejam previamente mapeadas, otimizando recursos humanos e financeiros no
428 mapeamento convencional. Além disso, possibilita-se a geração de informação em formatos

429 vetorial e matricial prontamente disponível para projetos ambientais e de adequabilidade de
430 uso.

431 O MDS brasileiro precisa experimentar-se mais na geração de informações sobre
432 propriedades do solo. O mapeamento de classes de solos atende uma gama muito grande de
433 aplicações. Contudo, o mapeamento de propriedades pode ser executado em uma abordagem
434 *demand driven*, focando em clientes e necessidades específicas. Essa abordagem irá mobilizar
435 os recursos estritamente necessários para o conhecimento de uma determinada propriedade do
436 solo, por exemplo, para a predição da matéria orgânica do solo, dispensando a descrição
437 completa de perfis. Também pode ocorrer por meio de verificação a campo, podendo-se gerar
438 mapas de precisão e acurácia associados, indicando o grau de confiabilidade dos dados a
439 serem disponibilizados.

440

441 REFERÊNCIAS

442

443 BALLABIO, C. Spatial prediction of soil properties in temperate mountain regions using
444 support vector regression, **Geoderma**, v.151, n.3-4, p.338-350, 2009. Disponível em:
445 <http://www.sciencedirect.com/science/article/pii/S0016706109001499>. Acesso em: 12 out.
446 2011. doi: 10.1016/j.geoderma.2009.04.022.

447 BEAUDETTE D.E.; O'GEEN A. T. An iPhone Application for On-Demand Access to Digital
448 Soil Survey Information. **Soil Science Society of America Journal**. v.74, n.5, p.1682-1684,
449 2010. Disponível em: [http://casoilresource.lawr.ucdavis.edu/dylan/publications/sssaj-
450 iphone_app.pdf](http://casoilresource.lawr.ucdavis.edu/dylan/publications/sssaj-iphone_app.pdf) Acesso em: 12 out. 2011. doi: 10.2136/sssaj2010.0144N.

451 BURROUGH, P.A.; BOUMA, J.; YATES, S.R. The state of the art in pedometrics.
452 **Geoderma**, v.62, p.311-326, 1994. Disponível em:
453 <http://www.sciencedirect.com/science/article/pii/0016706194900434> Acesso em: 12 out.
454 2011. doi: 10.1016/0016-7061(94)90043-4.

- 455 CARRE, F. et al. Digital soil assessments: Beyond DSM, **Geoderma**, v.142, n.1-2, p.69-79,
456 2007. Disponível em: <http://www.sciencedirect.com/science/article/pii/S0016706107002261>
457 Acesso em: 12 out. 2011. doi: 10.1016/j.geoderma.2007.08.015.
- 458 CARVALHO, C.C.N. et al. Mapa digital de solos: Uma proposta metodológica usando
459 inferência fuzzy. **Revista Brasileira de Engenharia Agrícola e Ambiental**, v.13, n.1, p.46-
460 55, 2009. Disponível em: [http://www.scielo.br/scielo.php?script=sci_arttext&pid=S1415-
461 43662009000100007&lng=en&nrm=iso](http://www.scielo.br/scielo.php?script=sci_arttext&pid=S1415-43662009000100007&lng=en&nrm=iso). Acesso em: 25 mai. 2011. doi: 10.1590/S1415-
462 43662009000100007.
- 463 CHAGAS, C. da S. et al. Atributos topográficos e dados do Landsat7 no mapeamento digital
464 de solos com uso de redes neurais. **Pesquisa Agropecuária Brasileira**, v.45, n.5, p.497-507,
465 2010. Disponível em: [http://www.scielo.br/scielo.php?script=sci_arttext&pid=S1415-
466 43662009000100007&lng=en&nrm=iso](http://www.scielo.br/scielo.php?script=sci_arttext&pid=S1415-43662009000100007&lng=en&nrm=iso). Acesso em: 25 mai. 2011. 10.1590/S1415-
467 43662009000100007.
- 468 COELHO, F.F.; GIASSON, E. Comparação de métodos para mapeamento digital de solos
469 com utilização de sistema de informação geográfica. **Ciência Rural**, v.40, n.10, p.2099-
470 2106, 2010. Disponível em: [http://www.scielo.br/scielo.php?script=sci_arttext&pid=S0103-
471 84782010001000008&lng=en&nrm=iso](http://www.scielo.br/scielo.php?script=sci_arttext&pid=S0103-84782010001000008&lng=en&nrm=iso). Acesso em: 20 fev. 2011. doi: 10.1590/S0103-
472 84782010005000156.
- 473 CONGALTON, R.G. A Review of Assessing the Accuracy of Classifications of Remotely
474 Sensed Data. **Remote Sensing of Environment**, v.37, p.35-46, 1991. Disponível em:
475 <http://www.sciencedirect.com/science/article/pii/003442579190048B> Acesso em: 12 out.
476 2011. doi: 10.1016/0034-4257(91)90048-B.
- 477 CPRM – Serviço geológico do Brasil. Disponível em: <http://www.cprm.gov.br/> Acessado em
478 10 ago. 2011.

- 479 CRIVELENTI et al. Mineração de dados para a inferência de relações solo-paisagem em
480 mapeamentos digitais de solo. **Revista Agropecuária Brasileira**, v.44, n.12, p.1707-1715,
481 2009. Disponível em: [http://www.scielo.br/scielo.php?script=sci_arttext&pid=S0100-](http://www.scielo.br/scielo.php?script=sci_arttext&pid=S0100-204X2009001200021&lng=en&nrm=iso)
482 [204X2009001200021&lng=en&nrm=iso](http://www.scielo.br/scielo.php?script=sci_arttext&pid=S0100-204X2009001200021&lng=en&nrm=iso). Acesso em: 25 mai. 2011. doi: 10.1590/S0100-
483 204X2009001200021.
- 484 FIGUEIREDO, S.R. et al. Uso de regressões logísticas múltiplas para mapeamento digital de
485 solos no planalto médio do RS. **Revista Brasileira de Ciência do Solo**, v.32, p.2779-2785,
486 2008. Disponível em: [http://www.scielo.br/scielo.php?script=sci_arttext&pid=S0100-](http://www.scielo.br/scielo.php?script=sci_arttext&pid=S0100-06832008000700023&lng=en&nrm=iso)
487 [06832008000700023&lng=en&nrm=iso](http://www.scielo.br/scielo.php?script=sci_arttext&pid=S0100-06832008000700023&lng=en&nrm=iso). Acesso em: 20 fev. 2011. doi: 10.1590/S0100-
488 06832008000700023.
- 489 FINKE, P.A. Quality assessment of digital soil maps: producers and users perspectives. In:
490 LAGACHERIE, P.; MCBRATNEY, A.B.; VOLTZ, M., Digital Soil Mapping: An
491 Introductory Perspective. *Developments in Soil Science*, v.31, cap.39, Elsevier, Amsterdam,
492 2007, p.523-541.
- 493 HARTEMINK, A.E. et al. Digital Soil Mapping with Limited Data. Ed. Springer, 2008. 445p.
- 494 HARTEMINK, A.E.; MCBRATNEY, A.B., A soil science renaissance. **Geoderma**, v.148,
495 n.2, p.123-129, 2008. Disponível em:
496 <http://www.sciencedirect.com/science/article/pii/S0016706108002802> Acesso em: 12 out.
497 2011. doi: 10.1016/j.geoderma.2008.10.006.
- 498 HENGL, T.; ROSSITER, D.G. Supervised Landform classification to enhance and replace
499 photo-interpretation in semi-detailed soil survey. **Soil Science Society of America Journal**.
500 v.67, p.1810-1822, 2003. Disponível em:
501 <https://www.soils.org/publications/sssaj/abstracts/67/6/1810>. Acesso em: 14 nov. 2011. doi:
502 10.2136/sssaj2003.1810.

- 503 GIASSON, E. et al. Decision trees for digital soil mapping on subtropical basaltic steepplands.
504 **Scientia Agrícola**, v.68, p.167-174, 2011. Disponível em:
505 http://www.scielo.br/scielo.php?script=sci_arttext&pid=S0103-90162011000200006&lng=en
506 [&nrm=iso](#). Acesso em: 25 mai. 2011. doi: 10.1590/S0103-90162011000200006.
- 507 GIASSON, E. et al. Digital soil mapping using multiple logistic regression on terrain
508 parameters in southern Brazil. **Scientia Agrícola**, v.63, p.262-268, 2006. Disponível em:
509 <http://www.scielo.br/scielo.php?script=sciarttext&pid=S0103-90162006000300008&lng=en>
510 [&nrm=iso](#). Acesso em: 20 fev. 2011. doi: 10.1590/S0103-90162006000300008.
- 511 GRUNWALD S. Multi-criteria characterization of recent digital soil mapping and modeling
512 approaches. **Geoderma**, v.152, p.195-207. 2009. Disponível em:
513 <http://www.sciencedirect.com/science/article/pii/S0016706109001827> Acesso em: 12 out.
514 2011. doi: 10.1016/j.geoderma.2009.06.003.
- 515 GRUNWALD, S. et al. Are Current Scientific Visualization and Virtual Reality Techniques
516 Capable to Represent Real Soil-Landscapes? In: LAGACHERIE, P.; MCBRATNEY, A.B.;
517 VOLTZ, M., Digital Soil Mapping: An Introductory Perspective. Developments in Soil
518 Science, v.31, cap.42, Elsevier, Amsterdam, 2007, p.571-580.
- 519 LAGACHERIE, P. & MCBRATNEY, A. B. Spatial soil information systems and spatial soil
520 inference systems: perspectives for digital soil mapping. In: LAGACHERIE, P.;
521 MCBRATNEY, A. & VOLTZ, M. ed. **Digital soil mapping: an introductory perspective**.
522 Amsterdam: Elsevier, 2007, p. 3-22.
- 523 KEMPEN, B.; et al. Updating the 1:50,000 Dutch soil map using legacy soil data: a
524 multinomial logistic regression approach. **Geoderma**, v.151, p.311-326, 2009. Disponível
525 em: <http://www.sciencedirect.com/science/article/pii/S0016706109001475> Acesso em: 20
526 fev. 2011. doi: 10.1016/j.geoderma.2009.04.023.

- 527 KHEIR, R.B. et al. Predictive mapping of soil organic carbon in wet cultivated lands using
528 classification-tree based models: The case study of Denmark. **Journal of Environmental**
529 **Management**. v.91, n.5, p. 1150-1160. 2010a. Disponível em:
530 <http://www.sciencedirect.com/science/article/pii/S0301479710000022> Acesso em: 12 out.
531 2011. doi: 10.1016/j.jenvman.2010.01.001.
- 532 KHEIR, R.B. et al. Spatial soil zinc content distribution from terrain parameters: A GIS-based
533 decision-tree model in Lebanon. **Environmental Pollution**, v.158, n.2, p.520-528, 2010b.
534 Disponível em: <http://www.sciencedirect.com/science/article/pii/S0269749109004163>
535 Acesso em: 12 out. 2011. doi: 10.1016/j.envpol.2009.08.009.
- 536 KER, J.C. ; NOVAIS, R.F. Fundamentos da pedologia e relação com a fertilidade do solo. In:
537 XXIX Congresso Brasileiro de Ciência do Solo, 2003, Ribeirão Preto, 2003.
- 538 KROL, B.G.C.M. Towards a Data Quality Management Framework for Digital Soil Mapping
539 with Limited Data. In: Digital Soil Mapping with Limited Data. HARTEMINK, A.E.;
540 MCBRATNEY, A.B.; MENDONÇA-SANTOS, M. de L. (Eds.), cap.11, Springer, 2008,
541 p.137-149.
- 542 MCBRATNEY, A. B.; MENDONCA SANTOS, M. L. & MINASNY, B. On digital soil
543 mapping. **Geoderma**, v.17, p.3-52, 2003. Disponível em:
544 <http://www.sciencedirect.com/science/article/pii/S0016706103002234>. Acesso em: 12 Out.
545 2011. doi: 10.1016/S0016-7061(03)00223-4.
- 546 MENDONÇA-SANTOS, M.L. & SANTOS, H.G. dos The state of the art of brazilian soil
547 mapping and prospects for digital soil mapping. In: Digital soil mapping: an introductory
548 perspective. LAGACHERIE P., MCBRATNEY A.B. & VOLTZ, M. (ed.), cap.3, p.39-54,
549 2007.
- 550 MINASNY, B. et al. Digital Soil Mapping Technologies for Countries with Sparse Data
551 Infrastructures. In: Digital Soil Mapping With Limited Data. HARTEMINK, A.E.;

- 552 MCBRATNEY, A.B.; MENDONÇA-SANTOS, M. de L. (Eds.), cap.2, Springer. 2008, p.15-
553 30.
- 554 QI, F.; ZHU, A. X. Knowledge discovery from soil maps using inductive learning.
555 **International Journal of Geographical Information Science**. v.17, n.8, p.771-795, 2003.
556 Disponível em: http://solimserver.geography.wisc.edu/pdfs/QiFeng_IJGIS2003.pdf Acesso
557 em: 12 out. 2011. doi: 10.1080/13658810310001596049.
- 558 RAMOS D.P. Desafios da Pedologia Brasileira frente ao novo milênio - (CNPS/EMBRAPA).
559 In: Simpósio do XXIX Congresso Brasileiro de Ciência do Solo, Ribeirão Preto, SP, 2003.
- 560 SACHS, J. et al. Monitoring the world's agriculture. **Nature**, v.466, p.558–560, 2010.
561 Disponível em: [http://www.css.cornell.edu/faculty/lehmann/publ/Nature%20466,%20558-
562 560,%202010%20Sachs.pdf](http://www.css.cornell.edu/faculty/lehmann/publ/Nature%20466,%20558-560,%202010%20Sachs.pdf) Acesso em: 12 out. 2011. doi:10.1038/466558a.
- 563 SANCHEZ, P.A. et al. Digital soil map of the world. **Science**, v.325, p.680-681, 2009.
564 Disponível em: <http://www.sciencemag.org/content/325/5941/680.full.pdf> Acesso em: 12 out.
565 2011. doi: 10.1126/science.1175084.
- 566 SARMENTO, E. C. et al. Sistema de informação geográfica como apoio ao levantamento
567 detalhado de solos do Vale dos Vinhedos. **Revista Brasileira de Ciência do Solo**, v.32,
568 p.2795-2803. 2008. Disponível em [http://www.scielo.br/scielo.php?script=sci_arttext&pid
569 =S0100-06832008000700025&lng=en&nrm=iso](http://www.scielo.br/scielo.php?script=sci_arttext&pid=S0100-06832008000700025&lng=en&nrm=iso). Acesso em: 23 ago. 2011. doi:
570 [http://dx.org/ 10.1590/S0100-06832008000700025](http://dx.org/10.1590/S0100-06832008000700025).
- 571 SCHMIDT, K. et al. Instance selection and classification tree analysis for large spatial
572 datasets in digital soil mapping. **Geoderma**, v.146, p.138-146, 2008. Disponível em:
573 <http://www.sciencedirect.com/science/article/pii/S0016706108001304> Acesso em: 12 out.
574 2011. doi: 10.1016/j.geoderma.2008.05.010.

- 575 SCULL, P. et al. Predictive soil mapping: a review. **Progress in Physical Geography**, v.27,
576 p.171-197, 2003. Disponível em: <http://ppg.sagepub.com/content/27/2/171>. Acesso em: 14
577 nov. 2011. doi: 10.1191/0309133303pp366ra.
- 578 SCULL, P.; et al. The application of classification tree analysis to soil type prediction in a
579 desert landscape. **Ecological Modelling**, v.181, p.1-15, 2005. Disponível em:
580 <http://www.sciencedirect.com/science/article/pii/S0304380004003540>. Acesso em: 14 nov.
581 2011. doi: 10.1016/j.ecolmodel.2004.06.036.
- 582 TEN CATEN, A. et al. Regressões logísticas múltiplas: fatores que influenciam sua aplicação
583 na predição de classes de solos. **Revista Brasileira de Ciência do Solo**, v.35, n.1, p.53-62,
584 2011a. Disponível em: [http://www.scielo.br/scielo.php?script=sci_arttext&pid=S0100-
585 06832011000100005&lng=en&nrm=iso](http://www.scielo.br/scielo.php?script=sci_arttext&pid=S0100-06832011000100005&lng=en&nrm=iso). Acesso em: 16 mai. 2011. doi: 10.1590/S0100-
586 06832011000100005.
- 587 TEN CATEN, A.; et al. Extrapolação das relações solo-paisagem a partir de uma área de
588 referência. **Ciência Rural**, v.41, n.5, p. 812-816, 2011b. Disponível em:
589 [http://www.scielo.br/scielo.php?script=sci_arttext&pid=S0103-84782011000500012&lng=
590 en&nrm=iso](http://www.scielo.br/scielo.php?script=sci_arttext&pid=S0103-84782011000500012&lng=en&nrm=iso). Acesso em: 12 out. 2011. doi: 10.1590/S0103-84782011000500012
- 591 TEN CATEN et al. Componentes principais como preditores no mapeamento digital de
592 classes de solos. **Ciência Rural**, v.41, n.7, p. 1170-1176, 2011c. Disponível em:
593 [http://www.scielo.br/scielo.php?script=sci_arttext&pid=S0103-84782011000700011&lng=
594 en&nrm=iso](http://www.scielo.br/scielo.php?script=sci_arttext&pid=S0103-84782011000700011&lng=en&nrm=iso). Acesso em: 12 out. 2011. doi: 10.1590/S0103-84782011000700011.
- 595 VALERIANO, M. M. ; ROSSETTI, D. F. Topodata: seleção de coeficientes geoestatísticos
596 para refinamento unificado de dados SRTM. São José dos Campos, SP: NPE: Coordenação de
597 Ensino, Documentação e Programas Especiais (INPE-16701-RPQ/853) (Boletim). 2010. 74p.
- 598 WEBSTER, R.; BURROUGH, P.A. Computer-based soil mapping of small areas from
599 sample data. 1. Multivariate classification and ordination. **Journal of Soil Science**, v.23,

600 p.210-221, 1972. Disponível em: <http://dx.doi.org/10.1111/j.1365-2389.1972.tb01654.x>
601 Acesso em: 12 out. 2011. doi.: 10.1111/j.1365-2389.1972.tb01654.x.
602 WITTEN, I.H. & FRANK, E. Data Mining: Practical Machine Learning Tools and
603 Techniques, 2nd ed., Morgan Kaufmann, 2005, 560p.
604 WOOD, J. Overview of Software Packages Used in Geomorphometry. Cap.10. In.:
605 Geomorphometry: Concepts, Software, Applications. HENGL, T. & REUTER, H.I. (ed).
606 Elsevier, 2009, 76.

ARTIGO 2 - Resolução espacial de um modelo digital de elevação definida pela função wavelet.⁽¹⁾

1 **Resumo** – O objetivo do presente estudo foi empregar a transformada wavelet para avaliar a
2 variabilidade de quatro atributos de terreno, possibilitando definir a resolução espacial mais
3 apropriada para representar a variabilidade desses. Os dados utilizados no estudo têm sua
4 origem em três transetos de 27 km posicionados em áreas do Planalto, Rebordo do Planalto e
5 Depressão Central, na região central do Estado do Rio Grande do Sul. Foram utilizados os
6 atributos elevação, declividade, curvatura em perfil e índice de umidade topográfica. Tais
7 atributos foram derivados de um modelo digital de elevação Topodata com resolução de 30
8 m. A avaliação da resolução em que ocorria a máxima variabilidade foi realizada pela
9 aplicação da wavelet mãe denominada Morlet. Os resultados foram avaliados a partir do
10 isograma e do escalograma dos coeficientes wavelets e indicam que sensores remotos com
11 resolução espacial próxima a 60 e 100 m podem ser utilizados em pesquisas que considerem
12 os atributos de terreno declividade, curvatura em perfil e índice de umidade topográfica, ou
13 ainda, fenômenos ambientais correlacionados a esses. No entanto, não foi possível estabelecer
14 um valor conclusivo para a resolução espacial mais adequada para o atributo elevação.
15 Termos para indexação: mapeamento digital de solos, wavelet mãe de Morlet, escala,
16 pedometria.

17

18

Spatial resolution of a digital elevation model defined by the wavelet function

19 **Abstract** – The aim of this study was to use wavelet transform to assess the variability of four
20 terrain attributes, allowing the appropriate definition of spatial resolution to represent the
21 variability of each attribute. The data used in the study had their origin in three 27 km
22 transects positioned in areas of the Planalto, Rebordo do Planalto and Depressão Central in
23 the central region of the Rio Grande do Sul State. The attributes elevation, slope, profile

24 curvature and topographic wetness index were derived from a Topodata digital elevation
25 model with resolution of 30 m. The assessment of the resolution with the highest spatial
26 variability was evaluated through the application of the Morlet mother wavelet. The results
27 were evaluated from wavelet coefficients plotted in isograms and scalograms. The study
28 results indicate that remote sensors with spatial resolution between 60 and 100 m should be
29 used in researches considering the terrain attributes slope, profile curvature and topographic
30 wetness index, or even environmental phenomena related to them. The study was unable to
31 establish a conclusive value for a more appropriate spatial resolution to the attribute elevation.
32 Index terms: digital soil mapping, Morlet mother wavelet, scale, pedometrics.

33

34 **INTRODUÇÃO**

35 Variáveis ambientais normalmente apresentam uma substancial variabilidade em uma
36 ampla gama de escalas espaciais. As propriedades do solo estão entre as variáveis que, em
37 geral, exibem um complexo padrão espacial fruto das interações entre processos físicos,
38 químicos e biológicos, os quais ocorrem em uma variedade de escalas espaciais. A elucidação
39 das escalas em que esses processos se desenvolvem é essencial para o entendimento e a
40 predição de processos hidrológicos, biológicos e químicos, que atuam diretamente na gênese
41 do solo (Si, 2008).

42 A resolução espacial está intimamente relacionada à escala (Bian, 1997; Hengl, 2006).
43 Entre as dúvidas dos pesquisadores está a que diz respeito à resolução espacial em que os
44 modelos baseados em parâmetros de terreno irão desempenhar seu máximo. Com a
45 diminuição da resolução espacial, os parâmetros primários de terreno derivados do modelo
46 digital de elevação (MDE) perdem detalhamento, e muitas nuances da paisagem também são
47 perdidas (Wu et al., 2008a).

48 Segundo Wu et al. (2008b), diferentes resoluções espaciais do MDE produzem distintas
49 informações sobre os atributos de terreno. À medida que o MDE é agregado em resoluções

50 cada vez mais grosseiras, a declividade diminui, e a área de contribuição aumenta. De acordo
51 com os autores, existe um valor ótimo para a resolução espacial, a partir do qual o modelo
52 hidrológico irá gerar os melhores resultados. Embora a resolução espacial dependa da
53 complexidade da paisagem e dos atributos de terreno utilizados na modelagem, Hengl (2006)
54 argumenta que muitos trabalhos, os quais utilizam dados em formato matricial (raster), optam
55 por uma determinada resolução espacial sem qualquer justificativa científica em favor da
56 opção adotada.

57 A estratégia automatizada de mapeamento do solo, denominada Mapeamento Digital de
58 Solos (MDS), emprega preditores derivados de MDE para a predição espacial de classes e
59 propriedades do solo (Sanchez et al., 2009). Um conjunto de preditores de larga aplicação no
60 MDS são os atributos de terreno (Scull et al., 2003). Esses podem ser derivados diretamente
61 de MDE, como é o caso da declividade, direção do fluxo e curvaturas. Outros atributos são de
62 definição um pouco mais complexa e constituem-se em atributos mais elaborados resultados
63 da análise geomorfométricas da paisagem. Esse é o caso do índice de posição topográfica
64 (IPT), que é utilizado para classificar a paisagem em classes morfológicas (Tagil & Jenness,
65 2008).

66 A metodologia utilizada influencia sobremaneira na qualidade dos valores gerados quando
67 da derivação do IPT (Judex et al., 2006). O IPT representa a diferença entre a elevação de um
68 pixel e a elevação média dos pixels vizinhos, possibilitando uma classificação simples e
69 rápida da paisagem em classes morfológicas. Como esse índice está relacionado à extensão e
70 à forma com que o algoritmo analisa pixels vizinhos, os parâmetros adotados durante a
71 geração do índice devem ser produto de uma criteriosa análise prévia das características
72 geomorfológicas da região para a qual os dados estão sendo gerados. Com isso, será possível
73 avaliar a extensão e os padrões da variabilidade espacial, os quais têm reflexo direto na
74 resolução espacial de trabalho (Tagil & Jenness, 2008).

75 Entre as alternativas de estudo das resoluções espaciais mais apropriadas para os diversos
76 fenômenos ambientais, tem-se a aplicação da análise wavelet (Cho & Chon, 2006; Dong et
77 al., 2008). A aplicação da análise wavelet consiste em mimetizar os dados originais a partir de
78 um sinal artificial produzido por uma pequena onda, a wavelet mãe. A wavelet é dilatada e
79 transladada a fim de buscar a melhor convolução entre os dados originais e a wavelet mãe.
80 Desse procedimento são gerados os coeficientes wavelet, os quais representam o grau de
81 correlação entre o sinal original e o sinal gerado. Esses coeficientes podem, então, ser
82 visualizados na forma do isograma local e escalograma global, que irão auxiliar na avaliação
83 das estruturas predominantes nos dados em questão (Torrence & Compo, 1998).

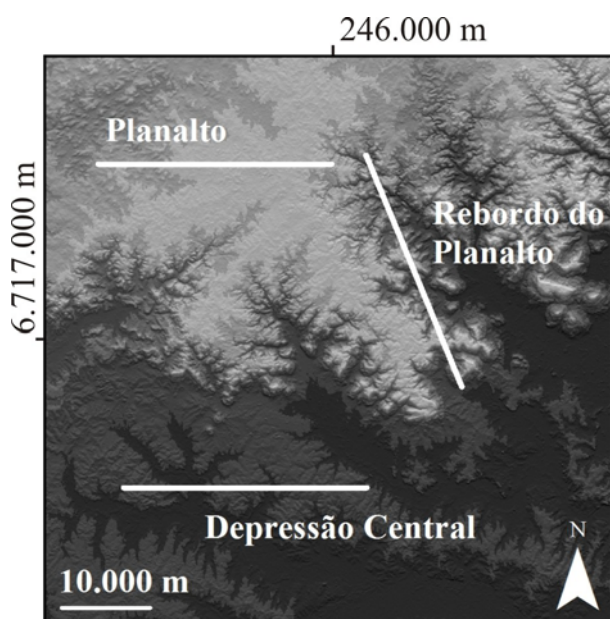
84 A wavelet é uma função matemática representada na forma de uma pequena onda. É
85 considerada pequena, porque contrasta com as funções seno e cosseno típicas da análise de
86 Fourier que se estendem para o infinito. A análise wavelet é executada por meio da
87 transformada wavelet, a qual decompõe a variância de um processo físico em uma série de
88 coeficientes que representam a distribuição da variância em diferentes resoluções espaciais e
89 posições da paisagem (Yates et al., 2006). A análise wavelet constitui-se em uma boa
90 alternativa para a análise de dados assimétricos e não estacionários comuns em dados
91 ambientais (Lark & Webster, 1999).

92 Este estudo teve como objetivo empregar a transformada wavelet para avaliar a
93 variabilidade dos atributos de terreno elevação, declividade, curvatura em perfil e índice de
94 umidade topográfica em três regiões fisiográficas distintas do Estado do Rio Grande do Sul. O
95 estudo buscou definir a resolução espacial mais apropriada para representar a variabilidade de
96 cada um dos quatro atributos nessas regiões, além de demonstrar a aplicação da análise
97 wavelet, que se constitui como uma ferramenta com enorme potencial de aplicação para
98 estudos espaço-temporais.

99 MATERIAL E MÉTODOS

100 Área de estudo

101 Os dados utilizados no estudo têm sua origem em três transetos posicionados em distintas
 102 regiões fisiográficas típicas da região central do Estado do Rio Grande do Sul (Figura 1).
 103 Essas três regiões são o Planalto da Serra Geral, o Rebordo do Planalto e a Depressão Central.
 104 A complexidade geomorfológica dessa região do Estado implica em um desafio para o
 105 mapeamento preditivo de classes e propriedades do solo devido à diversidade de nuances do
 106 relevo que, por sua vez, influencia sobremaneira a formação dos solos da região. Cada
 107 transeto tem um comprimento de aproximadamente 27 quilômetros, o que possibilitou, na
 108 atual resolução espacial dos atributos de terreno, um total de 900 pixels amostrados por
 109 transeto, bem superior às 128 amostras preconizadas como valor mínimo para o emprego da
 110 análise wavelet (Si, 2008).



111

112 **Figura 1.** Localização dos transetos em três regiões fisiográficas predominantes da região
 113 central do Estado do Rio Grande do Sul. Níveis de cinza claros representam cotas acima de
 114 400 metros, níveis escuros identificam locais com cota em torno de 60 m. Coordenadas UTM
 115 referenciadas ao datum SIRGAS 2000/Fuso22.

116

117 **Atributos de terreno**

118 Optou-se por utilizar atributos de terreno com significado físico e com relação à formação
 119 dos solos, visando a sua aplicação posterior em estudos relacionados ao MDS. A elevação, ou
 120 altitude acima de uma superfície de referência, tem importante papel na definição do clima
 121 local. A declividade, como medida da taxa de mudança da elevação na direção do declive,
 122 afeta a velocidade de fluxos superficiais e subsuperficiais. A curvatura em perfil, que
 123 representa a taxa de mudança da declividade na direção de uma linha de fluxo, é fundamental
 124 no comportamento da água sobre a paisagem. O índice de umidade topográfica, que é uma
 125 função do logaritmo natural da razão entre área de contribuição e a declividade, indica o
 126 controle da topografia sobre a umidade do solo (Wilson & Gallant, 2000).

127 Os quatro atributos foram derivados do MDE e extraídos do Topodata (Valeriano &
 128 Rossetti, 2010), que, por sua vez, consiste em uma interpolação dos dados SRTM (*Shuttle*
 129 *Radar Topography Mission*). Esse MDE apresenta resolução espacial relativamente elevada
 130 (30 m) e é obtido por uma única técnica, além de estar disponível de maneira contínua para
 131 extensas regiões. Os arquivos raster de atributos de terreno foram gerados a partir do
 132 programa TAPES-G de acordo com procedimento detalhado por Wilson & Gallant (2000).

133 Cada um dos atributos de terreno foi amostrado a partir do respectivo transeto pela função
 134 “criar gráfico de perfil” do programa ArcMap 9.3. Essa função permite a exportação direta de
 135 todos os valores dos pixels encobertos pela linha do transeto imediatamente para o formato
 136 ASCII.

137 **Análise wavelet**

138 A análise foi executada por meio da transformada wavelet, como segue:

$$139 \quad W(b, x_n) = \frac{1}{\sqrt{b}} \sum_{j=1}^n f(x_j) g\left(\frac{x_j - x_n}{b}\right)$$

140 onde: $W(b, x_n)$ refere-se aos coeficientes wavelet, positivos ou negativos de acordo com a
 141 correlação entre a wavelet e o sinal original; j representa a j -ésima observação; n é o total de

142 observações; b representa a resolução que está sendo testada centrada em x_n ; x_j é a posição
 143 onde a análise está sendo executada no transeto. A função $f(x_j)$ representa os dados originais,
 144 e a função $g(x_j)$ é a wavelet mãe. Neste estudo, foi utilizada a wavelet mãe denominada
 145 Morlet, como segue:

$$146 \quad g(x_j) = \pi^{-\left(\frac{1}{4}\right)} e^{i6\eta} e^{-\left(\frac{1}{2}\right)\eta^2}$$

147 em que

$$148 \quad \eta = \frac{b}{x_j} \text{ (Torrence \& Compo, 1998).}$$

149 Durante a transformada wavelet, valores elevados de variância indicam que houve um bom
 150 ajuste entre a wavelet e os dados originais. Esse processo é repetido para uma faixa de
 151 resoluções ao longo de todo o transeto, dando origem ao isograma local (*wavelet power*
 152 *spectrum*), em que os coeficientes wavelet são plotados ao quadrado, resultando apenas
 153 valores positivos (Torrence & Compo, 1998).

154 Os resultados da análise wavelet são altamente correlacionados com a escolha da wavelet
 155 mãe mais apropriada. Tal ponto merece uma especial atenção (Labat, 2005). A wavelet do
 156 tipo Morlet foi utilizada neste estudo devido à sua superior capacidade de detectar e localizar
 157 padrões nos dados (Mi et al., 2005).

158 A variância total $v(b)$ foi calculada por:

$$159 \quad v(b) = \frac{1}{n} \sum_{i=1}^n W^2(b, x_n)$$

160 A variância total é importante na medida em que seu valor será máximo quando ocorrer a
 161 melhor convolução entre as funções $f(x_j)$ e $g(x_j)$, ou seja, a variância será maximizada quando
 162 a resolução for igual ao tamanho médio das estruturas, como fragmentos e interrupções
 163 presentes nos dados (Si, 2008). O resultado da variância total será representado no
 164 escalograma global (*global wavelet power*). O espectro de ruído vermelho (Torrence &

165 Compo, 1998) ao nível de significância de 5% foi utilizado como hipótese nula para o teste de
166 significância das informações presentes no isograma e no escalograma. Dessa forma, regiões
167 do isograma local circundadas por isolinhas denotam regiões de elevada correlação entre o
168 sinal original e a wavelet mãe, o que indicará padrões relevantes nos dados. No caso do
169 escalograma, locais onde a variância total é superior à linha de significância, as resoluções
170 espaciais são relevantes para o presente estudo.

171 Sobre os isogramas dos coeficientes wavelets, foi representado o cone de influência. Os
172 cones de influência delimitam a presença de efeitos de bordas nos dados próximos às
173 extremidades de cada transeto. Informações nessa região dos isogramas não devem ser
174 consideradas para efeito de afirmativas sobre o padrão de resoluções presentes nos dados
175 originais (Torrence & Compo, 1998). A análise wavelet dos dados foi executada de acordo
176 com o algoritmo descrito e disponibilizado por Torrence & Compo (1998).

177 Gráficos de variância (isogramas e escalogramas) foram utilizados para avaliar os padrões
178 de resolução para cada atributo em cada transeto. Partindo-se do pressuposto de que as áreas
179 do Planalto e da Depressão Central (Figura 1) constituem-se em áreas morfologicamente mais
180 similares, os gráficos de variância foram agrupados em quatro figuras, considerando-se um
181 atributo em duas localidades distintas, a saber: elevação no Planalto e no Rebordo, curvatura
182 em perfil no Planalto e no Rebordo, declividade na Depressão e no Rebordo e índice de
183 umidade topográfica na Depressão e no Rebordo. Dessa forma, são executadas comparações
184 em um mesmo atributo entre duas áreas de maior contrastante, além de possibilitar a análise
185 dos quatro atributos na área do rebordo, onde há a maior diversidade de formas na paisagem.

186 Para definir a resolução espacial (tamanho de pixel) que melhor representa a variabilidade
187 espacial existente nos quatro atributos, foi adotado o critério de He et al. (2007). Para esses
188 autores, a resolução espacial do pixel deve ser de um quarto (1/4) do valor identificado pela

189 análise wavelet. Segundo os mesmos autores, com esse critério, há uma maior segurança para
190 evitar que padrões importantes associados com os atributos não sejam perdidos.

191

192 **RESULTADOS E DISCUSSÃO**

193 **Atributos de terreno**

194 A análise preliminar dos dados advindos da amostragem nos três transetos permite
195 visualizar características típicas das regiões fisiográficas onde os dados foram coletados
196 (Tabela 1). O menor (115,329 m) e o maior (491,970 m) valor para o atributo elevação, no
197 transeto localizado na área do Rebordo do Planalto, são reflexo da conformação do relevo
198 nessa área. Localizada na transição entre o Planalto Meridional Brasileiro e a Depressão
199 Periférica, essa região tem como características a presença de paisagem fortemente dissecada,
200 formando vales em forma de “V”. Essas formas são resultado dos processos erosivos que ao
201 longo de milhares de anos vêm esculpindo a região, gerando esse padrão de grande
202 variabilidade na geomorfologia (Uberti & Klamt, 1984). O desvio padrão desse atributo para
203 o Planalto ($\sigma = 32,329$ m) e para a Depressão Central ($\sigma = 20,352$ m) indica que essas áreas
204 são mais planas do que se comparadas a área do Rebordo do Planalto.

205 Quanto à declividade, o valor máximo (41,963 °), o valor médio (13,720 °) e o desvio
206 padrão (9,830 °) desse atributo em áreas do Rebordo indicam relevo acidentado. Comparando
207 os valores da declividade no Rebordo com os obtidos no Planalto e Depressão, verifica-se que
208 nessas áreas o relevo é mais suave. Segundo Sartori (2009), na Depressão Central
209 predominam as colinas e planícies aluviais formadas a partir de rochas sedimentares. Nas
210 áreas do Planalto, o relevo varia de ondulado a suave ondulado formado sobre uma sucessão
211 de derrames vulcânicos sobrepostos.

212 A curvatura em perfil indica a taxa de mudança do relevo ao longo do declive, valores
213 positivos indicam superfícies convexas com possível perda de material, que, por sua vez,
214 acumula-se nos locais côncavos com curvatura negativa (Wilson & Gallant, 2000). As áreas

215 de Rebordo possuem as maiores curvaturas entre os três transetos (Tabela 1), o que era
 216 esperado devido à complexidade do terreno. A suavidade do relevo na Depressão e Planalto é
 217 retratada pelo pequeno desvio padrão do atributo curvatura em perfil.

218 **Tabela 1.** Estatística descritiva dos quatro parâmetros de terreno nos três transetos utilizados
 219 no estudo.

Parâmetro de terreno	Transeto	Mínimo	Máximo	Média	Mediana	Desvio Padrão	Variância
Elevação (m)	Rebordo	115,329	491,970	272,617	240,220	106,159	11269,818
	Planalto	383,259	513,487	462,471	475,019	32,329	1045,153
	Depressão	57,715	140,803	97,904	100,173	20,352	414,187
Curvatura em perfil (m ⁻¹)	Rebordo	-1,241	0,996	0,015	0,028	0,221	0,049
	Planalto	-0,296	0,253	0,005	0,007	0,063	0,004
	Depressão	-0,182	0,450	0,004	-0,001	0,062	0,004
Declividade (°)	Rebordo	0,240	41,963	13,720	10,854	9,830	96,639
	Planalto	0,264	17,088	3,382	2,946	2,047	4,188
	Depressão	0,342	9,547	3,430	3,203	1,946	3,786
Índice de umidade topográfica	Rebordo	4,657	16,306	7,458	6,869	1,985	3,942
	Planalto	5,734	18,725	8,553	8,193	1,524	2,322
	Depressão	6,119	19,531	8,467	7,987	1,569	2,462

220

221 O índice de umidade topográfica, por se tratar de um atributo que considera a área de
 222 contribuição da bacia e a declividade do terreno, possibilita maior semelhança entre os valores
 223 do Planalto, Rebordo e Depressão (Tabela 1). Contudo, as maiores amplitudes entre áreas
 224 úmidas (valores elevados) e áreas drenadas (valores baixos) foram verificadas na Depressão.
 225 Para essa região, drenam todos os córregos que se originam nas encostas do Rebordo do
 226 Planalto (Sartori, 2009), possibilitando valores elevados de área de contribuição, associadas às
 227 menores declividades dessa região, o que acarretou em maiores valores do índice de umidade
 228 topográfica.

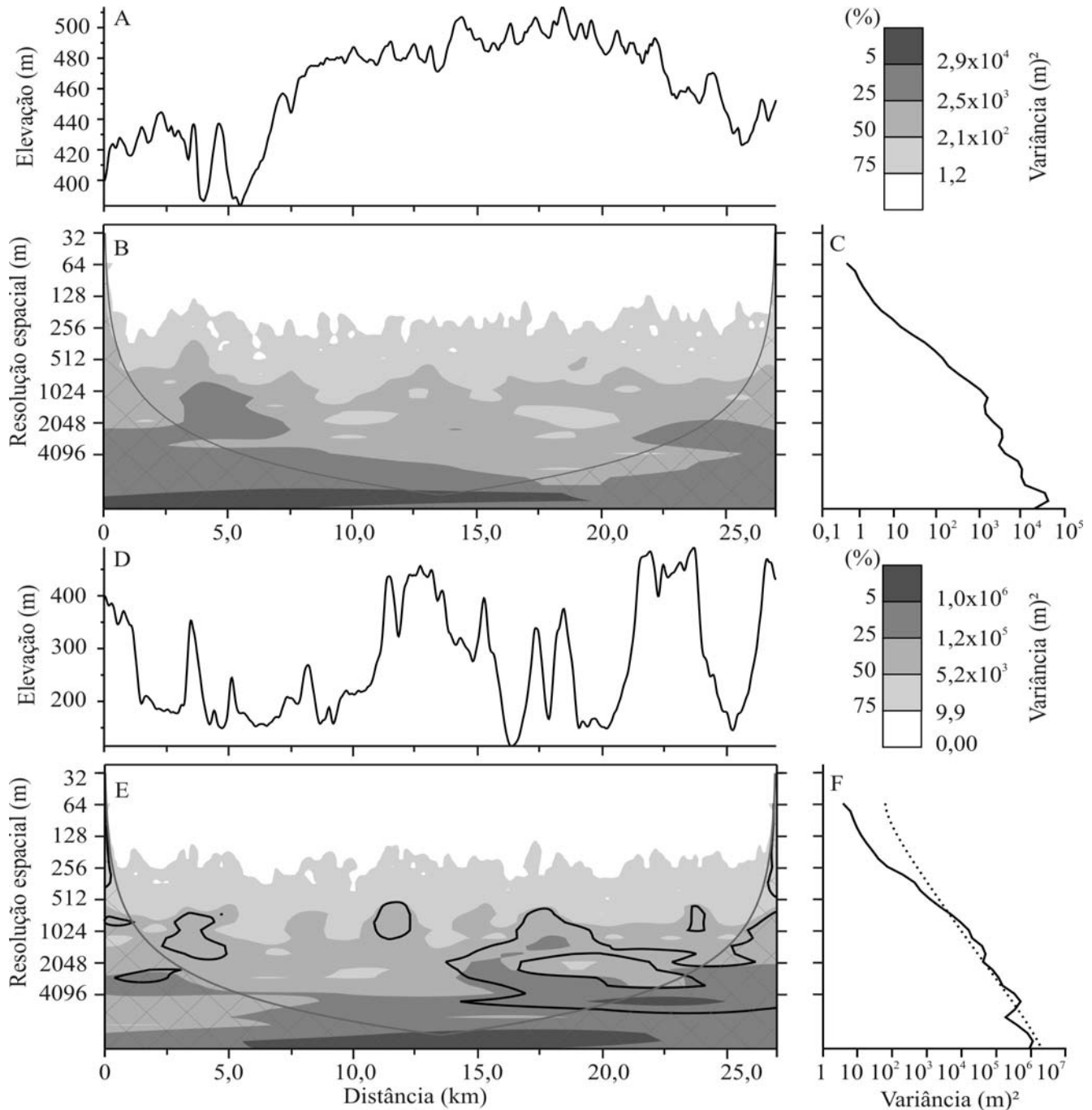
229 Transformada wavelet

230 Na região do Planalto, o atributo elevação possui um comportamento que, conforme
 231 aumenta a resolução espacial, também aumenta a variância, e a distribuição dos tons de cinza
 232 ocorre de maneira quase que horizontal ao longo de todo o transeto (Figura 2 B). Há uma

233 maior variabilidade entre 4 e 5 km do início do transeto, para valores de resolução de 900 a
234 2100 m, indicando a presença de pequenos talwegues e mudanças abruptas da elevação
235 (Figura 2 A). Esse é o padrão típico de um evento local, não havendo um ciclo repetitivo
236 desse padrão no restante do isograma e, portanto, não se pode atribuir uma resolução a este
237 atributo (Yates et al., 2006). O escalograma da Figura 2 C indica uma variância crescente com
238 o aumento da resolução, embora os valores não sejam estatisticamente significativos.

239 Três são os locais de variância significativa para os dados de elevação no transeto
240 localizado no Rebordo (Figura 2 E). O primeiro localizado entre 2,5 e 5 km do início do
241 transeto na faixa de resoluções de 700 a 1500 m; o segundo na região de 11,5 km para
242 resoluções de 512 a 1024 m; e o terceiro iniciando em 13,5 km, estendendo-se até o final do
243 transeto, o que mostra um grande padrão de variabilidade. Contudo, nesse último, existem
244 duas resoluções espaciais predominantes, uma em torno de 1024 m e outra localizada próxima
245 aos 4096 m. Todos esses eventos estão relacionados à grande variabilidade do atributo
246 elevação demonstrado na Figura 2 D. Nas distâncias em que ocorreram os maiores valores de
247 variância significativa, também estão localizadas mudanças abruptas no relevo.

248 As informações do isograma da Figura 2 E não revelam a presença de um padrão cíclico na
249 variabilidade. Nesse tipo de gráfico, resoluções associadas a eventos cíclicos ao longo de todo
250 o transeto são o foco deste estudo. Entretanto, a análise de significância para a variância
251 global (Figura 2 F) demonstra que existem resoluções associadas a padrões de variabilidade
252 entre 550 e 4200 m. Ressalta-se que essas informações devam ser tomadas com cautela, uma
253 vez que não ocorreram eventos cíclicos no isograma, e os dados de variância global sejam
254 apenas levemente superiores aos do espectro de ruído vermelho. Também é importante frisar
255 que áreas contidas sob o cone de influência não devam ser consideradas devido a efeitos de
256 bordas nos dados (Torrence & Compo, 1998). Observa-se que uma grande área de variância
257 significativa encontra-se na região na Figura 2 E.



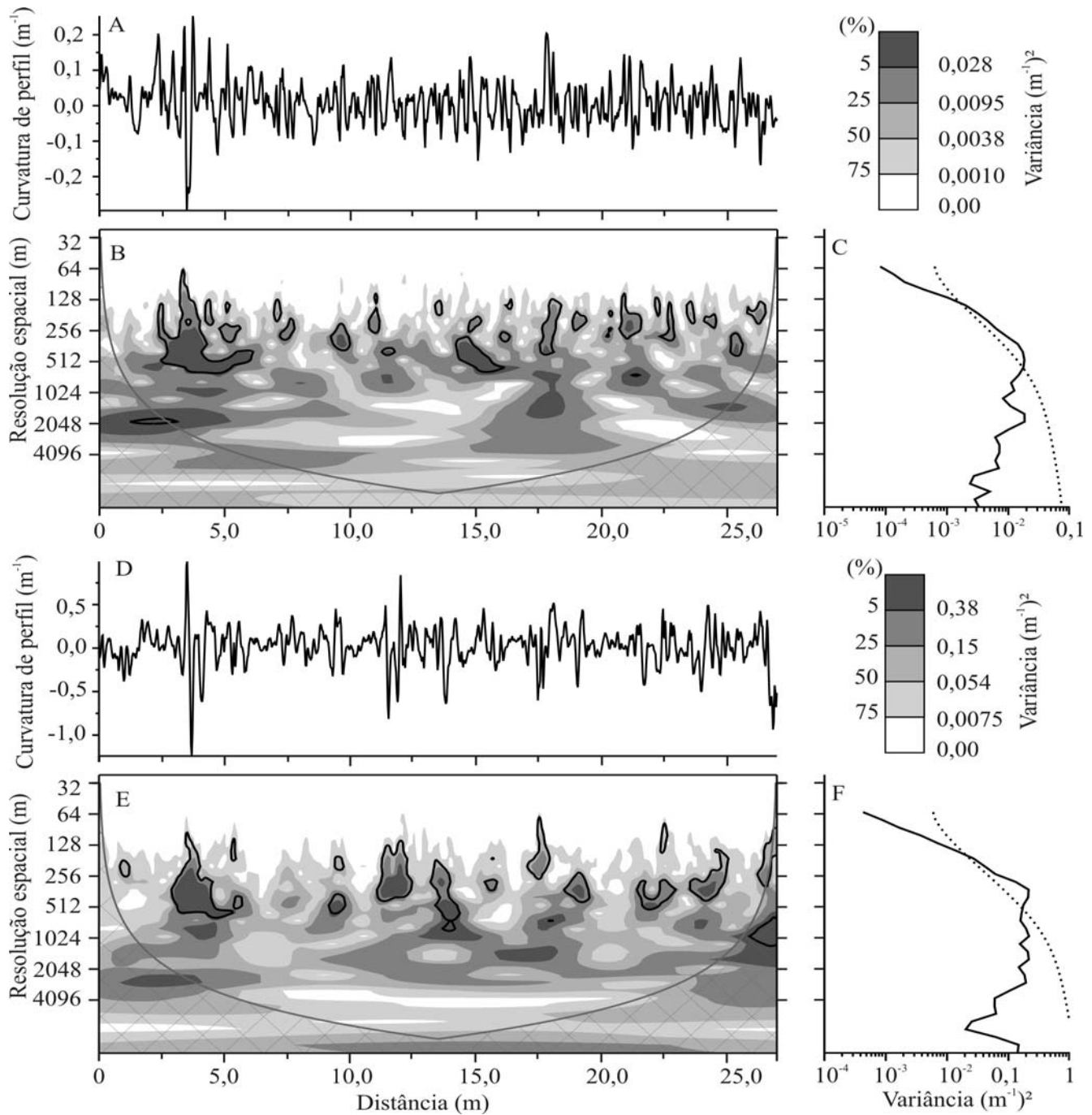
258
259

260 **Figura 2.** Elevação no Planalto (A) e Rebordo (D). Isogramas de variância wavelet no
 261 Planalto (B) e Rebordo (E). Regiões quadriculadas representam o cone de influência.
 262 Escalogramas da variância total no Planalto (C) e Rebordo (F). Análise de significância foi
 263 realizada com o ruído vermelho como espectro de fundo ao nível de 5% de significância nos
 264 isogramas (isolinhas) e escalogramas (linha tracejada). A legenda dos níveis de cinza indica o
 265 valor e a quantidade de variância relativa em cada nível acima no isograma.

266 A distribuição da variância do atributo curvatura em perfil apresentou um padrão cíclico ao
267 longo dos 27 km do transeto para a região do Planalto (Figura 3 B). Entre 2,5 e 5 km existe
268 uma grande variabilidade nos valores desse atributo (Figura 3 A), o que foi revelado pela
269 transformada wavelet. Para a sequência do transeto, o que se verifica é que as regiões de
270 variabilidade significativa se repetem praticamente em intervalos regulares. Grande parte
271 desse padrão está situado entre as resoluções de 128 e 512 m, que corresponde à variância
272 total significativa na Figura 3 C. Segundo Yates et al. (2006), nos casos em que há uma região
273 cíclica e contígua no padrão da variabilidade, a opção pela resolução espacial deve ser feita ao
274 centro desse padrão. Adotando esse critério, uma resolução importante para descrever o
275 padrão espacial da curvatura em perfil na região do Planalto seria em torno de 256 m.

276 O padrão da variância da transformada wavelet para a curvatura em perfil no transeto
277 localizado no Rebordo também apresentou ciclos de ocorrência repetitivos (Figura 3 E). Para
278 esse atributo é possível perceber que as variações de curvatura envolveram amplitudes
279 maiores, localizadas em 3; 12; 14; 17,5; 22,5; 24 e 27 km (Figura 3 D), do que aquelas
280 percebidas na área do Planalto (Figura 3 A), o que provocou um aumento no valor da
281 resolução espacial (variância total com pico em torno de 380 m – Figura 3 F), confirmando o
282 aumento da resolução para a curvatura em perfil em áreas do Rebordo.

283 Duas regiões de resoluções espaciais significativas podem ser verificadas na transformada
284 wavelet dos dados de declividade do transeto localizado na Depressão (Figura 4 B). Uma
285 primeira, com um padrão mais cíclico, situada entre 70 e 512 m, e outra, essa não cíclica,
286 entre as resoluções de 512 e 1100 m, aproximadamente. A observação do espectrograma da
287 variância total (Figura 4 C) indica que a primeira região tem um pico de variância situado por
288 volta de 256 m, o que indica ser essa uma resolução importante para a variabilidade da
289 declividade nessa área. Em estudo realizado por Si & Farrell (2004), foi verificado que o
290 atributo comprimento de rampa possuía uma resolução significativa em 180 m.



291

292

Figura 3. Curvatura em perfil no Planalto (A) e Rebordo (D). Isogramas de variância wavelet

293

no Planalto (B) e Rebordo (E). Regiões quadriculadas representam o cone de influência.

294

Escalogramas da variância total no Planalto (C) e Rebordo (F). Análise de significância foi

295

realizada com o ruído vermelho como espectro de fundo ao nível de 5% de significância nos

296

isogramas (isolinhas) e escalogramas (linha tracejada). A legenda dos níveis de cinza indica o

297

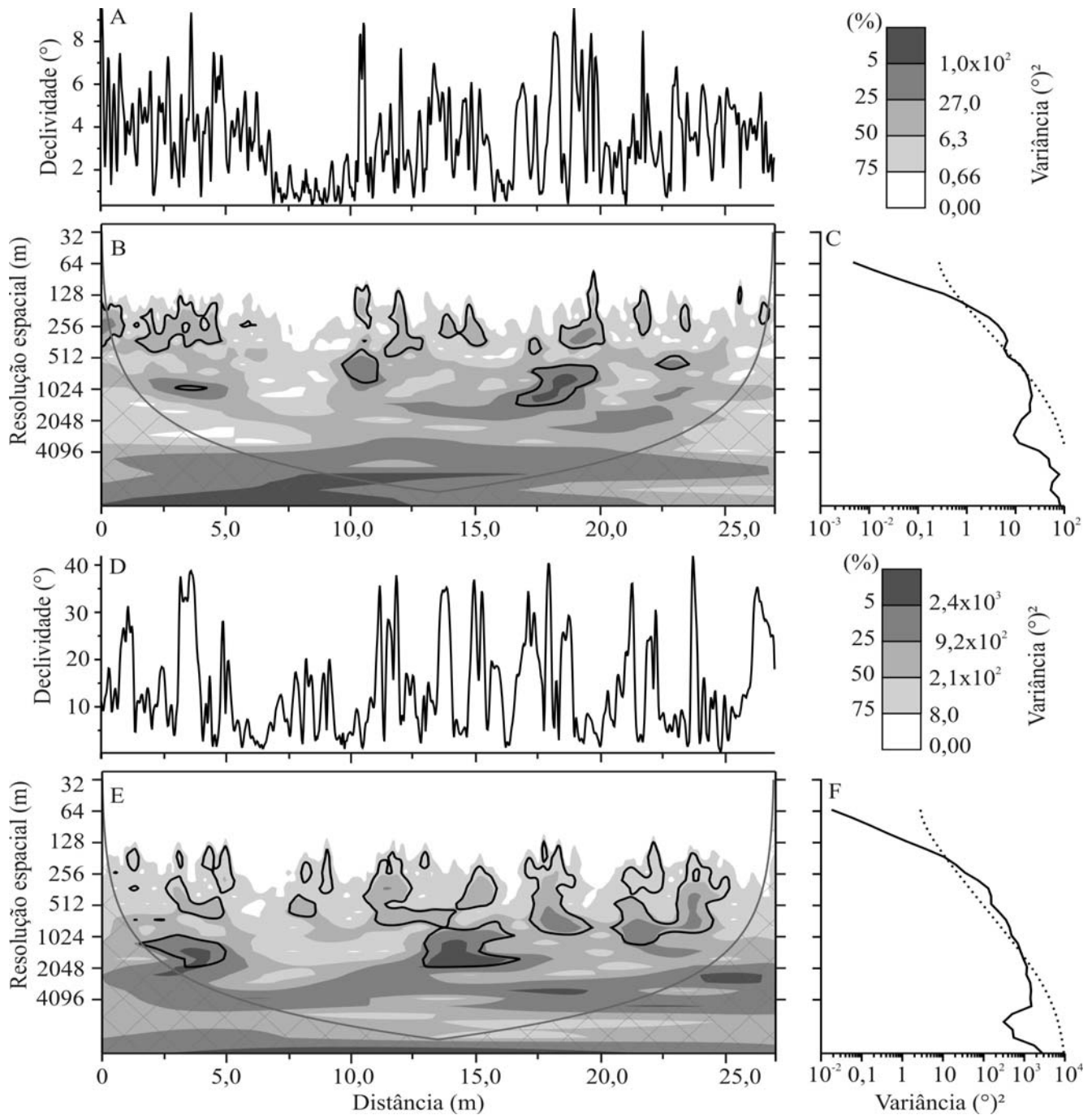
valor e a quantidade de variância relativa em cada nível acima no isograma.

298 A avaliação do somatório de todas as áreas tidas como significativas no isograma da
299 Figura 4 E indica ser esse o maior valor entre os gráficos analisados. Tal fato decorre das
300 grandes amplitudes de declividade presentes sistematicamente ao longo de todo o transeto do
301 Rebordo do Planalto. Há um padrão cíclico entre as resoluções de 128 a 1024 e outro entre as
302 resoluções de 1024 a 2048 m, embora esse último não se repita sistematicamente ao longo dos
303 27 km. O gráfico da variância total indica que há um pico na região situada em torno da
304 resolução de 380 m (Figura 4 F). Esse mesmo valor foi relatado como significativo para a
305 determinação do padrão espacial do atributo curvatura em perfil na região do Rebordo (Figura
306 3 E).

307 Como a região do Rebordo do Planalto serve como área de contribuição da drenagem para
308 a Depressão Central, os valores do índice de umidade topográfica foram maiores em áreas do
309 transeto localizado na Depressão do que aquele transeto localizado em área de Rebordo
310 (Figura 5 A e 5 D). Ambos os isogramas apresentam um padrão cíclico da variância
311 significativa ao longo dos 27 km (Figura 5 B e 5 E). A observação dos espectrogramas para
312 essas regiões indica que há uma faixa de resoluções, entre 200 e 800 m, responsável por
313 grande parte da variabilidade nesses transetos (Figura 5 C e 5 F). Se utilizarmos o critério de
314 Yates et al. (2006) e adotarmos a resolução espacial central entre esses dois valores como
315 sendo aquela determinante para explicar o padrão de distribuição espacial do índice de
316 umidade topográfica, esse valor seria de 400 m.

317 A análise feita nos isogramas indicou que três são as resoluções importantes, a saber: para
318 os atributos declividade e curvatura em perfil, tem-se o valor de 256 m em áreas da Depressão
319 e Planalto e de 380 m para áreas do Rebordo; o atributo índice de umidade topográfica teria
320 seu padrão de variabilidade espacial melhor representado em resoluções de 400 m nas áreas
321 da Depressão e no Rebordo. O critério do valor de um quarto dessas resoluções (He et al.,
322 2007) remete a tamanhos de pixels de 64, 95 e 100 m, respectivamente. Quanto ao atributo

323 elevação não foi possível precisar com segurança a resolução espacial que melhor representa
 324 sua variabilidade, contudo, os resultados apontam que valores maiores do que 550 m
 325 poderiam ser utilizados.



326

327 **Figura 4.** Declividade na Depressão (A) e Rebordo (D). Isogramas de variância wavelet na

328 Depressão (B) e Rebordo (E). Regiões quadriculadas representam o cone de influência.

329 Escalogramas da variância total na Depressão (C) e Rebordo (F). Análise de significância foi

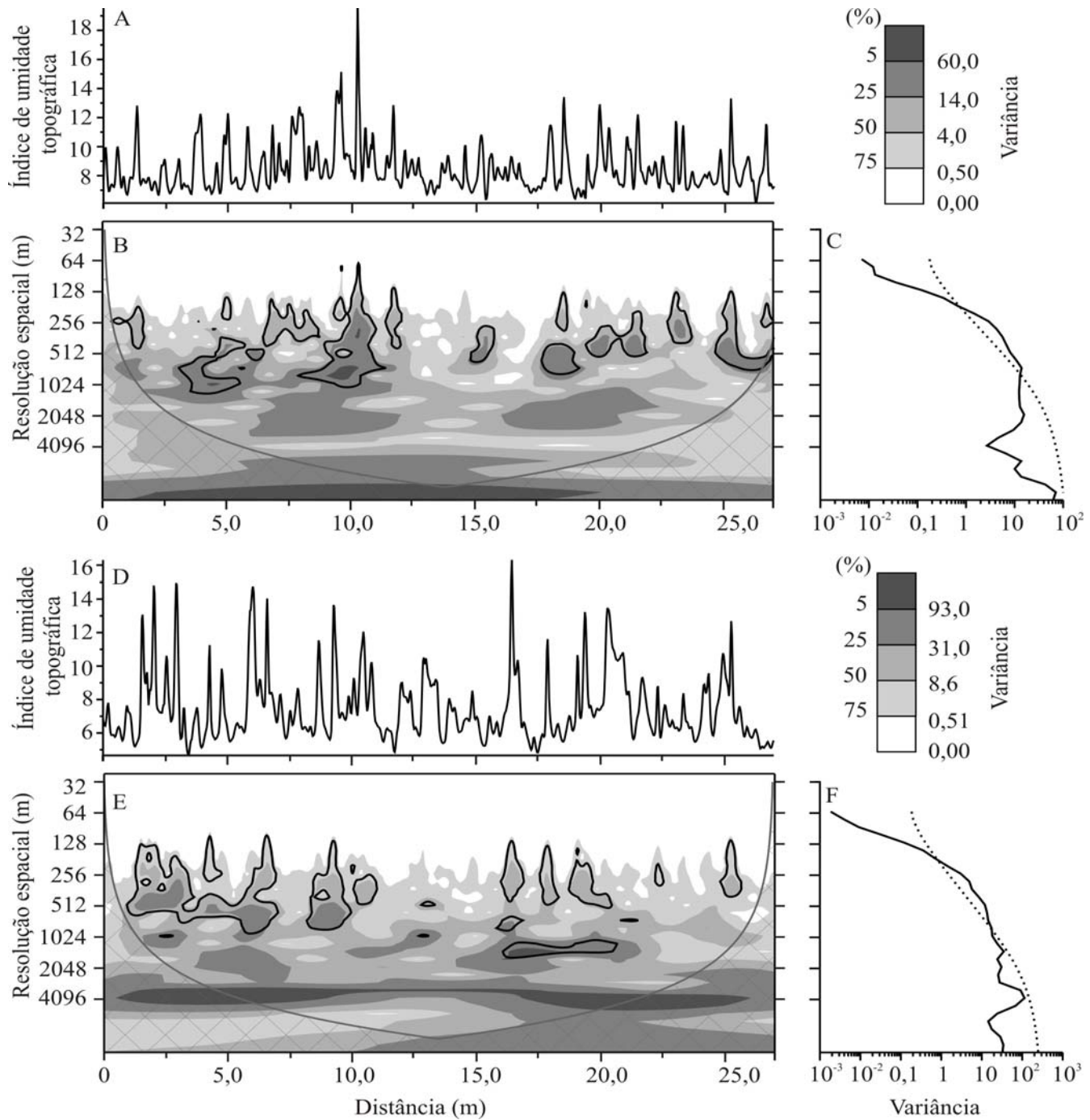
330 realizada com o ruído vermelho como espectro de fundo ao nível de 5% de significância nos
331 isogramas (isolinhas) e escalogramas (linha tracejada). A legenda dos níveis de cinza indica o
332 valor e a quantidade de variância relativa em cada nível acima no isograma.

333

334 A definição da resolução espacial e do tamanho de pixel desses atributos de terreno pode
335 ainda ser utilizada para definir a densidade de amostragem de outras covariáveis ambientais.
336 Si & Farrell (2004) identificaram um padrão de variabilidade espacial similar entre o índice
337 de umidade topográfica e o rendimento do trigo em áreas de estudo do Canadá para
338 resoluções menores do que 140 m. Para He et al. (2007), uma resolução espacial apropriada,
339 determinada pela transformada wavelet para o estudo de padrões associados a atributos de
340 terreno e índices de área foliar, seria de 20 m de resolução.

341 Os resultados do nosso estudo indicam que sensores remotos com tamanhos de pixel
342 próximos a 60 e 100 m podem ser utilizados em pesquisas que considerem os atributos de
343 terreno declividade, curvatura em perfil e índice de umidade topográfica, ou ainda fenômenos
344 ambientais correlacionados a esses.

345 O emprego do MDE Topodata (30 m) (Valeriano & Rossetti, 2010) para MDS, conforme
346 informações obtidas neste trabalho, pode representar um grau de detalhamento acima do
347 necessário, o que deverá ser testado em estudos futuros. A aplicação da informação original,
348 oriunda do MDE SRTM com resolução de 90 m, como satisfatória para MDS também deve
349 ser avaliada, assim como a aplicação da transformada wavelet para definição da melhor
350 resolução espacial das plataformas Landsat e Cbers aplicadas ao MDS.



351

352 **Figura 5.** Índice de umidade topográfica na Depressão (A) e Rebordo (D). Isogramas de
 353 variância wavelet na Depressão (B) e Rebordo (E). Regiões quadriculadas representam o cone
 354 de influência. Escalogramas da variância total na Depressão (C) e Rebordo (F). Análise de
 355 significância foi realizada com o ruído vermelho como espectro de fundo ao nível de 5% de
 356 significância nos isogramas (isolinhas) e escalogramas (linha tracejada). A legenda dos níveis
 357 de cinza indica o valor e a quantidade de variância relativa em cada nível acima no isograma.

CONCLUSÕES

1. Os atributos declividade e curvatura em perfil oriundos de relevos mais suaves têm seu padrão de variabilidade associados à resolução de 256 m. Em relevos mais acidentados, a resolução apropriada é de 380 m. Tamanhos dos pixels em torno de 64 e 95 m são indicados para utilização desses atributos na região de estudo.

2. A resolução de 400 m representou o padrão de variabilidade do índice de umidade topográfica e foi independente da conformação do terreno. O tamanho de pixel de 100 m seria adequado para análises em torno desse atributo.

3. Não foi possível estabelecer para o atributo elevação um valor conclusivo e adequado de resolução espacial e tamanho de pixel.

4. A utilização da wavelet do tipo Morlet permitiu visualizar padrões de variabilidade local e global em atributos de terreno.

AGRADECIMENTOS

Ao CNPq pela bolsa de produtividade em pesquisa do segundo autor e aporte financeiro para auxílio deste trabalho.

REFERÊNCIAS

BIAN, L. Multiscale nature of spatial data in scaling up environmental models. In: QUATTROCHI, D.A; GOODCHILD, M.F. (Ed.) **Scale in remote sensing and GIS**. CRC Press, 1997. p. 13-26.

CHO, E.; CHON, T. Application of wavelet analysis to ecological data. **Ecological Informatics**, v.1, n.3, p.229-233, 2006.

DONG, X.; NYREN, P.; PATTON, B.; NYREN, A.; RICHARDSON, J.; MARESCA, T. Wavelets for agriculture and biology: a tutorial with applications and outlook. **BioScience**, v.58, n.5, p.445-453, 2008.

- HE, Y.; GUO, X.; SI, B.C. Detecting grassland spatial variation by a wavelet approach. **International Journal of Remote Sensing**, v.28, n.7, p.1527-1545, 2007.
- HENGL, T. Finding the right pixel size. **Computers & Geosciences**, v.32, p.1283-1298, 2006.
- JUDEX, M.; THAMM, H.; MENZ, G. Improving land-cover classification with a knowledge based approach and ancillary data. In: Workshop of the EARSeL SIG on Land Use and Land Cover. Braun, M. (Ed.), 2nd., 2006, Bonn. **Proceedings**. Bonn: Center for Remote Sensing of Land Surfaces, 2006. p.184-191.
- LABAT, D. Recent advances in wavelet analyses: Part 1. A review of concepts. **Journal of Hydrology**, v.314, p.275-288, 2005.
- LARK, R.M; WEBSTER, R. Analysis and elucidation of soil variation using wavelets. **European Journal of Soil Science**, v. 50, p. 185-206, 1999.
- MI, X.; REN, H.; OUYANG, Z.; WEI, W.; MA, K. The use of the Mexican Hat and the Morlet wavelets for detection of ecological patterns. **Plant Ecology**, v.179, p.1-9, 2005.
- SANCHEZ, P.A.; AHAMED, S.; CARRÉ, F.; HARTEMINK, A.E.; HEMPEL, J.; HUISING, J.; LAGACHERIE, P.; MCBRATNEY, A.B.; MCKENZIE, N.J.; MENDONÇA-SANTOS, M. de L.; MINASNY, B.; MONTANARELLA, L.; OKOTH, P.; PALM, C.A.; SACHS, J.D.; SHEPHERD, K.D.; VÅGEN, T.; VANLAUWE, B.; WALSH, M.G.; WINOWIECKI, L.A.; ZHANG, G. Digital soil map of the world. **Science**, v.325, p.680-681, 2009.
- SARTORI, P.L.P. Geologia e geomorfologia de Santa Maria. **Ciência & Ambiente**, v.38, p.17-42, 2009.
- SCULL, P.; FRANKLIN, J.; CHADWICK, O. A.; MCARTHUR, D. Predictive soil mapping: a review. **Progress in Physical Geography**, v.27, p.171-197, 2003.
- SI, B.C. Spatial scaling analyses of soil physical properties: A review of spectral and wavelet methods. **Vadose Zone Journal**, v.7, n.2, p.547-562, 2008.

SI, B.C.; FARRELL R.E. Scale-dependent relationship between wheat yield and topographic indices: A wavelet approach. **Soil Science Society of America Journal**, 2004, v.68, n.2, p.577-587, 2004.

TAGIL, S.; JENNESS, J. GIS-based automated landform classification and topographic, landcover and geologic attributes of landforms around the Yazoren Polje, Turkey. **Journal of Applied Sciences**, v.8, n.6, p.910-921, 2008.

TORRENCE C.; COMPO, G.P. A practical guide to wavelet analysis. **Bulletin of the American Meteorological Society**, v.79, n.1, p.61-78, 1998.

UBERTI, A.A.; KLAMT, E. Relações solo-superfícies geomórficas na encosta inferior do nordeste do Rio Grande do Sul. **Revista Brasileira de Ciência do Solo**, v.8, p.124-132, 1984.

VALERIANO, M. M. ; ROSSETTI, D. F. Topodata: seleção de coeficientes geoestatísticos para refinamento unificado de dados SRTM. São José dos Campos, SP: NPE: Coordenação de Ensino, Documentação e Programas Especiais (INPE-16701-RPQ/853) (Boletim). 2010. 74p.

WILSON, J.P.; GALLANT, J.C. **Terrain Analysis: Principles and Applications**. New York: John Wiley & Sons, 2000. 479p.

WU, W.; FAN, Y.; WANG, Z.; LIU, H. Assessing effects of digital elevation model resolutions on soil-landscape correlations in a hilly area. **Agriculture, Ecosystems & Environment**, v.126, p.209-216, 2008a.

WU, S.; LIB, J.; HUANG, G. A study on DEM-derived primary topographic attributes for hydrologic applications: Sensitivity to elevation data resolution, **Applied Geography**, v.28, p. 210-223, 2008b.

YATES, T.T.; SI, B.C.; FARRELL, R.E.; PENNOCK, D.J. Wavelet spectra of nitrous oxide emission from hummocky terrain during spring snowmelt. **Soil Science Society American Journal**, v.70, p.1110-1120, 2006.

ARTIGO 3 - Mapeamento digital de solos: estratégia de pré-processamento de dados¹

1

2 **RESUMO:** Mapas de solos têm na borda dos polígonos a região de maior variabilidade, o
3 que leva os pedólogos a divergirem quanto ao delineamento das classes de solos nesses locais.
4 O objetivo deste estudo foi propor uma estratégia de pré-processamento de dados aplicada ao
5 mapeamento digital de solos. Polígonos de solos em um mapa de treinamento foram
6 deslocados para seu interior em 100 e 160 m. Essa estratégia possibilitou que covariáveis
7 localizadas próximas à borda das classes de solos não fossem utilizadas para a geração dos
8 modelos de Árvore Decisão (AD). Três AD, geradas a partir de oito covariáveis preditoras,
9 ligadas aos fatores relevo e organismos, amostradas por um mapa de solos completo e com
10 polígonos deslocados em 100 e 160 m, foram utilizadas para prever classes de solos. O
11 modelo de AD a partir de observações distantes 160 m da borda dos polígonos no mapa
12 original é menos complexo e tem melhor desempenho preditivo.

13 **Termos de indexação:** mapa cloroplético, pedometria, levantamento de solos.

14 **Digital soil mapping: strategy for data preprocessing**

15 **SUMMARY:** Soil maps have at its polygons edge the region of highest variability, which
16 leads pedologists to disagree about the proper delineation of soil classes at those locations.
17 The aim of this study was to propose a preprocessing strategy applied to digital soil mapping.
18 Soil polygons on a training map were displaced in its inward direction by 100 and 160 m.
19 This strategy has enabled that data covariates located near at the border of soil classes were
20 not used for Decision Tree (DT) model adjusting. Three DT models derived from eight
21 predictors covariate, related with factors relief and organisms, sampled by a complete soil
22 map and by polygons displaced 100 and 160 m were used to predict soil classes. The DT
23 model derived from observations distant 160 m of the boundary between polygons in the
24 original map is less complex and has a better predictive performance.

25 **Index terms:** choropleth map, pedometrics, soil survey.

26

INTRODUÇÃO

28 Mapas de solos podem ser preditos a partir de informações pré-existentes de solos. O
29 modelo 'scorpan' (McBratney et al., 2003) contempla informações existentes sobre classes e
30 propriedades do solo que possam auxiliar na predição e no mapeamento digital de solos
31 (MDS) em áreas onde a informação espacial sobre o solo não existe ou está indisponível na

32 escala demandada. De acordo com Qi & Zhu (2003), a ideia básica da formalização das
33 relações solo-paisagem contidas em mapas cloropléticos de solos consiste em se reverter o
34 levantamento de solos pela aplicação de técnicas de mineração de dados. As informações
35 contidas nos polígonos de solos em mapas cloropléticos têm demonstrado sua aplicabilidade
36 ao MDS (Crivellini et al., 2009; Giasson et al., 2011; ten Caten et al., 2011a).

37 Nos mapas de solos, o delineamento das unidades de mapeamento é feito através da
38 interpretação visual de pares estereoscópicos, em escala compatível conforme o objetivo do
39 mapa (Dalmolin et al. 2004). A posição espacial dos polígonos de solos implica relações entre
40 as diferentes classes de solos e as condições ambientais presentes na paisagem. Ao delinear as
41 bordas dos polígonos, o pedólogo é conduzido pelo seu conhecimento tácito das relações
42 entre os múltiplos planos de informação que impactam na gênese do solo local, como a
43 geologia, o relevo e o uso da terra (Qi & Zhu, 2003).

44 Intrínseco ao próprio método está presente a subjetividade. Variações tênues ou graduais
45 nas condições ambientais são de difícil localização através desse método (Zhu et al., 2001), o
46 que gera uma incerteza quanto à real localização na paisagem das transições entre classes de
47 solos. A utilização desses polígonos como referência para o treinamento de modelos
48 preditivos implicará a adição de informações desviadas. Essas, por sua vez, terão diferentes
49 impactos na qualidade preditiva dos modelos (Qi, 2004).

50 Na busca por melhorar a qualidade do banco de dados utilizado para prever classes de
51 solos, Qi & Zhu (2003) utilizaram apenas as informações que se posicionassem próximas da
52 moda de cada covariável preditora. Através da construção de um histograma para cada
53 covariável, os autores descartaram aqueles dados que estavam fora da moda do conjunto de
54 dados. A acurácia dos modelos desenvolvidos por árvore de decisão no conjunto de dados
55 filtrados chegou a alcançar 86%, ao passo que modelos utilizando o conjunto total de dados
56 obtiveram uma média de 75%. Para esses autores, a metodologia de filtragem foi efetiva
57 como estratégia de pré-tratamento de dados aplicados ao MDS. Contudo, Schmidt et al.,
58 (2008) advertem que essa metodologia por histograma tem como desvantagem a necessidade
59 de se construir um histograma para cada covariável preditora, gerando dificuldade em
60 estudos com número elevado de covariáveis e observações.

61 Métodos para a seleção das observações são utilizados quando uma grande quantidade de
62 amostras está disponível. Estratégias apropriadas para a seleção das observações conduzem a
63 melhores resultados do que quando os modelos são ajustados ao total de observações
64 disponíveis (Schmidt et al., 2008). O desafio está em fazer mais com menos (Liu & Motoda,
65 2002), ou seja, extrair um conjunto de dados representativo dos dados originais, utilizando um

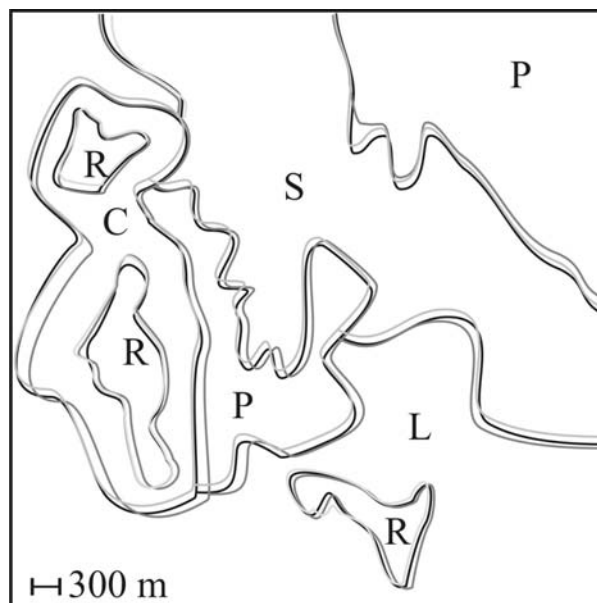
66 conjunto menor de dados, facilmente manipulável pelos algoritmos. Nesse caso, há ganho de
67 acurácia e velocidade de processamento dos dados.

68 Estudos em MDS buscam avaliar a qualidade das covariáveis e das observações, visando
69 potencializar a acurácia das predições (Schmidt et al., 2008). Entre os estudos que consideram
70 as observações utilizadas no MDS estão os que buscam definir a melhor densidade de
71 amostragem (Moran & Bui, 2002; Scull et al., 2005; Grinand et al., 2008; Schmidt et al.,
72 2008), embora sejam raros os estudos que se preocupam com o padrão das observações dentro
73 de cada classe de solos a ser predita (Qi & Zhu, 2003; Qi, 2004).

74 Avaliando a aplicação de regressões logísticas múltiplas para a predição de classes de
75 solos, ten Caten et al. (2011b) observaram a confusão dos modelos entre classes próximas na
76 paisagem. Observou-se que os maiores percentuais de confusão da predição ocorreram entre
77 as quatro distintas subordens de Argissolos, uma vez que essas classes de solo ocupam
78 posições muito semelhantes na paisagem. Para os autores, essa dificuldade por parte dos
79 modelos pode ter origem no próprio delineamento do mapa original que serviu de
80 treinamento, já que o solo não tem uma transição abrupta, como as classes de solo no mapa
81 original (com polígonos cloropléticos) ou, ainda, devido a diferenças muito tênues entre os
82 atributos do terreno (covariáveis ambientais), que podem não apresentar nenhum tipo de
83 gradiente na borda dos polígonos das classes de solo.

84 A influência das áreas de transição entre diferentes classes de solos também foi registrada
85 em estudo realizado por Carvalho et al. (2009). Segundo esses autores, devido à prática da
86 cartografia baseada em polígonos (Booleana) adotada no modelo convencional de mapas de
87 solos, ocorre a tendência, na construção de mapas digitais gerados pela aplicação da
88 metodologia da lógica nebulosa, ao aparecimento de áreas que estariam relacionadas a zonas
89 de transição entre duas ou mais unidades de solo. Isso acarretaria o surgimento de
90 delineamentos inexistentes no mapa convencional utilizado para treinamento dos modelos.

91 A definição da exata posição das bordas dos polígonos é questão controversa entre os
92 pedólogos. Legros (2005) avaliou a qualidade do delineamento das unidades de mapeamento
93 executadas por 20 diferentes pedólogos em uma mesma área. Foi identificado que os
94 polígonos tendem a se sobreporem. Contudo, pequenos desvios no delineamento são
95 ocasionados pela percepção de cada pedólogo das variações das informações contidas nas
96 fotografias aéreas. Esses desvios do delineamento possibilitam definir uma região de incerteza
97 quanto à mais adequada localização para a borda dos polígonos (Figura 1). Nessa região de
98 incerteza, as covariáveis preditoras podem estar apresentando uma contribuição duvidosa à
99 qualidade dos modelos preditivos.



100

101 Figura 1. Classes de solos delineadas por três diferentes pedólogos, indicando a região de
 102 incerteza na definição do posicionamento da borda dos polígonos. Cambissolo (C), Neossolo
 103 (R), Argissolo (P) e Latossolo (L). Adaptado de Legros (2005).

104

105 O objetivo deste estudo foi avaliar o impacto de uma estratégia de pré-processamento de
 106 dados aplicada ao MDS. Foi verificado o efeito da não utilização de informações derivadas de
 107 covariáveis preditoras presentes nas bordas dos polígonos de solos em modelos por árvore de
 108 decisão para a predição de classes de solos. Este estudo propõe uma metodologia de
 109 amostragem que exclui amostras desviadas e seleciona observações que contém as
 110 características das covariáveis preditoras em cada classe de solos.

111

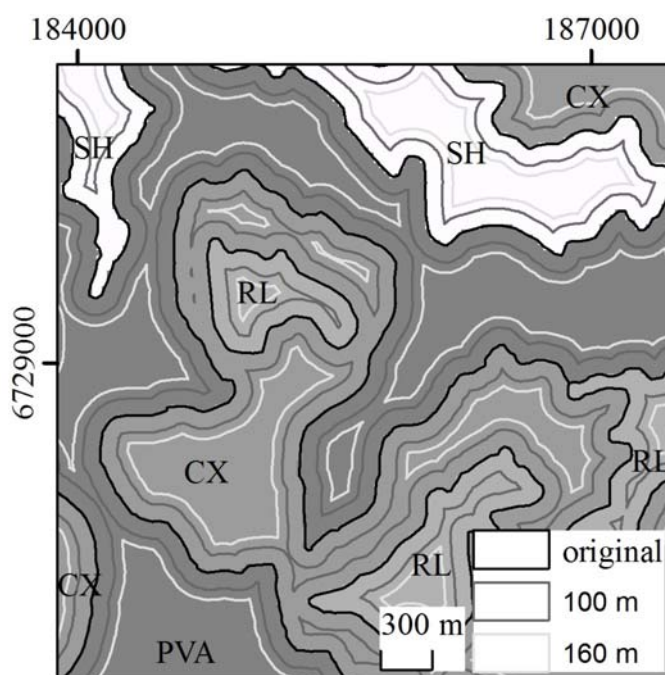
112 MATERIAL E MÉTODOS

113 Mapa de solos para treinamento

114 Para a realização do presente estudo, foi utilizado como teste o levantamento de solos
 115 semidetalhado na escala 1:50.000, realizado por Klamt et al. (2001) para o município de São
 116 Pedro do Sul, localizado na região central do Estado do Rio Grande do Sul, com uma
 117 extensão de 874 km². As classes de solos ao nível de subordem foram vetorizadas no
 118 programa ArcGIS 9.3 (ESRI, 2008).

119 A proposta deste estudo consiste em reposicionar cada classe de solo presente no mapa
 120 original para dentro do polígono originalmente definido pelo pedólogo. Com isso, as regiões
 121 de maior incerteza quanto à real posição das bordas dos polígonos de classes de solos não são
 122 amostradas para a utilização nos modelos preditivos. Para verificar o efeito da metodologia no

123 conjunto de dados a serem gerados, foram utilizadas as distâncias de 100 e 160 m em relação
 124 à posição original de cada classe de solo (Figura 2).



125
 126 Figura 2. Estratégias de amostragem das covariáveis preditoras adotadas no estudo, usando o
 127 mapa de solos original, a 100 m e 160 m da borda dos polígonos. Solos Hidromórficos (SH -
 128 Planossolos e Gleissolos), Cambissolo Háplico (CX), Neossolo Litólico (RL) e Argissolo
 129 Vermelho-Amarelo (PVA).

130
 131 Como todos os polígonos de solos são reposicionados na sua direção interna, a faixa
 132 de dados efetivamente descartada foi de 200 e 320 m nas bordas de polígonos vizinhos. Os
 133 valores de 100 e 160 m foram definidos a partir de uma análise visual dos gradientes de
 134 valores ocorrentes nas covariáveis preditoras próximas à borda dos polígonos. Os novos
 135 polígonos foram criados a partir da função Buffer do programa ArcGIS 9.3. A metodologia de
 136 amostragem proposta não implicou que alguma classe de solo não fosse amostrada em virtude
 137 de representar uma diminuição da área total compreendida por cada classe de solo.

138 **Covariáveis preditoras**

139 Neste estudo foram utilizadas covariáveis do modelo 'scorpan' (McBratney et al., 2003),
 140 relacionadas aos fatores de formação do solo relevo (*r*) e organismos (*o*). O fator organismo
 141 foi representado pela covariável desvio padrão do índice de vegetação por diferença
 142 normalizada (NDVI), doravante denominada DEPA. A covariável DEPA foi usada em
 143 detrimento ao NDVI pelo fato dessa última apresentar diferentes valores ao longo do ano em
 144 função dos diferentes usos agrícolas. Para o cálculo da covariável DEPA, foram utilizadas

145 imagens obtidas no período de fevereiro de 2004 a janeiro de 2005. Esse período foi
146 selecionado devido à maior disponibilidade de imagens com ausência de nuvens durante
147 aquele ano. Todos os cálculos de NDVI foram realizados conforme Jensen (2009). O preditor
148 DEPA foi gerado a partir de dados de oito distintas datas daquele período, obtidas pela
149 plataforma Landsat 5 sensor TM com resolução espacial de 30 m. Cada data teve seu valor de
150 NDVI calculado individualmente, em seguida foi executado o cálculo do desvio padrão do
151 valor do NDVI entre as oito datas para cada pixel no programa ArcGIS 9.3, na função *Raster*
152 *Calculator*.

153 O fator relevo foi representado pelas covariáveis elevação (ELEV), declividade (DECL),
154 índice de umidade topográfica (IUT), capacidade de transporte de sedimento (CTS), curvatura
155 planar (PLAN), curvatura de perfil (PERF) e índice de rugosidade do terreno (IRT). As
156 covariáveis foram geradas de acordo com Wilson & Gallant (2000) a partir de um modelo
157 digital de elevação (MDE). O MDE com resolução espacial de 30 m foi derivado de curvas de
158 nível de cartas topográficas com escala 1:50.000. A interpolação das curvas de nível para o
159 formato matriz ocorreu na ferramenta *topo do raster* do programa ArcGIS 9.3, usando a
160 técnica *spline* (Wahba, 1990). No programa SAGA-GIS (Olaya, 2004), os atributos ELEV,
161 DECL, PLAN, PERF, e IRT foram gerados na ferramenta *standard terrain analysis*, e os
162 atributos IUT e CTS com a ferramenta *grid calculator*.

163 Esses atributos foram tabulados a partir dos polígonos de solos, gerando três distintos
164 conjuntos de dados para construção dos modelos por árvore de decisão: original, 100 m e 160
165 m.

166 **Árvore de Decisão (AD)**

167 Nos três conjuntos de dados utilizados para desenvolver os modelos, todas as classes de
168 solos foram amostradas proporcionalmente a sua área. O desenvolvimento das AD foi
169 realizado no programa de mineração de dados WEKA 3.6.3 (Hall et al., 2009). Para o
170 processamento dos dados, foi utilizado o algoritmo J48 que apresentou os melhores resultados
171 em estudo realizado por Giasson et al. (2011). O número mínimo de observações por folha
172 (*minNumObj* - WEKA), para cada conjunto de dados foi determinado após uma análise do
173 percentual de observações erroneamente classificadas. Essa análise foi executada com uma
174 série de valores para o número mínimo de observações em cada um dos três conjuntos de
175 dados. O método de poda que mitigava o erro na árvore gerada também foi selecionado
176 (*reducedErrorPruning = True* - WEKA). Durante a fase de geração da árvore, cada conjunto

177 de dados foi particionado em 70% para geração do modelo e 30% para validação da árvore
178 (*Percentage split* - WEKA).

179 **Mapa de solos**

180 Após a análise da complexidade das AD geradas e do número de observações
181 erroneamente classificadas, as AD foram implementadas no programa ArcGIS 9.3, com a
182 função *Raster Calculator*. As informações derivadas da árvore foram convertidas na função
183 condicional ‘con(teste, verdadeiro, falso)’ do programa. Essa função permite que os arquivos
184 raster com as covariáveis ambientais sejam processados de acordo com o conjunto de regras
185 derivadas da AD. Como a escala de publicação pretendida é de 1:50.000, em cada mapa de
186 solos gerado foram extraídas regiões de píxeis isolados com área mínima mapeável menor do
187 que um hectare.

188 **Qualidade dos modelos e mapas**

189 A qualidade dos modelos AD foi avaliada a partir do valor percentual de observações
190 erroneamente classificadas em toda a árvore. Este valor é uma das saídas do programa WEKA
191 após o modelo gerado ser validado no conjunto de 30% dos dados reservados. Para seu
192 cálculo, o programa soma todas as observações erroneamente classificadas e as divide pelo
193 total de observações do conjunto teste, multiplicando o resultado por 100 para gerar o valor
194 percentual (Hall et al., 2009). Os mapas de solos gerados a partir dos três distintos conjuntos
195 de dados foram comparados com a geração do índice kappa. O índice kappa é um indicador
196 utilizado para atestar a qualidade dos mapeamentos preditivos (Giasson et al., 2011). A matriz
197 dos erros para o cálculo do índice kappa foi executada conforme Congalton (1991).

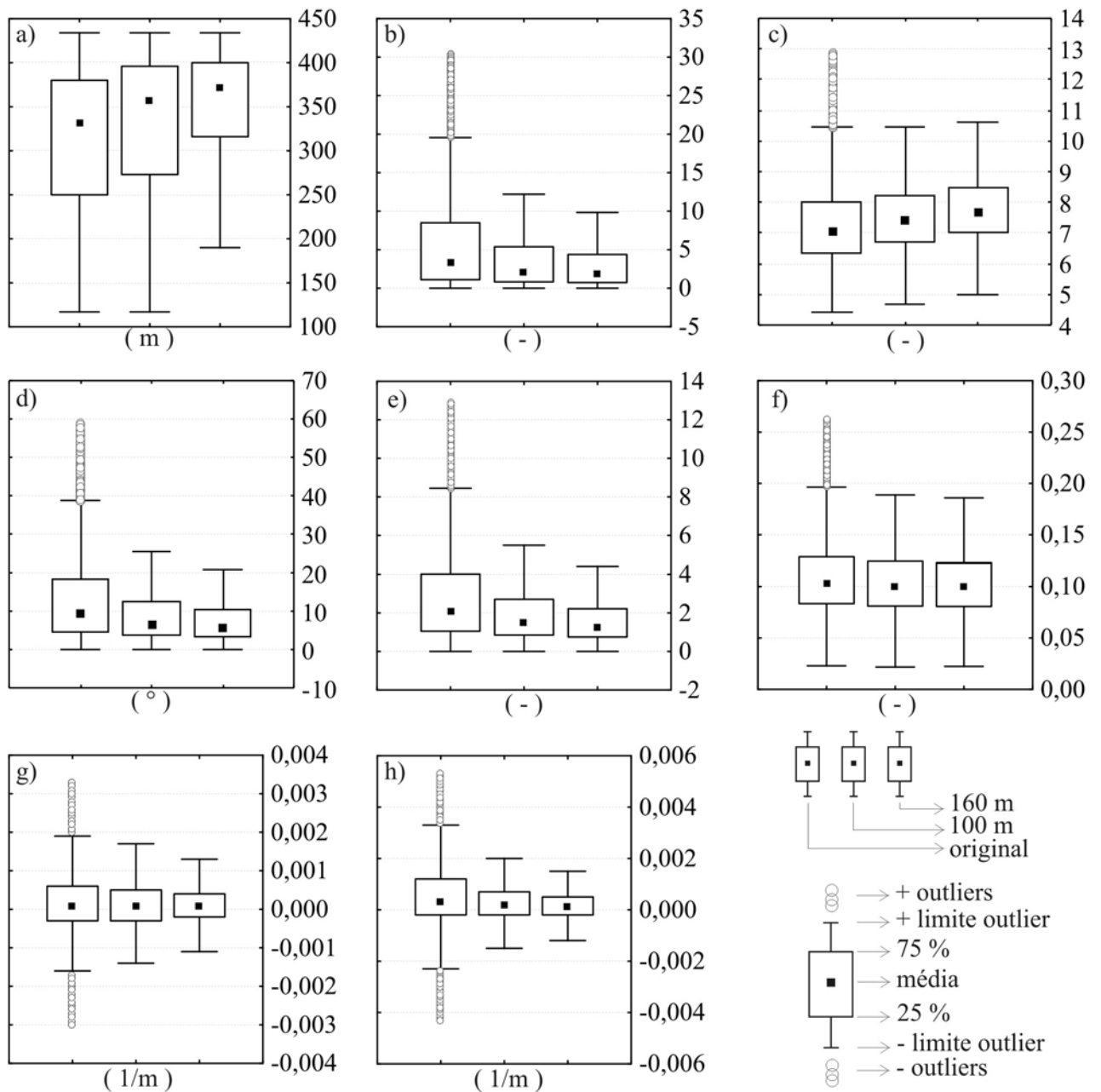
198

199 **RESULTADOS E DISCUSSÃO**

200 A metodologia de pré-processamento proposta implicou a redução do número total de
201 observações disponíveis para gerar os modelos por árvore de decisão. No conjunto de dados
202 derivados do mapa original, estão presentes 100% das informações utilizadas neste estudo. O
203 conjunto de dados gerado a partir do deslocamento dos polígonos em 100 m reteve 60% dos
204 dados originais. Por ocasião de um deslocamento de 160 m nos polígonos de solos, 43% dos
205 dados originais foram tabulados para gerar a árvore de decisão. Esses percentuais de dados
206 são superiores aos 30% utilizados por Grinand et al. (2008) e aos 25% utilizados por Moran &
207 Bui (2002) para o ajuste de árvores de decisão aplicadas ao mapeamento digital de solo.

208 Os conjuntos de dados derivados de polígonos de Neossolos Litólicos deslocados em 100
209 e 160 m têm medidas descritivas distintas do conjunto de dados original (Figura 3). Entre as

210 observações mais distantes das características centrais de cada covariável preditora estão os
 211 outliers. Com exceção do atributo elevação, as demais covariáveis continham outliers entre os
 212 dados amostrados a partir do polígono original. Com a aplicação dos deslocamentos nos
 213 polígonos de solos, observações próximas das bordas dos polígonos foram descartadas, e a
 214 presença de valores distantes em cada covariável não ocorreu. O pré-processamento
 215 significou mudanças similares no padrão dos dados originados a partir das demais classes de
 216 solos, embora esses não estejam aqui representados.



217
 218 Figura 3. Boxplot das covariáveis preditoras gerados a partir dos dados amostrados pela classe
 219 dos Neossolos Litólicos nas três situações de posição das bordas dos polígonos. a) elevação
 220 (m), b) capacidade de transporte de sedimento, c) índice de umidade topográfica, d)

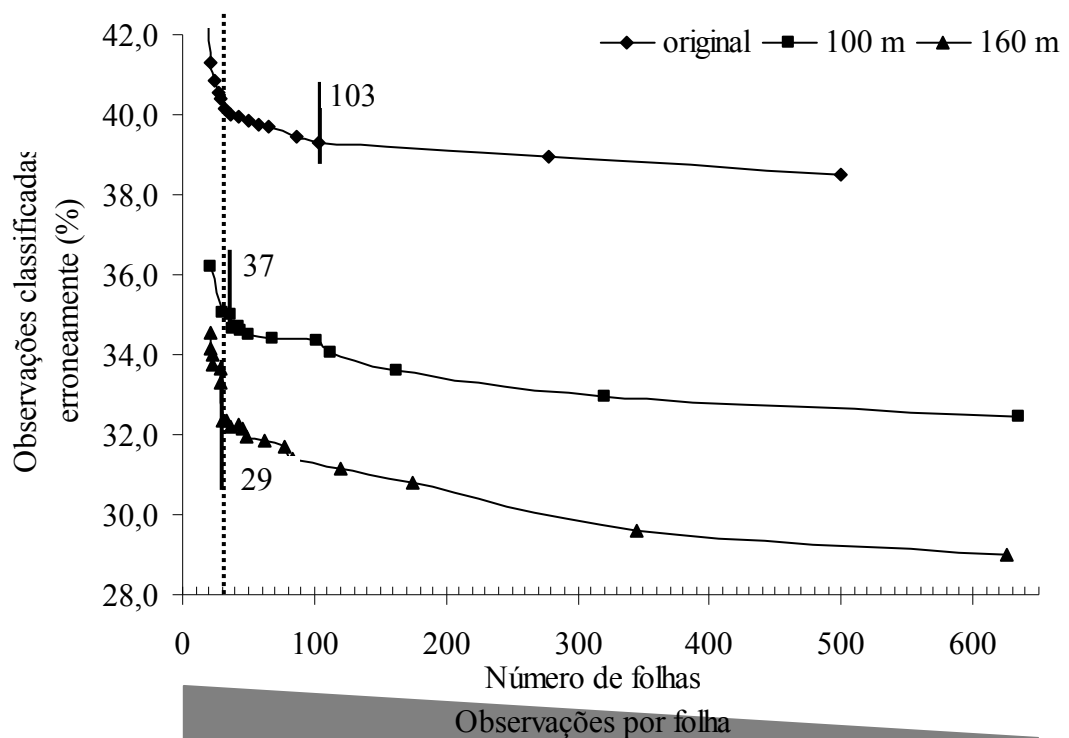
221 declividade ($^{\circ}$), e) índice de rugosidade do terreno, f) desvio padrão do NDVI, g) curvatura
222 planar (1/m) e h) curvatura de perfil (1/m). Unidades de cada covariável estão entre
223 parênteses. O limite outlier foi calculado a partir da altura da caixa central do boxplot
224 multiplicado por 1,5. Foram tomados como outliers os valores além do limite outlier.

225

226 A análise visual dos gráficos indica que o pré-processamento não alterou a distribuição
227 dos dados. A posição da média em relação ao conjunto de 50% do total de dados indica que
228 esses não têm distribuição normal. Esse padrão se repetiu para os dados oriundos dos
229 deslocamentos de 100 e 160 m. Embora tenha ocorrido uma diminuição da dispersão dos
230 dados, esses passaram a se situar menos distantes da média. Esse comportamento é menos
231 marcante nos dados oriundos da covariável desvio padrão do NDVI (Figura 3f). Isso pode ser
232 atribuído ao fato de que na área de estudo os Neossolos Litólicos estejam situados em áreas de
233 encosta, as quais se encontram cobertas por mata nativa. Da posição dos polígonos no mapa
234 original até os polígonos deslocados 160 m não ocorrem mudanças marcantes do uso da terra
235 no local.

236 AD obtidas pelos três conjuntos de dados têm um comportamento distinto quanto ao
237 número de observações erroneamente classificadas (Figura 4). As AD geradas a partir do
238 conjunto total de dados obtiveram uma classificação das observações com um erro mínimo de
239 38,5%. Os demais conjuntos de dados alcançaram erros mínimos de 32 e 28% para
240 deslocamentos de 100 e 160 m, respectivamente. Isso indica que houve um benefício do pré-
241 processamento dos dados para o melhor ajuste do modelo de AD. Contudo, as árvores
242 ajustadas aos dados com menor presença de observações desviadas não classificaram
243 adequadamente cerca de um terço dos dados. Isso pode estar ocorrendo em decorrência da
244 complexidade da distribuição espacial dos solos no local, que não foi representada pelas oito
245 covariáveis preditoras escolhidas para este estudo. Além disso, pode estar ligado a
246 características do mapa de referência como a escala e o delineamento dos polígonos.

247 Nos três conjuntos de dados, o menor percentual de erros foi alcançado com árvores com
248 um maior número de nós e folhas terminais. O número de observações erroneamente
249 classificadas permaneceu praticamente inalterado na medida em que era permitido ao
250 programa agrupar um maior número de observações nas folhas finais. Contudo, próximo ao
251 número de 30 folhas, os três conjuntos de dados apresentaram uma tendência de aumento das
252 observações erroneamente classificadas. Acredita-se que a partir desse limite a árvore seja
253 demasiadamente simplificada e torna-se incapaz de prever a complexidade presente nos
254 dados.



255

256 Figura 4. Relação entre o número de folhas e observações classificadas erroneamente por
 257 árvore de decisão. Conjuntos de dados sem deslocamento (original) e com deslocamento de
 258 100 m e 160 m da borda dos polígonos. Os números próximos às barras verticais indicam o
 259 número mínimo de folhas para que todas as classes de solos fossem preditas com aquele
 260 conjunto de dados. A linha tracejada vertical indica o número de 30 folhas como sendo aquele
 261 a partir do qual ocorre um súbito incremento do número de observações erroneamente
 262 classificadas nos três conjuntos de dados.

263

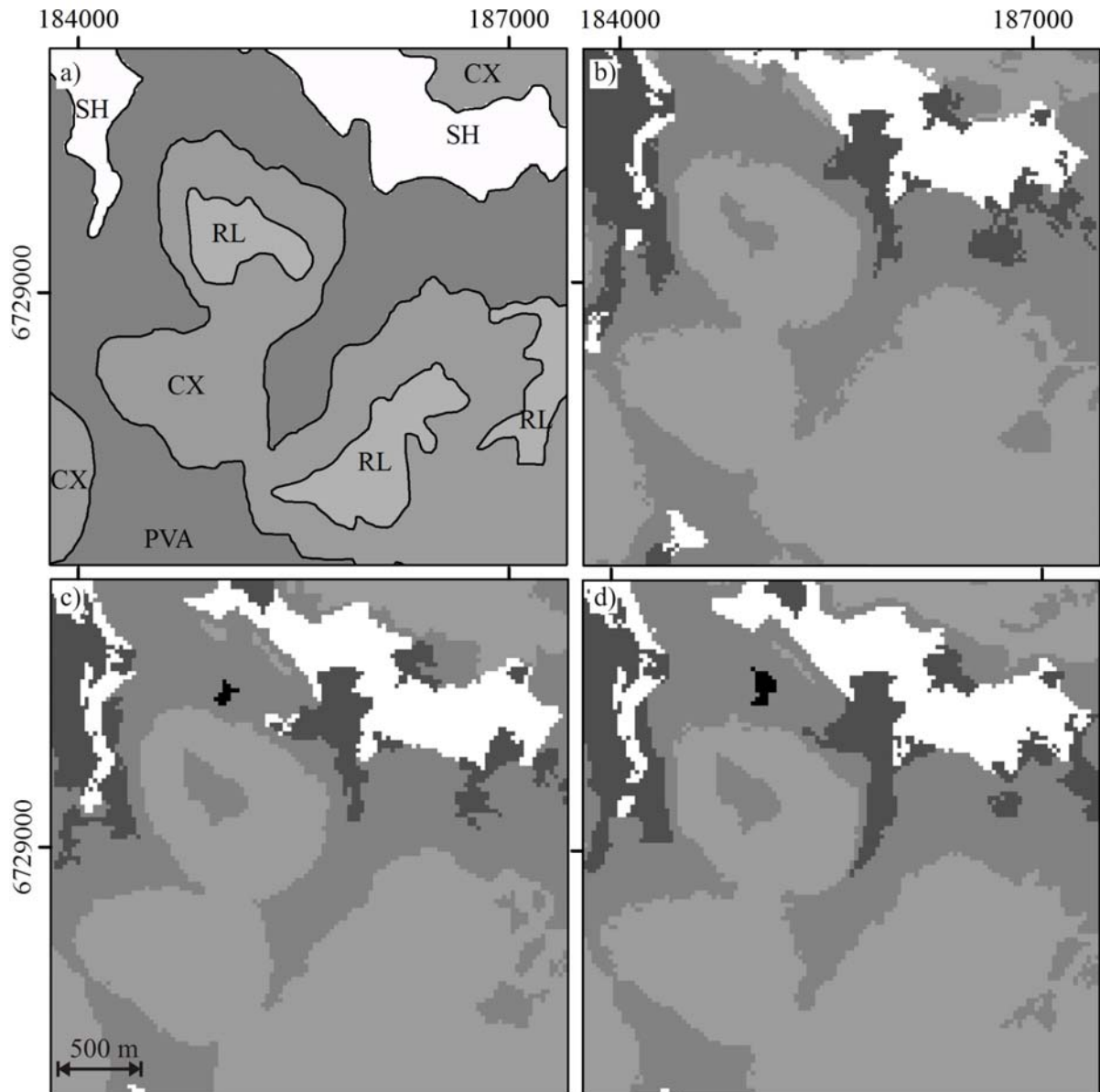
264 AD geradas a partir dos dados afastados das bordas dos polígonos de solos têm mais
 265 poder preditivo. Observando o número mínimo de folhas para que as árvores de decisão ainda
 266 mapeassem as sete classes de solos presentes na área de estudo, verifica-se que, para os dados
 267 oriundos do deslocamento de 100 m, o número mínimo de folhas foi de 37. Quando o
 268 deslocamento efetuado foi de 160 m, o número mínimo de folhas foi de 29. Por outro lado,
 269 com a utilização da totalidade das observações presentes na área de estudo, somente com
 270 árvores de maior complexidade, acima de 103 folhas, foi possível predizer todas as classes de
 271 solos. Esses resultados indicam que a presença de observações desviadas demanda árvores
 272 mais complexas para capturar o padrão presente entre os dados. Em estudo realizado por
 273 Giasson et al. (2011), a partir de 1.333 observações aleatoriamente distribuídas, os autores

274 verificaram que entre as árvores testadas apenas as mais complexas eram capazes de prever
275 todas as classes de solos.

276 Os resultados não permitem inferir qual seria a distância adequada de deslocamento das
277 bordas dos polígonos de solos. Essa distância será uma função de vários fatores, como a
278 escala, a complexidade dos elementos que governam a distribuição do solo, a experiência e
279 habilidade do pedólogo que mapeou a região que está sendo utilizada para gerar os modelos.
280 Possivelmente, em um mesmo mapa, poderiam ser adotadas distâncias distintas, uma vez que
281 as classes de solos apresentam um gradiente de complexidade nas suas transições. No entanto,
282 a análise da figura 4 permite afirmar que o pré-processamento pelo deslocamento dos
283 polígonos de solos, evitando utilizar as observações localizadas nas regiões de maior incerteza
284 do mapa, possibilita AD mais precisas e menos complexas.

285 Mapas de solos foram gerados por AD com um número mínimo de 30 folhas, a partir dos
286 três conjuntos de dados (Figura 5). Esse número foi adotado em virtude de ser o valor que
287 permite a melhor relação custo benefício entre complexidade e poder preditivo dos modelos
288 como verificado na figura 4. Embora, como anteriormente descrito, esse número mínimo de
289 folhas implique que as árvores geradas a partir do conjunto original e com deslocamento de
290 100 m não predizem todas as classes de solos.

291 Visualmente os Solos Hidromórficos (Planossolos e Gleissolos) foram espacializados de
292 maneira similar ao mapa de treinamento. Nos três mapas preditos, a classe Argissolo Bruno-
293 Acinzentado foi espacializada em uma posição intermediária na paisagem entre as áreas de
294 várzea e áreas mais elevadas das colinas, onde estão localizados os Argissolos Vermelho-
295 Amarelos. O modelo gerado pelo conjunto de dados oriundos do deslocamento de 160 m
296 possibilitou a espacialização da classe Argissolo Acinzentado, embora essa classe, presente
297 em outras regiões do mapa original, não esteja localizada na área dos extratos utilizados para
298 demonstrar os resultados (Figura 5).



299

300 Figura 5: a) Extrato do mapa de solos utilizado para treinamento. b) Predito a partir da
 301 totalidade dos dados. c) Predito a partir de deslocamento de 100 m nos polígonos. d) Predito a
 302 partir de deslocamento de 160 m nos polígonos. Solos Hidromórficos (SH), Cambissolo
 303 Háplico (CX), Neossolo Litólico (RL) e Argissolo Vermelho-Amarelo (PVA).

304

305 A comparação entre os três mapas de solos gerados e o mapa original permite verificar
 306 um satisfatório grau de concordância entre esses (Quadro 1). Aproximadamente, em 60% dos
 307 pontos (pixels) sobre a paisagem, a mesma classe existente no mapa original foi atribuída
 308 pelos modelos aos mapas preditos. Esse valor é superior aos valores encontrados na literatura
 309 de 43,0% (Crivelenti et al., 2009), 38,6% (Coelho & Giasson, 2010) e 51,8% (Giasson et al.,
 310 2011). Apesar de uma qualidade de predição superior obtida neste trabalho, em relação aos

311 trabalhos encontrados na literatura, observa-se que, com o deslocamento de 100 m, a classe
 312 dos Argissolos Acinzentados não foi predita. Esse fato provavelmente decorre de que a
 313 limitação de um número máximo de 30 folhas seja uma restrição muito forte à complexidade
 314 das informações existentes entre os dados advindos de posições na paisagem onde ocorram os
 315 Argissolos Acinzentados. Nesse caso, há a necessidade de que, para os conjuntos de dados
 316 mais dispersos, árvores mais complexas sejam geradas para a predição de todas as classes de
 317 solos.

318 Quadro 1. Índice kappa indicando o grau de concordância entre os mapas gerados a partir dos
 319 três conjuntos de dados deste estudo e o mapa de treinamento.

	Mapa de treinamento	Original*	100 m*
Original*	60,64	-	-
100 m*	60,16	90,5	-
160 m	59,34	87,6	88,85

320 * Classe de solos dos Argissolos Acinzentados não foi predita.

321

322 A comparação entre os três mapas de solos indica um elevado grau de similaridade entre
 323 esses (Quadro 1). Os três mapas foram idênticos em torno de 90% dos pontos na
 324 paisagem, embora tenha de se considerar que uma das classes de solos não tenha sido mapeada
 325 pela AD gerada por dois dos conjuntos de dados. Se considerarmos que o pré-processamento
 326 dos dados implicou uma redução do volume de observações, o conjunto de dados derivado do
 327 deslocamento de 160 m na borda dos polígonos, proposta de se fazer mais com menos (Liu &
 328 Motoda, 2002), foi alcançada. O índice kappa de 59,34% em relação ao mapa original foi
 329 alcançado com a totalidade das classes mapeadas, e com um volume de dados a ser
 330 manipulado, aproximadamente 60% menor do que o volume original.

331 O pré-processamento de dados apresentado difere da metodologia empregada por Qi &
 332 Zhu (2003), a qual emprega o conceito do histograma e da moda em dois aspectos principais.
 333 O primeiro deles está ligado ao fato de que a seleção das observações a partir de sua relação
 334 com a moda de uma covariável poderá representar a exclusão de informações relevantes.
 335 Podem ocorrer situações em que uma covariável apresente características multimodais em
 336 uma mesma classe de solo. Assim, é difícil determinar quantas são e quais são as modas
 337 presentes nos dados. Outro aspecto diz respeito ao volume de processamento necessário. O
 338 deslocamento dos polígonos de solos é automatizado, ao passo que a construção e análise dos
 339 histogramas irão demandar mais tempo para a análise.

340

341 **CONCLUSÕES**

- 342 1. O pré-processamento das observações reduz o volume de dados a ser manipulado no
343 MDS.
- 344 2. Observações presentes na borda dos polígonos de solos aumentam o número de
345 observações erroneamente classificadas por árvore de decisão.
- 346 3. Árvores de decisão geradas a partir de observações afastadas das bordas dos polígonos de
347 solos são menos complexas e têm maior poder preditivo em MDS.
- 348 4. Mapa de solos obtido por árvores de decisão a partir de observações não desviadas tem
349 mais similaridade ao mapa de treinamento.
- 350 5. A metodologia apresentada é de simples implementação e fácil aplicação para MDS.

351

352 **AGRADECIMENTOS**

353 À CAPES pelos recursos utilizados na execução do trabalho e ao CNPq pela bolsa de
354 produtividade em pesquisa do segundo autor e aporte financeiro para auxílio desse trabalho.

355

356 **LITERATURA CITADA**

357 CARVALHO, C.C.N.; FRANCA-ROCHA, W. & UCHA, J.M. Mapa digital de solos: Uma
358 proposta metodológica usando inferência fuzzy. Rev. Bras. Eng. Agríc. Ambient., 13:46-55,
359 2009.

360 CRIVELENTI, R. C.; COELHO, R. M.; ADAMI, S. F. & OLIVEIRA, S. R. de M. Mineração
361 de dados para a inferência de relações solo-paisagem em mapeamentos digitais de solo. Rev.
362 Agro. Bras. 44:1707-1715, 2009.

363 CONGALTON, R.G. A review of assessing the accuracy of classification of remotely sensed
364 data. Remote Sens. Environ. 37: 35-46, 1991.

365 COELHO, F.F. & GIASSON, E. Comparação de métodos para mapeamento digital de solos
366 com utilização de sistema de informação geográfica. Cienc. Rural, 40:2099-2106, 2010.

367 DALMOLIN, R. S. D.; KLAMT, E.; PEDRON, F. de A.; AZEVEDO, A.C. de. Relação entre
368 as características e o uso das informações de levantamentos de solos de diferentes escalas.
369 Cienc. Rural, 34:1479-1486, 2004.

370 ESRI Environmental Systems Research Institute, Inc., Redlands, CA, 2008.

371 GIASSON, E.; SARMENTO, E.C.; WEBER, E.; FLORES, C.A. & HASENACK, H.
372 Decision trees for digital soil mapping on subtropical basaltic steep slopes. Sci. Agríc., 68:167-
373 174, 2011.

- 374 GRINAND, C.; ARROUAYS, D.; LAROCHE, B. & MARTIN, M. Extrapolating regional
375 soil landscapes from an existing soil map: Sampling intensity, validation procedures, and
376 integration of spatial context. *Geoderma.*, 143:180-190, 2008.
- 377 HALL, M.; FRANK, E.; HOLMES, G.; PFAHRINGER, B.; REUTEMANN, P. & WITTEN,
378 I.H. The WEKA Data Mining Software: An Update. *SIGKDD Explorations Newsletter*,
379 11:10-18, 2009.
- 380 JENSEN, J.R. Sensoriamento Remoto do Ambiente: uma perspectiva em recursos terrestres.
381 São José dos Campos. Ed. Parêntese. 2009. 598p.
- 382 KLAMT, E.; FLORES, C. A. & CABRAL, D. R. Solos do Município de São Pedro do Sul.
383 Departamento. de Solos/CCR/UFSM. Santa Maria, 2001, 96p.
- 384 LEGROS, J. P. Mapping of the soil. Enfield: Science Publisher, 2005, 411p.
- 385 LIU, H. & MOTODA, H. On Issues of Instance Selection, *Data Min. Knowl. Disc.*, 6:115-
386 130, 2002.
- 387 MCBRATNEY, A. B.; MENDONCA SANTOS, M. L. & MINASNY, B. On digital soil
388 mapping. *Geoderma*, 117:3-52, 2003.
- 389 MORAN, C.J. & BUI, E.N. Spatial data mining for enhanced soil map modeling. *Int. J.*
390 *Geogr. Inf. Sci.*, 16:533-549, 2002.
- 391 OLAYA, V. A gentle introduction to SAGA GIS. The SAGA user group, Gottingen,
392 Germany, 208p.
- 393 QI, F. Knowledge discovery from area-class resource maps: data preprocessing for noise
394 reduction. *Transactions in GIS*. 8:297–308, 2004.
- 395 QI, F. & ZHU, A. X. Knowledge discovery from soil maps using inductive learning. *Int. J.*
396 *Geogr. Inf. Sci.*, 17:771-795, 2003.
- 397 TEN CATEN, A.; DALMOLIN, R.S.D.; PEDRON, F. de A. & MENDONÇA-SANTOS, M
398 de L. Extrapolação das relações solo-paisagem a partir de uma área de referência. *Cienc.*
399 *Rural.*, 41:812-816, 2011a.
- 400 TEN CATEN, A.; DALMOLIN, R.S.D.; PEDRON, F. de A. & MENDONÇA-SANTOS, M
401 de L. Regressões logísticas múltiplas: fatores que influenciam sua aplicação na predição de
402 classes de solos. *R. Bras. Ci. Solo*, 35:53-62, 2011b.
- 403 SCHMIDT, K; BEHRENS, T. & SCHOLTEN, T. Instance selection and classification tree
404 analysis for large spatial datasets in digital soil mapping. *Geoderma*, 146:138-146, 2008.
- 405 SCULL, P.; FRANKLIN, J.; CHADWICK, O.A. & MCARTHUR, D. Predictive soil
406 mapping: a review. *Prog. Phys. Geog.*, 27:171-197, 2003.

- 407 ZHU, A.X.; HUDSON, B.; BURT, J.; LUBICH, K. & SIMONSON, D. Soil mapping using
408 GIS, expert knowledge, and fuzzy logic. *Soil Sci. Soc. of Am. J.*, 65:1463-1472, 2001.
- 409 WAHBA, G. Spline models for Observational data. In: *CBMS-NSF Regional Conference*
410 *Series in Applied Mathematics*. Philadelphia: Soc. Ind. Appl. Maths., v.59, 1990, 169p.
- 411 WILSON, J. P. & GALLANT, J. C. Digital terrain analysis. In: _____ ed. *Terrain analysis:*
412 *principles and applications*. New York: Wiley & Sons, 2000, p. 1-27.

ARTIGO 4 - Mapeamento de classe de solos por árvore de decisão: impacto do volume de dados.¹

1

2 **Resumo:** Informações digitais possibilitam ao Mapeamento Digital de Solos (MDS) um
3 elevado grau de redundância para o ajuste de modelos preditivos de classes e propriedades de
4 solo. Entre esses modelos, a técnica de Árvores de Decisão (AD) tem aplicação crescente
5 devido a sua robustez no tratamento de grandes volumes de dados. Esse trabalho tem como
6 objetivo avaliar o impacto do volume de dados utilizados para gerar os modelos por AD na
7 qualidade dos mapas de solos preditos. A área de estudo, com 889,33 km², está localizada na
8 região do Planalto Médio do Rio Grande do Sul. As relações solo-paisagem foram obtidas a
9 partir de reambulação da área de estudo e delineamento das unidades de mapeamento em
10 cartas topográficas de escala 1:50.000. Seis covariáveis preditoras ligadas aos fatores de
11 formação do solo relevo e organismos, juntamente com os conjuntos de dados de um, três,
12 cinco, 10, 15, 20 e 25% do volume total de dados foram utilizados para gerar os modelos
13 preditivos por AD no programa WEKA. Nesse estudo, densidades de amostragem menores do
14 que 5% implicam modelos preditivos com menor poder de capturar a complexidade da
15 distribuição espacial do solo. Amostragens entre cinco e 15% provocaram melhor relação
16 custo-benefício quanto ao volume de dados a ser manipulado e à capacidade preditiva dos
17 modelos gerados. Dados coletados a campo indicaram uma acurácia dos mapas preditos
18 próxima a 70% para os modelos oriundos dessas densidades de amostragem.

19

INTRODUÇÃO

21 O avanço tecnológico em áreas como sensoriamento remoto, velocidade de
22 processamento dos computadores, possibilidade de manipular grandes volumes de dados,
23 métodos quantitativos para descrever padrões espaciais e visualização tridimensional tem
24 possibilitado uma nova oportunidade para prever processos e propriedades do solo

25 (Grundwald, 2009). Essa oportunidade foi denominada por McBratney et al. (2003) de
26 Mapeamento Digital de Solos (MDS).

27 A aplicação de informações digitais ao MDS possibilita um elevado grau de
28 redundância para o ajuste de modelos. Por outro lado, o volume de informações disponíveis
29 irá requerer a manipulação e o processamento de grandes volumes de dados (McBratney et
30 al., 2003). A busca pela melhor relação custo-benefício entre a densidade de amostragem no
31 banco de dados original e a acurácia preditiva dos modelos tem motivado pesquisas neste
32 sentido.

33 Entre as metodologias de mineração de dados aplicadas ao MDS, a técnica de Árvores
34 de Decisão (AD) tem aplicação crescente devido a sua robustez no tratamento de grandes
35 volumes de dados (Witten e Frank, 2005). A abordagem por AD tem sido empregada por
36 apresentar a vantagem de possibilitar a expressão das relações solo-paisagem de maneira
37 explícita (Kheir et al., 2010a). As AD foram empregadas em estudos relacionados à erosão
38 superficial e subterrânea de solos (Geissen et al., 2007), à predição de classes de solos
39 (Giasson et al., 2011) e à espacialização de propriedades do solo (Lemercier et al., 2011).

40 Além de possibilitar o agrupamento e a busca por padrões, a AD possibilita o
41 entendimento de como estes dados são inter-relacionados (Kheir et al., 2010b). A AD não
42 necessita que seja especificada a forma do modelo que será ajustado aos dados, bem como, à
43 medida em que um modelo por AD é construído, ele pode ser convertido em algoritmos que
44 podem ser facilmente implementados por linguagens de programação (Kheir et al., 2010a).

45 Modelos por AD foram utilizados por Scull et al. (2005) com 39.877 amostras para a
46 predição de classes de solos, em uma área de deserto de 2.590 km², correspondendo a,
47 aproximadamente, um por cento do total de dados disponíveis. Em estudo realizado por
48 Giasson et al. (2011), foram utilizadas 1.333 amostras para a predição de classes de solos em
49 uma área aproximadamente 6,1 km². A densidade de amostras correspondeu a menos de um

50 por cento do volume de dados original. Qi e Zhu (2003) amostraram, em uma área de 4 km²,
51 cerca de 12% dos dados originais para o ajuste de modelos por AD.

52 Frente à disponibilização de bancos de dados com um crescente número de
53 observações e a diversidade de densidade de amostragens relatada nos estudos realizados, este
54 trabalho tem como objetivo avaliar a influência da proporção do conjunto de dados originais
55 na qualidade do mapa de solos gerado pela técnica de árvore de decisão.

56

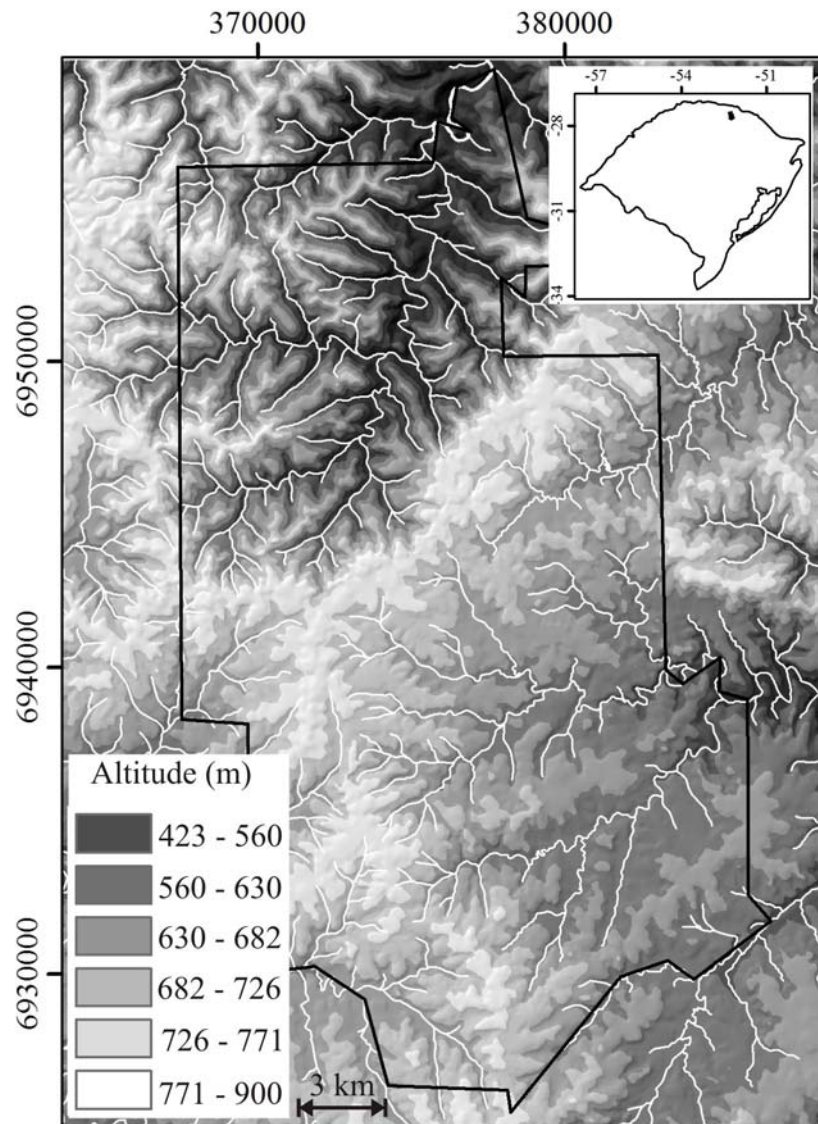
57 **MATERIAL E MÉTODOS**

58 **Área de estudo**

59 A área de estudo está definida pelo polígono do município de Erechim, no Planalto
60 Médio do Rio Grande do Sul, entre as coordenadas SIRGAS2000 52°06' e 52°24' de
61 longitude Oeste e 27°28' e 27°47' de latitude Sul. A área compreende 889,33 km² e possui
62 uma elevação local que varia entre 423 a 900 m (Figura 1). A região foi escolhida por possuir
63 mapas de solos apenas em nível de reconhecimento e exploratório (ten Caten et al., 2007),
64 fazendo dessa uma região carente por informação espacial em solos.

65 Visando gerar as AD, as relações solo-paisagem foram obtidas a partir de reambulação
66 da área de estudo com informações obtidas no campo por um pedólogo experiente. Foram
67 descritos perfis e coletadas amostras de solos para análise em laboratório. Essas informações
68 possibilitaram delinear as unidades de mapeamento em cartas topográficas de escala 1:50.000
69 e definir as classes de solos de ocorrência significativa de acordo com o Sistema Brasileiro de
70 Classificação de Solos - SiBCS (EMBRAPA, 2006) (Tabela 1). Essas informações, em
71 conjunto com as covariáveis ambientais, foram utilizadas como dados de treinamento para as
72 AD.

73



74

75 Figura 1: Área do estudo com o perímetro do município de Erechim sobreposto à rede de
 76 drenagem e ao mapa hipsométrico da região. Imagem interna, no canto superior direito,
 77 localiza a área do estudo no Estado do Rio Grande do Sul.

78 **Covariáveis ambientais**

79 Neste estudo, foram utilizadas covariáveis do modelo 'scorpan' (McBratney et al,
 80 2003), relacionadas aos fatores de formação do solo relevo (r) e organismos (o). As demais
 81 covariáveis do modelo não estão disponíveis ou não possuem a resolução espacial demandada
 82 para este estudo. O fator organismos foi representado pela covariável Índice de Vegetação por
 83 Diferença Normalizada (NDVI) do mês de outubro (NOUT) e pelo desvio padrão do NDVI

84 (DEPA). A opção pelo NDVI do mês de outubro foi devido a este ser o mês de maior
 85 contraste entre os diferentes usos da terra na região. Para o cálculo de ambas as covariáveis,
 86 as imagens utilizadas são todas do ano de 2004 devido à maior disponibilidade de imagens
 87 com ausência de nuvens durante aquele ano. Todos os cálculos de NDVI foram realizados
 88 conforme Jensen (2009).

89 Tabela 1 – Classes de solos encontradas na área de estudo.

Classe SiBCS - legenda	Área (ha)	Fração total (%)
GLEISSOLO HÁPLICO - GX	31,28	12,52
LATOSSOLO VERMELHO - LV	122,99	49,22
NEOSSOLO LITÓLICO - RL	26,26	10,51
Associação NEOSSOLO LITÓLICO e CAMBISSOLO HÁPLICO – RL/CX	18,27	7,31
Associação NITOSSOLO VERMELHO e CAMBISSOLO HÁPLICO – NV/CX	7,05	2,82
Associação NITOSSOLO VERMELHO e CHERNOSSOLO ARGILÚVICO fase altitude elevada – NV/MT fae	26,01	10,41
Associação NITOSSOLO VERMELHO e CHERNOSSOLO ARGILÚVICO fase altitude baixa – NV/MT fab	18,03	7,22

90

91 O preditor DEPA foi gerado a partir de dados de oito distintas datas do ano de 2004,
 92 obtidas pela plataforma Landsat 5 sensor TM com resolução espacial de 30 m. Cada data teve
 93 seu valor de NDVI calculado individualmente. Em seguida, foi executado o cálculo do desvio
 94 padrão do valor do NDVI entre as oito datas para cada pixel no programa ArcGIS (ESRI,
 95 2008), na função Raster Calculator. O desvio padrão foi calculado como segue:

$$96 \quad DEPA = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (NDVI_i - \overline{NDVI})^2}$$

97 onde

98 $DEPA$ = desvio padrão do NDVI;

99 n = número de distintas datas utilizadas para gerar o DEPA;

100 $NDVI_i$ = i-ésimo NDVI;

101 \overline{NDVI} = valor médio do NDVI.

102 O fator relevo foi representado pelas covariáveis elevação (ELEV), distância vertical
 103 da rede de drenagem (DIST), declividade (DECL), índice de umidade topográfica (IUT),
 104 capacidade de transporte de sedimento (CTS), curvatura planar (PLAN), curvatura de perfil
 105 (PERF) e rugosidade (RUGO). As covariáveis foram geradas de acordo com Wilson e Gallant
 106 (2000), a partir de um modelo digital de elevação (MDE). O MDE com resolução espacial de
 107 30 m foi derivado de curvas de nível de cartas topográficas com escala 1:50.000. A
 108 interpolação das curvas de nível para o formato matriz ocorreu na ferramenta *topo do raster*
 109 do programa ArcGIS 9.3 (ESRI, 2008), usando a técnica *spline* (Wahba, 1990). No programa
 110 SAGA-GIS (Olaya, 2004), os atributos ELEV, DIST, DECL, PLAN, PERF, e RUGO foram
 111 gerados na ferramenta *standard terrain analysis*, e os atributos IUT e CTS com a ferramenta
 112 *grid calculator*.

113 A tabulação dos dados para tratamento estatístico ocorreu a partir da amostragem do
 114 mapa de solos de treinamento sobreposto aos dez planos de informação das covariáveis
 115 ambientais. Foi gerada uma tabela com 277.491 linhas e 11 colunas, na qual estão
 116 representados 249,80 km² de superfície, ou seja, 28% da área total do município de Erechim.

117 A tabela foi submetida à análise estatística multivariada para a verificação de presença
 118 de multicolinearidade entre os preditores, tendo sido executado o teste de correlação de

119 Pearson e Análise de Componentes Principais (ACP) (Hair et al., 2006). A análise de
120 correlação e a ACP foram executadas no programa R (R DEVELOPMENT CORE TEAM,
121 2011).

122 **Árvores de decisão**

123 A tabela com 28% dos dados originais foi reamostrada em sete conjuntos de dados. Os
124 conjuntos corresponderam às densidades de amostragem de um, três, cinco, 10, 15, 20 e 25%
125 da área total do estudo. Todas as classes de solos foram amostradas proporcionalmente a sua
126 área, conforme Scull et al. (2005). O desenvolvimento das AD foi realizado no programa de
127 mineração de dados WEKA (Hall et al., 2009). Para o processamento dos dados, foi utilizado
128 o algoritmo J48 devido a esse ter tido o melhor desempenho em estudo realizado por Giasson
129 et al. (2011). O número mínimo de observações (*minNumObj* no WEKA) por folha para cada
130 conjunto de dados foi determinado após uma análise do percentual de observações
131 erroneamente classificadas. Essa análise foi executada com uma série de valores para o
132 número mínimo de observações em cada um dos sete conjuntos de dados. O método de poda,
133 que mitigava o erro na árvore gerada, também foi selecionado (*reducedErrorPruning = True*
134 no WEKA). Durante a fase de geração da árvore, foi utilizado um conjunto de dados
135 independente de cinco por cento da área total (*Supplied test set* no WEKA), com o intuito de
136 verificar a qualidade da árvore gerada.

137 **Mapa de solos**

138 As AD que apresentaram a melhor relação custo-benefício entre a complexidade do
139 modelo e o número de observações erroneamente classificadas foram implementadas no
140 programa ArcGIS, na função *raster calculator*. As informações derivadas da árvore foram
141 convertidas na função condicional ‘con(teste, verdadeiro, falso)’ do programa. Essa função
142 permite que os arquivos raster com as covariáveis ambientais sejam processados de acordo
143 com o conjunto de regras derivadas da AD. Como a escala de publicação pretendida é de

144 1:50.000, em cada mapa de solos gerado, foram extraídas regiões de píxeis isolados com área
145 mínima mapeável menor do que um hectare.

146 **Qualidade dos modelos e mapas**

147 A qualidade dos modelos AD foi avaliada a partir do valor percentual de observações
148 erroneamente classificadas em toda a árvore. Esse valor é uma das saídas do programa
149 WEKA após o conjunto de dados teste ser processado no modelo gerado. Para seu cálculo, o
150 programa soma todas as observações erroneamente classificadas e as divide pelo total de
151 observações do conjunto teste, por fim multiplica por 100 para gerar o valor percentual (Hall
152 et al., 2009).

153 A precisão do modelo ajustado foi avaliada pelo conjunto de classes corretamente
154 classificadas e divididas pelo somatório do conjunto de classes correta e erroneamente
155 classificadas. Em seguida, esse valor é multiplicado por 100 para se ter um valor relativo
156 (Hall et al., 2009).

157 O índice kappa foi utilizado para verificar a qualidade da informação contida no mapa
158 (Hengl et al., 2007). A matriz dos erros para o cálculo do índice kappa foi executada
159 conforme Congalton (1991), a partir de reambulação a campo para identificação de 50 locais,
160 por classe, do solo ocorrente. No total, foram identificados 350 locais de ocorrência das
161 classes de solos mapeadas neste estudo.

162

163 **RESULTADOS E DISCUSSÃO**

164 **Estudo da multicolinearidade entre os preditores**

165 A análise de correlação nas 277.491 linhas de dados pode ser visualizada na Tabela 2.
166 Apenas a correlação entre PERF e DECL mostrou-se não significativa. As covariáveis DEPA
167 e NOUT mostraram ter correlação negativa (-0,480). A existência dessa correlação está ligada
168 ao fato de que o NDVI do mês de outubro (NOUT) foi utilizado para o cálculo do desvio
169 padrão do ano de 2004 (DEPA). Além disso, no mês de outubro, ocorrem grandes contrastes

170 dos valores de NDVI na paisagem. As áreas com mata irão apresentar elevados índices de
 171 vegetação, ao passo que as áreas cultivadas com a cultura do trigo senescida ou colhida irão
 172 apresentar os menores valores para o índice de vegetação. Estes locais de grande variação do
 173 NDVI ao longo do ano geram os maiores valores da covariável DEPA. Os menores valores da
 174 covariável NOUT encontram-se espacialmente nos mesmos locais onde DEPA apresenta seus
 175 maiores valores, por isso a correlação negativa.

176 Tabela 2 – Matriz de correlação entre covariáveis ambientais.

	DEPA	NOUT	CTS	IUT	ELEV	PLAN	PERF	DECL	DIST	RUGO
DEPA	1									
NOUT	-0,480 ^a	1								
CTS	-0,197 ^a	0,172 ^a	1							
IUT	0,049 ^a	-0,066 ^a	0,122 ^a	1						
ELEV	0,129 ^a	-0,113 ^a	-0,209 ^a	-0,152 ^a	1					
PLAN	0,057 ^a	-0,058 ^a	-0,345 ^a	-0,336 ^a	0,060 ^a	1				
PERF	0,076 ^a	-0,044 ^a	-0,131 ^a	-0,170 ^a	0,290 ^a	0,286 ^a	1			
DECL	-0,335 ^a	0,302 ^a	0,345 ^a	-0,515 ^a	-0,175 ^a	0,052 ^a	0,005 ^b	1		
DIST	-0,132 ^a	0,142 ^a	-0,014 ^a	-0,416 ^a	0,278 ^a	0,244 ^a	0,428 ^a	0,368 ^a	1	
RUGO	-0,291 ^a	0,245 ^a	0,278 ^a	-0,294 ^a	-0,150 ^a	0,018 ^a	-0,021 ^a	0,715 ^a	0,250 ^a	1

177 ^a: $p < 0.001$; ^b: não significativo. NDVI do mês de outubro de 2004 (NOUT), desvio padrão
 178 do NDVI no ano de 2004 (DEPA), elevação (ELEV), distância vertical da rede de drenagem
 179 (DIST), declividade (DECL), índice de umidade topográfica (IUT), capacidade de transporte
 180 de sedimento (CTS), curvatura planar (PLAN), curvatura de perfil (PERF) e rugosidade
 181 (RUGO).

182

183 A maior correlação foi verificada entre as covariáveis RUGO e DECL (0,715). A
 184 existência de correlação entre essas variáveis se deve ao fato de que a declividade é uma das
 185 variáveis utilizadas para o cálculo da rugosidade, assim como os maiores valores de ambas as
 186 covariáveis foram verificados espacialmente nos mesmos locais, disto a correlação positiva.

187 O resultado da rotação dos eixos originais, visando potencializar a variabilidade em
 188 um novo conjunto de variáveis, pode ser visualizado na Tabela 3. Cerca de 78% da variância
 189 original fica distribuída entre as cinco primeiras componentes, com percentuais consideráveis
 190 de variabilidade retida entre a sexta e a nona componente. Valores residuais de variância nas
 191 componentes mais elevadas é um indicativo de um relativo grau de independência entre as
 192 variáveis originais.

193 Tabela 3 – Resultado da distribuição da variabilidade nas dez componentes principais.

Componente Principal	Autovalor	Variância (%)	Variância acumulada (%)
1	2,728	27,276	27,276
2	2,189	21,892	49,168
3	1,105	11,046	60,215
4	1,032	10,324	70,539
5	0,779	7,794	78,333
6	0,537	5,367	83,700
7	0,510	5,104	88,804
8	0,483	4,831	93,635
9	0,441	4,408	98,043
10	0,196	1,957	100,000

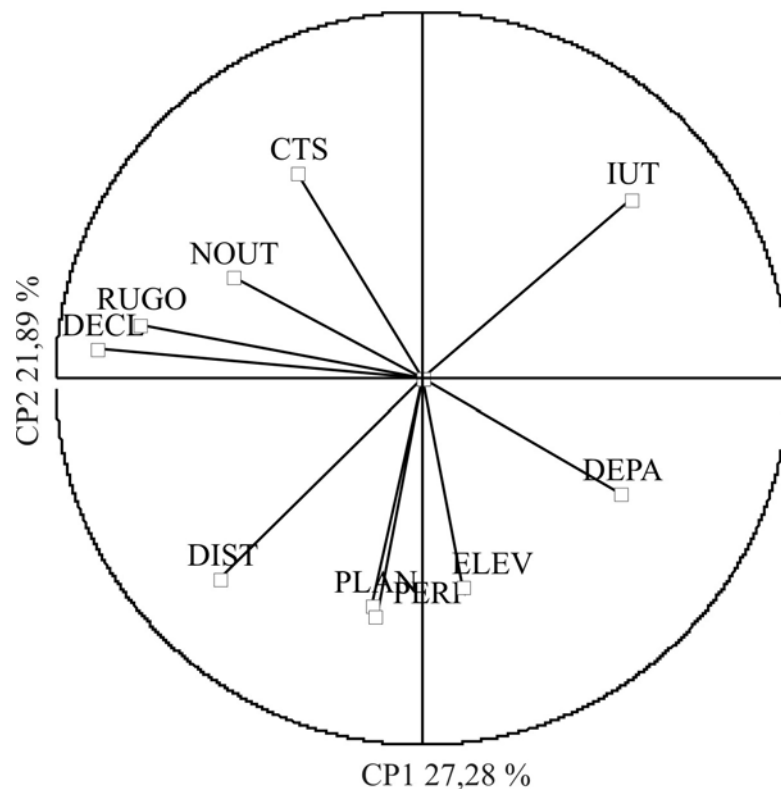
194

195 Resultados similares foram identificados por Hengl e Rossiter (2003), que utilizaram
 196 de ACP para mitigar o efeito de multicolinearidade entre atributos do terreno. Os autores
 197 identificaram 80% da variabilidade retida nas cinco primeiras componentes geradas a partir da
 198 rotação ortogonal de um conjunto de nove atributos do terreno. Para os autores, esse resultado
 199 demonstrou que as covariáveis eram menos redundantes do que se havia inicialmente
 200 pressuposto. Em estudo realizado por ten Caten et al. (2011), as três primeiras componentes
 201 principais, as quais reteram 65% da variabilidade original dos dados, possibilitaram a geração
 202 de modelos preditivos com desempenho similar àquelas que se utilizaram da totalidade das

203 covariáveis predictoras, com a vantagem de demandar que um menor número de dados fosse
 204 manipulado.

205 Uma visão mais detalhada das correlações entre as covariáveis pode ser obtida
 206 analisando-se a distribuição dessas no diagrama de ordenação unitário para a duas primeiras
 207 componentes principais (Figura 2). Nessas componentes, foi retida cerca de 50% da
 208 variabilidade original dos dados. A proximidade entre as covariáveis DECL e RUGO, bem
 209 como entre PLAN, PERF e ELEV indica haver entre elas uma redundância mais intensa do
 210 que em meio às covariáveis espacialmente mais distantes no gráfico. A posição do IUT
 211 reproduz o que pode ser verificado na Tabela 1, ou seja, a covariável é inversamente
 212 correlacionada com praticamente todas as covariáveis, à exceção da CTS e DEPA.

213



214

215 Figura 2 – Diagrama de ordenação unitário com as coordenadas das covariáveis ambientais
 216 para as duas primeiras componentes principais. NDVI do mês de outubro de 2004 (NOUT),
 217 desvio padrão do NDVI no ano de 2004 (DEPA), elevação (ELEV), distância vertical da rede

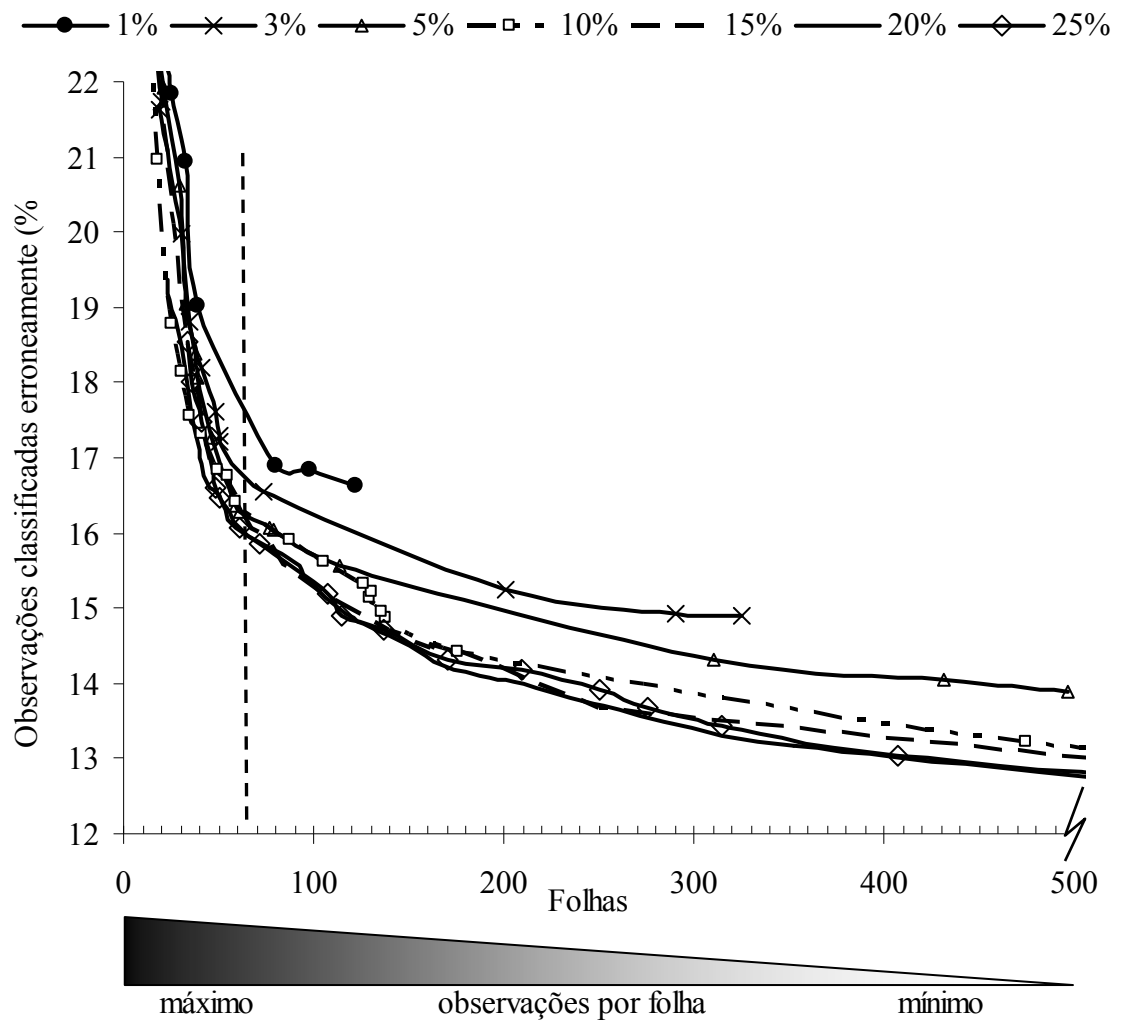
218 de drenagem (DIST), declividade (DECL), índice de umidade topográfica (IUT), capacidade
219 de transporte de sedimento (CTS), curvatura planar (PLAN), curvatura de perfil (PERF) e
220 rugosidade (RUGO).

221

222 A partir do estudo do padrão da variabilidade presente nas covariáveis preditoras,
223 decidiu-se simplificar o conjunto de dados utilizado para gerar os modelos por árvore de
224 decisão. Foram descartadas da sequência do estudo as covariáveis PLAN, PERF, RUGO e
225 NOUT.

226 **Complexidade do modelo por árvore de decisão**

227 AD com um grande número de folhas se superajustam aos dados e atingem seu
228 máximo potencial preditivo (Figura 3). Disso decorre que, durante o desenvolvimento dos
229 modelos por árvore de decisão, faz-se necessário avaliar o ganho preditivo dos modelos na
230 medida em que esses se tornam mais complexos. Para os sete conjuntos de amostras, verifica-
231 se que na medida em que as AD tornam-se mais complexas, com um maior número de folhas,
232 menor é o número de observações erroneamente classificadas pelos modelos. No entanto,
233 mesmo para árvores muito complexas e maiores volumes de dados para ajustar os modelos, o
234 percentual de observações erroneamente classificadas permanece em torno de 13% (Figura 3).
235 Possivelmente esse fato decorra da impossibilidade de se explicar a complexidade da
236 distribuição espacial das classes de solos na região apenas com as seis covariáveis utilizadas.
237 Scull et al. (2005) observaram que em torno de 18,52% das observações localizadas na área
238 de treinamento e 24,34% daquelas localizadas na área de validação eram classificadas
239 erroneamente na fase de geração de modelos de AD para o nível de ordem.



240

241 Figura 3 – Relação entre o número de folhas na árvore de decisão e o percentual de amostras
 242 erroneamente classificadas.

243 Para os conjuntos de dados de um e três por cento, mesmo com o menor número de
 244 observações por folha em árvores mais complexas, não são alcançados resultados melhores do
 245 que 16,5 e 15%, respectivamente. Entre os efeitos de se utilizar um número pequeno de
 246 amostras para a geração das AD, pode estar a não representatividade desses em relação ao
 247 conjunto total de dados disponíveis, acarretando em um menor poder preditivo.

248 À medida que as AD são podadas, ocorre um aumento no número de observações
 249 erroneamente classificadas (Figura 3). No entanto, a partir de um determinado número de
 250 folhas, a taxa de incremento dos erros modifica-se sensivelmente e passa a aumentar

251 linearmente para AD com um número de folhas menor do que 65 (linha vertical tracejada).
 252 Esse comportamento também é verificado nos menores conjuntos de dados de um e três por
 253 cento, embora aconteça de maneira não tão marcante como nos demais conjuntos de dados.

254 Buscando utilizar modelos que associassem simplicidade e poder preditivo nos
 255 diferentes conjuntos de dados, optou-se por utilizar o valor de 65 folhas como limite mínimo
 256 para o número de folhas nas AD. A partir desse valor, foram geradas AD para os conjuntos de
 257 dados de cinco, 10, 15, 20 e 25% (Tabela 4). Os conjuntos de dados de um e três por cento
 258 foram descartados por apresentarem os piores resultados quanto ao número de observações
 259 erroneamente classificadas na região de AD, com 65 folhas. Grinand et al. (2008) afirmam
 260 que volumes de amostras menores do que 10% podem levar a uma redução substancial da
 261 qualidade dos modelos.

262 Tabela 4 – Informações de saída do programa WEKA a partir do teste dos cinco modelos no
 263 conjunto de dados independentes (*Supplied test set*).

Conjuntos de dados (%)	Observações classificadas erroneamente (%)	Kappa (%)	Precisão (%)
5	16,27	76,50	83,00
10	16,23	76,80	83,20
15	16,14	76,94	83,30
20	15,79	77,36	83,70
25	16,02	77,06	83,60

264

265 Os resultados da Tabela 4 indicam que houve uma tendência para uma pequena
 266 diminuição do número de observações erroneamente classificadas à medida que o conjunto de
 267 dados utilizado para gerar os modelos era aumentado. O fato de o número de observações
 268 incorretas permanecer em torno de 16% para todos os cinco conjuntos de dados,
 269 possivelmente se deva à qualidade e quantidade das covariáveis preditoras utilizadas neste
 270 estudo, as quais foram incapazes de predizer a total complexidade da distribuição espacial da
 271 classes de solos na área de estudo. Os valores do índice Kappa e de precisão dos modelos
 272 gerados pelo programa WEKA também apresentaram sensível melhora à medida que um

273 maior volume de dados foi utilizado para gerar os modelos (Tabela 4). Contudo, esses
274 resultados indicam que a melhoria do ajuste dos modelos gerados por AD deverá vir
275 associada a outros fatores, tais como a utilização de uma maior quantidade e diversidade de
276 covariáveis predictoras, e não somente a utilização de um maior volume de dados para ajustar
277 os modelos.

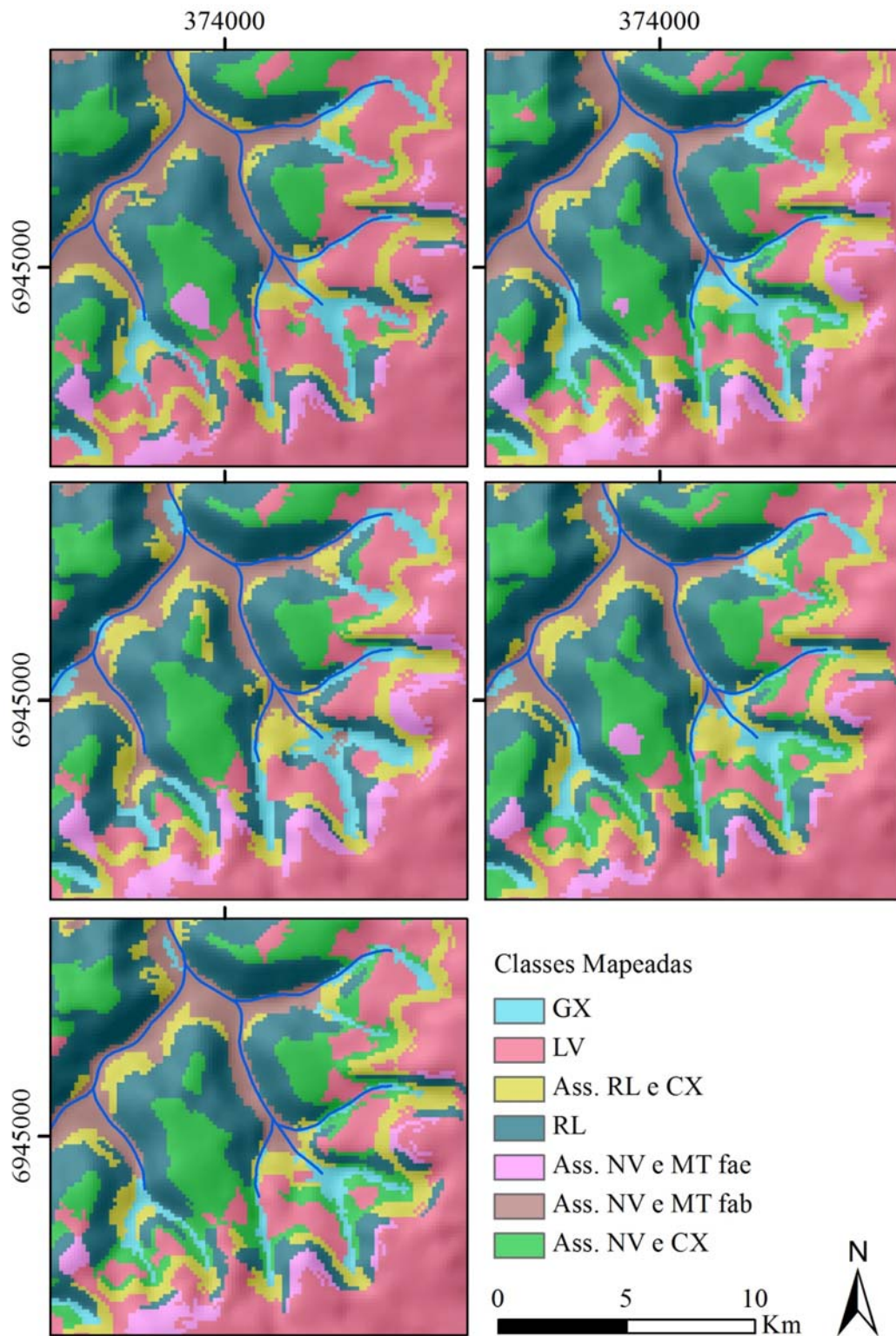
278 Em estudo realizado por Giasson et al. (2011), os autores identificaram um índice
279 kappa de 57,1% para uma AD com 79 folhas, sendo utilizada para mapear seis classes de
280 solos. De acordo com esses autores, AD muito simplificadas implicavam um reduzido poder
281 preditivo dos modelos. Naquele estudo, foi avaliada uma AD com apenas cinco folhas,
282 implicando um índice kappa de aproximadamente 43,9% e com apenas duas classes de solos
283 sendo preditas.

284 **Mapas de solos**

285 As cinco AD foram utilizadas para gerar mapas de solos para o município de Erechim
286 (Figura 4). Todos os cinco mapas de solos reproduziram a relação solo-paisagem, formalizada
287 pelo pedólogo de maneira muito similar. Nas colinas localizadas em locais de maior altitude,
288 foram posicionados os Latossolos Vermelhos (Figura 4). Nesses locais, as condições de
289 relevo possibilitaram a formação desses solos profundos, homogêneos e bem intemperizados.

290 Em posições mais baixas na toposequência, foram espacializadas as associações entre
291 Neossolos Litólicos e Cambissolos Háplicos (Figura 4). Essa associação encontra-se ao longo
292 de toda a parte superior do vale do Rio Dourado. Nesses locais, a declividade possibilita a
293 remoção constante de material. Como consequência, tem-se a presença de solos rasos e com
294 grande presença de pedregosidade. Os Neossolos Litólicos foram espacializados,
295 predominantemente, nas rampas mais longas e declivosas da toposequência (Figura 4). A
296 associação Nitossolo Vermelho e Chernossolo Argilúvico fase altitude baixa foi espacializada
297 ao longo da drenagem dos rios nas posições encaixadas dos vales. Nesses locais, encontram-

298 se áreas planas que possibilitaram a formação desses solos mais profundos e desenvolvidos
 299 (Figura 4).



300

301 Figura 4 – Sessões dos mapas de solos preditos a partir dos cinco conjuntos de amostras a)

302 5%, b) 10%, c) 15%, d) 20%, e) 25%. Coordenadas UTM Fuso 22 / SIRGAS2000.

303 A análise visual dos cinco mapas gerados não fornece subsídios suficientes para a
 304 decisão de qual densidade de dados é a mais adequada. Na sequência do estudo, fez-se a
 305 avaliação em campo da capacidade preditiva de cada uma das AD geradas a partir dos cinco
 306 conjuntos de dados.

307 **Qualidade dos mapas de solos**

308 A partir dos cinco mapas de solos gerados e dos dados coletados a campo, foi possível
 309 avaliar a qualidade dos modelos preditivos (Tabela 5). Optou-se, primeiramente, por avaliar a
 310 semelhança entre os mapas gerados a partir da comparação entre eles. Verifica-se que, à
 311 medida que o volume de dados utilizado aumenta, maior também é a concordância entre os
 312 mapas gerados. Quando mapas com volumes de dados de 5 e 10% são comparados, o índice
 313 kappa é menor do que o índice obtido da comparação entre mapas com volumes de dados de
 314 20 e 25%. Esse resultado é um indicativo de que volumes de dados maiores possibilitam
 315 capturar melhor a complexidade a ser mapeada, embora as semelhanças encontradas entre os
 316 mapas preditos com volumes menores de dados, neste estudo, foram também elevadas.

317 Tabela 5 – Índice kappa entre os mapas gerados com diferentes densidades de amostragem,
 318 com a área de treinamento e com amostras no campo.

	5%	10%	15%	20%	25%
5%	-	-	-	-	-
10%	79,4	-	-	-	-
15%	78,4	81,1	-	-	-
20%	81,1	84,1	84,1	-	-
25%	82,4	85,3	85,6	90,4	-
treinamento	77,3	77,3	77,1	77,7	77,6
campo	69,7	68,5	71,1	71,4	71,1

319

320 A reprodutibilidade dos cinco modelos por AD empregados neste estudo pode ser
 321 visualizada na sexta linha da Tabela 5. Observa-se que todos os modelos reproduziram, com

322 um valor superior a 77%, a classe de solo que havia sido definida pelo pedólogo nas áreas de
323 treinamento dos modelos. Esse valor é superior aos 46,12% encontrados por ten Caten et al.
324 (2011), que utilizaram de regressões logísticas múltiplas como modelos preditivos.

325 Embora o presente estudo tenha uma reprodutibilidade elevada, em aproximadamente
326 23% da área de treinamento, os modelos divergiram do que havia sido definido como sendo a
327 classe de solo ocorrente. Esse fato pode estar ligado a erros na descrição do solo a campo ou
328 na incapacidade dos modelos em capturar a complexidade inerente à distribuição espacial das
329 classes de solos no local.

330 A acurácia dos modelos preditivos está representada na última linha da Tabela 5.
331 Todos os modelos possibilitaram uma acurácia próxima a 70% no mapeamento das classes de
332 solos locais. Como era esperado, a acurácia dos modelos preditivos é um pouco inferior à
333 reprodutibilidade dos modelos comentada nos parágrafos anteriores. Contudo, ambos os
334 indicadores da qualidade do mapeamento indicam que os cinco modelos testados
335 desempenham de maneira satisfatória o mapeamento de classes de solos.

336 Quanto à influência do volume de dados utilizados para gerar as AD, o que se observa
337 é que há uma tendência de melhor desempenho dos modelos preditivos quando o volume de
338 dados utilizado para treinar os modelos é maior (Tabela 5). O melhor desempenho preditivo
339 foi atingido pelos volumes de dados de 15, 20 e 25%. Estudo realizado por Grinand et al.
340 (2008) testou amostras de 10 a 90% dos dados originais. Os autores verificaram que a grande
341 variação da acurácia estava entre 10 e 20% dos dados, não sendo verificado significativo
342 aumento de qualidade dos modelos a partir de 30% do volume total de amostras. De acordo
343 com os resultados aqui encontrados, volumes de dados de 15% aliariam um menor volume de
344 dados, reprodutibilidade e acurácia no mapeamento.

345 O mapeamento de classes de solos na carta topográfica de Dois Córregos (SP) foi
346 executado por Crivelenti et al. (2009) com um volume de dados de 714.000 amostras em uma

347 área total do estudo de 772 km², e embora não tenha sido possível a determinação exata do
348 volume de dados utilizados para a geração das árvores de decisão naquele estudo, acredita-se
349 que mais de 80% dos dados originais tenham sido utilizados para gerar os modelos. Os
350 autores reportaram um índice kappa de 43% em relação ao mapa utilizado para gerar o
351 modelo.

352 Estudos realizados por Bui e Moran (2003) demonstraram que a capacidade preditiva
353 dos modelos experimentou um significativo incremento a partir da utilização de 10% dos
354 dados originais. A qualidade dos modelos preditivos teve um valor máximo quando as
355 amostras utilizadas representaram 25% do total de dados, com um índice kappa de 64% em
356 relação ao mapa original. Contudo, os autores afirmam que a menor densidade de
357 amostragem, de 10%, já seria suficiente para capturar grande parte da diversidade contida nos
358 dados para a construção dos modelos preditivos por AD.

359

360 **CONCLUSÃO**

361 O presente estudo demonstrou que conjuntos de dados pouco representativos (menores
362 do que 5%), para a geração de modelos de Árvore de Decisão, têm dificuldades em capturar a
363 complexidade existente na distribuição espacial de solos. Por outro lado, conjuntos de dados
364 superiores a 20% podem implicar em redundância excessiva para a geração de modelos
365 preditivos, implicando manipulação desnecessária de dados, o que não irá refletir em ganho
366 preditivo pelos modelos derivados desses conjuntos de dados. Nas condições em que o
367 presente estudo foi realizado, volumes de dados entre cinco e 15% implicaram na melhor
368 relação custo-benefício quanto ao volume de dados a ser manipulado e à capacidade preditiva
369 dos modelos gerados. Valores de reprodutibilidade em torno de 77% e acurácia próxima a
370 70% foram alcançados com esses volumes de amostras utilizados para gerar os modelos por
371 árvore de decisão.

372

373 AGRADECIMENTOS

374 À CAPES pelos recursos utilizados na execução do trabalho e ao CNPq pela bolsa de
375 produtividade em pesquisa do segundo autor e aporte financeiro para auxílio desse trabalho.

376

377 REFERÊNCIAS

378 Bui, E. N.; Moran, C. J. 2003. A strategy to fill gaps in soil survey over large spatial extents:
379 an example from the Murray-Darling basin of Australia. *Geoderma* 111: 21-44.

380 Crivelenti, R.C. et al. 2009. Mineração de dados para a inferência de relações solo-paisagem
381 em mapeamentos digitais de solo. *Revista Agropecuária Brasileira*. 44: 1707-1715.

382 Congalton, R.G. 1991. A review of assessing the accuracy of classification of remotely sensed
383 data. *Remote Sensing of Environment* 37: 35-46.

384 Empresa Brasileira de Pesquisa Agropecuária [EMBRAPA]. 2006. Sistema brasileiro de
385 classificação de solos. Centro Nacional de Pesquisa em Solos, Rio de Janeiro, RJ, Brasil.

386 Environmental Systems Research Institute [ESRI], Inc., Redlands, CA, 2008.

387 Giasson, E.; Sarmiento, E.C.; Weber, E.; Flores, C.A.; Hasenack, H. 2011. Decision trees for
388 digital soil mapping on subtropical basaltic steeplands. *Scientia Agrícola* 68: 167-174.

389 Geissen, V.; Kampichler, C.; Lopez-de Llergo-Juarez, J.J.; Galindo-Acantara, A. 2007.
390 Superficial and subterranean soil erosion in Tabasco, tropical Mexico: Development of a
391 decision tree modeling approach. *Geoderma* 139: 277-287.

392 Grinand, C.; Arrouays, D.; Laroche, B.; Martin, M. 2008. Extrapolating regional soil
393 landscapes from an existing soil map: Sampling intensity, validation procedures, and
394 integration of spatial context. *Geoderma* 143: 180-190

395 Grunwald, S. 2009. Multi-criteria characterization of recent digital soil mapping and
396 modeling approaches. *Geoderma* 152: 195-207.

397 Hair, J.F.; Anderson, R.E.; Tatham, R.L.; Black, W. C. 2006. *Análise Multivariada de Dados*.
398 5ed. Bookman, Porto Alegre, RS, Brasil.

- 399 Hall, M.; Frank, E.; Holmes, G.; Pfahringer, B.; Reutemann, P.; Witten, I.H. 2009. The
400 WEKA Data Mining Software: An Update; SIGKDD Explorations, 11.
- 401 Hengl, T.; Rossiter, D.G. 2003. Supervised Landform classification to enhance and replace
402 photo-interpretation in semi-detailed soil survey. *Soil Science Society of America Journal* 67:
403 1810-1822.
- 404 Hengl, T.; Toomanian, N.; Reuter, H.I.; Malakouti, M.J. 2007. Methods to interpolate soil
405 categorical variables from profile observations: Lessons from Iran. *Geoderma* 140: 417-427.
- 406 Jensen, J.R. 2009. *Sensoriamento Remoto do Ambiente: uma perspectiva em recursos*
407 *terrestres*. Ed. Parêntese, São José dos Campos, SP, Brasil.
- 408 Kheir, R.B. et al. 2010a. Predictive mapping of soil organic carbon in wet cultivated lands
409 using classification-tree based models: The case study of Denmark. *Journal of Environmental*
410 *Management* 91: 1150-1160.
- 411 Kheir, R.B. et al. 2010b. Spatial soil zinc content distribution from terrain parameters: A GIS-
412 based decision-tree model in Lebanon. *Environmental Pollution* 158: 520-528.
- 413 Lemercier, B.; Lacoste, M.; Loum, M; Walter, C. 2011. Extrapolation at regional scale of
414 local soil knowledge using boosted classification trees: A two-step approach, *Geoderma*, In
415 Press, Corrected Proof.
- 416 McBratney, A. B.; Mendonça-Santos, M. L.; Minasny, B. 2003. On digital soil mapping.
417 *Geoderma* 117: 3-52.
- 418 Olaya, V. 2004 A gentle introduction to SAGA GIS. The SAGA user group, Göttingen,
419 Niedersachsen, Germany.
- 420 Qi, F.; Zhu, A. X. 2003. Knowledge discovery from soil maps using inductive learning.
421 *International Journal of Geographical Information Science* 17: 771-795.
- 422 R Development Core Team. 2011. *R: A Language and Environment for Statistical*
423 *Computing*, R Foundation for Statistical Computing, Vienna, Vienna, Austria.

- 424 Scull, P.; Franklin, J.; Chadwick, O.A. 2005. The application of classification tree analysis to
425 soil type prediction in a desert landscape. *Ecological Modelling* 181: 1-15.
- 426 ten Caten, A; Dalmolin, R.S.D.; Pedron, F.A.; Mendonça-Santos, M. de L. 2011. Estatística
427 multivariada aplicada à diminuição do número de preditores no mapeamento digital do solo.
428 *Pesquisa Agropecuária Brasileira* 46: 554-562.
- 429 ten Caten, A. ; Silva, C.A. ; Dalmolin, R. S. D. 2007. Mapas disponíveis e a demanda por
430 levantamentos de solos em grandes escalas. In: Congresso Brasileiro de Ciência do Solo. 31.
431 Anais. Gramado. CD-ROM.
- 432 Wilson, J. P. e Gallant, J. C. 2000. Digital terrain analysis. p. 1-27. *In: Wilson, J. P. e Gallant,*
433 *J. C. eds. Terrain analysis: principles and applications. Wiley & Sons, New York, NY, USA.*
- 434 Witten, I.H. e Frank, E. 2005. *Data Mining: Practical Machine Learning Tools and*
435 *Techniques. 2ed. Morgan Kaufmann, San Francisco, Califórnia, USA.*
- 436 Wahba, G. 1990. Spline models for Observational data. Paper presented at CBMS-NSF
437 Regional Conference Series in Applied Mathematics. Society for Industrial and Applied
438 Mathematics, Philadelphia, PA, USA.

DISCUSSÃO

Como esta tese foi estruturada em capítulos na forma de artigos, os quais constituem-se em si próprios em trabalhos completos, buscou-se nesta discussão traçar comentários gerais a respeito dos artigos, assim como a respeito do estado da arte do MDS no Brasil.

Em virtude do curto período de atividades, a pesquisa em MDS brasileira ainda não é resultado de grupos consolidados nessa linha de trabalho (Tabela 01 – Artigo 01). Disso decorre que, até o presente momento, pesquisadores têm se limitado a empregar estratégias automatizadas de mapeamento de solos como forma de se convencer e de convencer seus pares cientistas de que esta abordagem é capaz de gerar resultados similares àqueles alcançados pela metodologia convencional de mapeamento. A maioria das iniciativas de mapeamento de classes de solos tem se preocupado em avaliar a capacidade da abordagem automatizada em reproduzir o mapa que foi produzido pelo pedólogo, utilizando a abordagem convencional por fotointerpretação.

Como exemplo de mapeamento digital de classes de solos a partir de uma demanda, tem-se o estudo realizado por Silva (2011). O estudo detalha o mapeamento de classes de solos para o município de Erechim, a partir da necessidade de se espacializar as áreas aptas para o cultivo da erva mate. O autor utilizou de um mapa de classes de solos gerado a partir de modelos por Árvores de Decisão para gerar um mapa de aptidão agrícola das terras e de áreas propícias ao cultivo da erva mate. Segundo esse autor, a adequação dos usos da terra para o município de Erechim indicou a aptidão para atividades com lavouras, silvicultura e pecuária leiteira, com reduzido percentual de restrição de uso, que é reservado à preservação de fauna e flora.

Com os primeiros estudos publicados no país, a fase de ‘nova abordagem’ já pode ser tida como superada. Na pesquisa em MDS a ser realizada no país e especificamente para o mapeamento automatizado de classes de solos, três importantes frentes de estudos podem ser exploradas, a saber: (i) a identificação e o dimensionamento das demandas, otimizando os recursos disponíveis para gerar a informação espacial em solos; (ii) pontos metodológicos que influenciam a execução de mapeamentos de classes e propriedades de solos; (iii) a disponibilização das informações com a qualidade requerida e na forma que possa ser utilizada pelos usuários.

Analisando a identificação e dimensionamento de demandas, verifica-se que o MDS possui um problema em comum à abordagem convencional: a desinformação dos gestores e

tomadores de decisão. Esses desconhecem a aplicabilidade do mapa de solos para a melhoria da alocação e uso dos recursos naturais, econômicos e humanos. Por outro lado, a abordagem automatizada possui como vantagem o maior apelo visual e o potencial ilustrativo dos Sistemas de Informação Geográfica (Burrough et al., 1994). Assim, poderão auxiliar na divulgação e na demonstração de que mapas de classes de solos são ferramentas que qualificam a gestão e o planejamento (Bouma, 2005).

Se de um lado existe um grupo de usuários que até ignora a existência de mapas de solos, de outro existem usuários extremamente qualificados e ávidos pelo conhecimento da informação espacial em solos. Para pesquisadores de diversas áreas de conhecimento, a informação sobre o solo é fundamental para a execução de suas atividades científicas. Estudos relacionados à produção de sedimentos (Minella e Merten, 2011), impactos das mudanças do uso da terra no sequestro de carbono (Tornquist et al., 2009) e tomada de decisão em cenários de pouca disponibilidade de dados (Bacic et al., 2008) são alguns exemplos de linhas de pesquisa ávidas pelo conhecimento da distribuição espacial de propriedades e classes de solos. No entanto, na falta dessa informação, o que tem sido visto é que os estudos generalizam as informações existentes, muitas vezes em escalas inadequadas (Dalmolin et al., 2004). Também ocorre que pesquisadores de áreas não diretamente ligadas à ciência do solo buscam viabilizar as informações necessárias por conta própria, o que pode implicar em erros substanciais nas conclusões desses estudos. O pedólogo precisa viabilizar as informações necessárias à pesquisa na velocidade e qualidade demandadas, sob pena de ser dado como um profissional desnecessário em vários projetos.

Esta tese esteve focada na segunda frente de estudos sugerida ao MDS brasileiro, as quais são as questões metodológicas do mapeamento de classes de solos. Uma linha de trabalho muito importante para o MDS é o estudo da variabilidade espacial do solo, e nesse aspecto a escala e a resolução espacial refletem diretamente na capacidade de se compreender e estudar a variabilidade (Behrens et al., 2010). A análise Wavelet demonstra-se como uma abordagem promissora para o estudo e melhor entendimento do padrão de variabilidade de atributos de terreno. Trabalhos nessa linha devem explorar melhor não somente a relação entre as covariáveis preditoras e classes de solos, mas também como essa relação ocorre nas diferentes escalas. Como a contribuição de cada um dos fatores de formação do solo é distinta nas diversas escalas de estudo e representação da informação solo, assim também uma covariável preditora adequada para um estudo em uma escala 1:25.000 poderá não ser relevante para estudos em escala de 1:250.000. Sem a intenção de ser um trabalho conclusivo, o estudo realizado nesta tese buscou trazer e divulgar uma forma diferenciada de explorar a

variabilidade espacial. A análise Wavelet tem um enorme potencial para estudos de escalas espaço-temporais, especialmente onde um grande volume de dados esteja disponível e as relações sejam de complexo estabelecimento.

O grande potencial para coleta de dados ambientais faz necessário avaliar o melhor conjunto de amostras a ser utilizado para gerar o modelo preditivo (*instance selection*) (Schmidt et al., 2008). Na revisão de literatura, foi verificado que os estudos até aqui realizados no país não têm discriminado as amostras quanto ao seu potencial preditivo. Este estudo buscou, então, demonstrar que existem diferenças qualitativas entre as observações de acordo com a sua localização espacial relativa aos polígonos de solos definidos pelo conhecimento tácito do pedólogo. A estratégia de pré-processamento de dados apresentada é de simples implementação e possibilita que observações de covariáveis ambientais localizadas nas áreas de maior variabilidade espacial (bordas dos polígonos) sejam descartadas da fase de geração dos modelos.

Além da seleção das observações com maior potencial preditivo, a escolha entre as covariáveis que mais contribuem para a predição implica modelos mais eficientes e simplificados (*feature selection*) (Behrens et al., 2010). A Análise de Componentes Principais possibilita revelar as relações entre covariáveis predictoras e selecionar entre aquelas em que não ocorram efeitos de multicolinearidade (ten Caten et al., 2011a; ten Caten et al., 2011b).

A variabilidade da densidade de dados revelada pela revisão de literatura indicou que existe, ainda, falta de um número de referência para o volume de dados a ser manipulado na fase de ajuste dos modelos preditivos (Tabela 02 – Artigo 01). Este estudo indicou que valores entre cinco e 15% do universo total das observações (pixels) disponíveis seriam proporções de dados adequadas para o treinamento dos modelos preditivos. Sob a pena de, caso o valor seja menor, a complexidade das relações solo-paisagem não serem capturadas pelos modelos ou, caso a densidade de amostragem seja muito elevada, poder-se-á estar manipulando um excesso de dados. É importante ter em mente que esses percentuais dizem respeito a um universo total de 277.491 observações, totalizando 249,80 km², que estão distribuídos nos 889,33 km² daquele estudo. Esses valores de densidade de amostragem precisam ser mais bem avaliados nos casos em que a dimensão espacial do estudo não seja tão expressiva.

Ainda nessa linha de pesquisa, visando elucidar questões metodológicas para a implementação do MDS, existem muitas dúvidas a serem sanadas. Entre elas está a influência das diferentes fontes de Modelos Digitais de Elevação (MDE) na qualidade do fator relevo, e, conseqüentemente, nos atributos de terreno a serem gerados. O relevo tem sido um dos fatores

de formação mais utilizados pelo MDS (McBratney et al., 2003). A pesquisa precisa avaliar as novas tecnologias que irão possibilitar um ganho de resolução espacial nos atributos derivados do fator relevo. Isso ocorrerá especialmente em locais mais planos, onde o fator relevo não tem o mesmo poder preditivo do que nas áreas mais declivosas (Rivero et al., 2007).

A terceira frente de estudo, listada no segundo parágrafo dessa discussão, diz respeito à avaliação da qualidade da informação gerada. A abordagem mais comumente adotada tem sido a de comparar o mapa gerado com o mapa disponível pelas técnicas convencionais de levantamento. O índice kappa entre os dois mapas tem sido a forma mais frequente de indicar a qualidade do mapa gerado. A fidelidade da reprodução do mapa convencional pelo mapa gerado pelo MDS (índice kappa próximo de 100%) vem sendo utilizada como um indicativo de que a técnica de MDS é aceitável. Contudo, os erros inerentes ao mapa convencional não têm sido reportados, tão pouco a influência dos erros previamente existentes é estudada nos mapas gerados (Brus et al., 2011). Finalmente, a avaliação da qualidade do mapa gerado deverá ser testada a campo e reportada junto dele.

Quanto à forma de disponibilizar a informação, está claro que passa pelas tecnologias digitais (Grunwald, 2009). A produção de uma nova informação por vias automatizadas e sua posterior disponibilização somente pela impressão em meio analógico significa condenar essa informação ao mesmo destino de vários outros levantamentos de solos, hoje praticamente inacessíveis nas gavetas das instituições de ensino e pesquisa pelo país. Ao mesmo tempo em que é importante divulgar e tornar a nova informação acessível é importante que isso seja feito de maneira padronizada. Caso não sejam somados esforços em torno de uma normatização para a geração e a disponibilização dessa informação, corre-se o risco de, no futuro, existirem formatos de dados incompatíveis para a reunião em um mesmo projeto. O desafio de padronizar e uniformizar a forma de disponibilizar os dados não deve ser assumido por um único grupo de pesquisa, mas por vários pesquisadores, como os participantes da RedeMDS ou pela Divisão Especializada Solo no Espaço e no Tempo da Sociedade Brasileira de Ciência do Solo.

O presente estudo teve como foco o mapeamento de classes de solos, contudo o mapeamento de propriedades do solo deve ser foco maior de pesquisas. Devido ao caráter generalista dos mapas de classes de solos, pode estar no mapeamento de propriedades uma aproximação maior com o usuário final dessa informação. O mapeamento de propriedades como pH, textura, matéria orgânica e fertilidade pode possibilitar uma aplicação mais imediata da informação gerada, do que se comparado à informação presente em mapas de

classes de solos, as quais são por natureza generalizadas. A possibilidade inerente ao MDS de mapear a paisagem na forma de uma matriz regular de pontos (pixels) permite conhecer a informação solo de uma forma muito mais detalhada do que pode ser feito com os mapas de classes de solos (Zhu et al., 2001). Esse potencial do MDS precisa ser mais bem explorado no país, como também uma forma de divulgar mais a abordagem automatizada, tornando a informação sobre o solo mais clara e disponível ao usuário.

Decisiva para o futuro do MDS no país é a formação dos cientistas do solo, a qual necessariamente deverá abarcar uma nova gama de conhecimentos. Para Grunwald et al. (2011), além do conhecimento da teoria e prática da pedologia, o cientista do solo moderno deve (i) ser capaz de acessar e manipular as ferramentas geoespaciais, como o sensoriamento proximal e remoto do solo e de outras covariáveis ambientais, e (ii) ter fortes habilidades quantitativas, particularmente, em estatísticas e análise espacial.

REFERÊNCIAS BIBLIOGRÁFICAS

- BACIC, I.L.Z.; ROSSITER, D.G.; MANNAERTS, C.M. Applicability of a distributed watershed pollution model in a data-poor environment in Santa Catarina State, Brazil. **Revista Brasileira de Ciência do Solo**, Viçosa, v. 32, p. 1699-1712, aug. 2008.
- BEHRENS, T. et al. Multi-scale digital terrain analysis and feature selection for digital soil mapping. **Geoderma**, Amsterdam, v. 155, p. 175-185, mar. 2010.
- BOUMA, J. Soil scientists in a changing world. **Advances in Agronomy**, Newark, v. 88, p. 67-96, nov. 2005.
- BRUS, D.J.; KEMPEN, B.; HEUVELINK, G.B.M. Sampling for validation of digital soil maps. **European Journal of Soil Science**, v. 62, n. 3, p. 1365-2389, apr. 2011.
- BURROUGH, P.A.; BOUMA, J.; YATES, S.R. The state of the art in pedometrics. **Geoderma**, Amsterdam, v. 62, p. 311-326, 1994.
- DALMOLIN, R.S.D et al. Relação entre as características e o uso das informações de levantamentos de solos de diferentes escalas. **Ciência Rural**, Santa Maria, v. 34, n. 5, p. 1479-1486, set./out. 2004.
- GRUNWALD, S. Multi-criteria characterization of recent digital soil mapping and modeling approaches. **Geoderma**, Amsterdam, v. 152, n. 3-4, p. 195-207, jul. 2009.
- MCBRATNEY, A. B.; MENDONCA SANTOS, M. L.; MINASNY, B. On digital soil mapping. **Geoderma**, Amsterdam, v. 117, n. 1-2, p. 3-52, jun. 2003.
- MINELLA, J.P.G.; MERTEN, G.H. Monitoramento de bacias hidrográficas para identificar fontes de sedimentos em suspensão. **Ciência Rural**, Santa Maria, v. 41, n. 3, p. 424-432, mar. 2011.
- RIVERO, R.G.; S. GRUNWALD, S.; BRULAND, G.L. Incorporation of spectral data into multivariate geostatistical models to map soil phosphorus variability in a Florida wetland. **Geoderma**, Amsterdam, v. 140, p. 428-443, jun. 2007
- SCHMIDT, K; BEHRENS, T.; SCHOLTEN, T. Instance selection and classification tree analysis for large spatial datasets in digital soil mapping. **Geoderma**, Amsterdam, v. 146, p. 138-146, jun. 2008.
- SILVA, C. A. da **Zoneamento Pedoclimático da erva mate *Ilex paraguariensis* para o município de Erechim-RS**. 2011. 160f. Tese (Doutorado em Ciência do Solo)-Universidade Federal de Santa Maria, Santa Maria, 2011.

ten CATEN, A. et al. Componentes principais como preditores no mapeamento digital de classes de solos. **Ciência Rural**, Santa Maria, v. 41, n. 7, p. 1170-1176, jul. 2011a.

ten CATEN, A. et al. Estatística multivariada aplicada à diminuição do número de preditores no mapeamento digital do solo. **Pesquisa Agropecuária Brasileira**, Brasília, v. 46, n. 5, p. 554-562, maio 2011b.

TORNQUIST, C.G. et al. Spatially explicit simulations of soil C dynamics in Southern Brazil: Integrating century and GIS with i_Century. **Geoderma**, Amsterdam, v. 150, p. 404-414, apr. 2009.

ZHU, A.X. et al. Soil mapping using GIS, expert knowledge, and fuzzy logic. **Soil Science Society of America Journal**, Madison, v. 65, p. 1463-1472, set./oct. 2001.