

Estudo do potencial hídrico em cafeeiros utilizando técnicas de Aprendizado de Máquina

Pedro Henrique Nunes
Departamento de Automática
Universidade Federal de Lavras
Lavras, Brasil
phnunes95@gmail.com

Danton Diego Ferreira
Departamento de Automática
Universidade Federal de Lavras
Lavras, Brasil
danton@ufla.com.br

Eduardo Vilela Pierangeli
Departamento de Automática
Universidade Federal de Lavras
Lavras, Brasil
evp.pierangeli@gmail.com

Margarete Marin Lordelo Volpato
Empresa de Pesquisa Agropecuária de Minas Gerais
EPAMIG
Lavras, Brasil
margarete@epamig.ufla.br

Vânia Aparecida Silva
Empresa de Pesquisa Agropecuária de Minas Gerais
EPAMIG
Lavras, Brasil
vania.silva@epamig.ufla.br

Resumo—O potencial hídrico é um dos principais indicadores utilizados no estudo das relações hídricas nas plantas e mostra o nível de hidratação do tecido das mesmas. É medido de forma direta por um equipamento denominado bomba de Scholander, entretanto, este processo é complexo e demorado. Existem na literatura diversas variáveis numéricas capazes de descrever algumas propriedades das plantas por meio dos índices de refletância das folhas, que apresentam relações diretas e indiretas com o potencial hídrico. Este trabalho tem como objetivo explorar tais variáveis para estudar o potencial hídrico em cafeeiros através do uso de ferramentas de inteligência computacional e reconhecimento de padrões, utilizando amostras espectrais de duas lavouras de café, na região de Santo Antônio do Amparo e Diamantina, ambas localizadas no estado de Minas Gerais, nos anos de 2014 a 2018. Dessa forma, foram projetadas e treinadas, utilizando o conjunto de ferramentas *NNTool* do software *MatLab*, uma rede neural do tipo MLP (*Multi-Layer Perceptron*) para estimar o potencial hídrico e outra para fazer a classificação das amostras de acordo com o seu status hídrico, baseado em variáveis espectrais. Foram desenvolvidos ainda um classificador e um estimador baseados em árvore de decisão utilizando a ferramenta *Decision and Regression Tree*, também do software *MatLab*. Os resultados mostraram que as redes neurais artificiais foram superiores como estimador, com um índice de confiança médio de 0,8550, porém as árvores de decisão obtiveram um desempenho ligeiramente superior como classificador, com uma acurácia global de 88,8% e um coeficiente de *Kappa* de 70,07%.

Index Terms—Redes Neurais Artificiais. Árvores de Decisão. Potencial Hídrico Foliar. Cafeeiros.

I. INTRODUÇÃO

O uso racional de recursos hídricos no setor agrícola, visando a economia de água e uma boa produtividade, constitui um aspecto de grande relevância para a produção sustentável, tornando cada vez mais necessária a busca por tecnologias para o auxílio nessa área.

O conhecimento das relações hídricas das plantas é fundamental para o estabelecimento de boas práticas agrícolas. Exis-

tem diversas variáveis capazes de descrever o *status* hídrico da planta, dentre elas destaca-se o potencial hídrico foliar (ψ_{am}). Este índice é mensurado de forma direta por meio de um equipamento denominado bomba de Scholander, no qual amostras de folhas recolhidas das plantas são submetidas a diferentes níveis de pressão, determinando assim o valor de ψ_{am} [1].

O processo de medição direta do potencial hídrico foliar, entretanto, apresenta algumas dificuldades de execução como elevado tempo de medição e risco de explosão, além de ser um método destrutivo, que pode causar danos a planta. Dessa forma, muitos estudos têm sido realizados visando sua determinação de forma indireta, principalmente com base nos índices espectrais relativos à refletância foliar das plantas alvos [2].

A obtenção dos índices espectrais pode ser realizada de diversas formas, destacando-se o sensoriamento remoto e os espectrômetros de campo. Esses métodos funcionam por meio da medição do percentual de luz refletida pelas folhas das plantas para os diferentes comprimentos de onda da luz, e com base nesses valores são calculados os índices utilizando equações já existentes na literatura.

Os índices de refletância foliar fornecem informações diversas sobre a saúde das plantas, tais como vivacidade, amarelidão, vermelhidão, entre outras características que são interpretadas por fisiologistas a fim de obter informações relevantes sobre a amostra em questão. Esses índices possuem correlação com o *status* hídrico da planta, o que do ponto de vista de inteligência computacional, pode ser usado para estimar o potencial hídrico de forma indireta por meio de técnicas de natureza estocástica [3].

Dentre os algoritmos estocásticos conhecidos, destacam-se as redes neurais artificiais (RNAs), que possuem um rico embasamento teórico na literatura e são utilizadas em diversas aplicações de reconhecimento de padrões, tais como classificação, predição e estimação. Devido a sua alta capaci-

dade de generalização, as RNAs quando bem projetadas, são capazes de mapear o comportamento de funções, mesmo que estas sejam não-lineares e/ou descontínuas.

As redes neurais artificiais podem ser projetadas em plataformas de desenvolvimento como *Matlab*, *Weka*, *Scilab*, entre outras, sendo as maiores dificuldades do projeto a determinação da quantidade de camadas e de neurônios, bem como a seleção das melhores características a serem aplicadas como entradas para a RNA. Outro ponto importante para a escolha das redes neurais artificiais como metodologia, é que uma vez projetadas e treinadas elas podem ser facilmente implementadas em softwares mais genéricos como o *Excel* por exemplo e, geralmente, apresentam rápida execução na fase operacional [4].

Outra ferramenta aplicável para a solução deste problema são as árvores de decisão, que são capazes de realizar a classificação de amostras em grupos, e a estimação de seu comportamento de acordo com os valores de suas características. Essa ferramenta pode ser usada para classificar o potencial hídrico dos cafeeiros, identificando se seu *status* hídrico está adequado, eliminando a necessidade do uso da bomba de Scholander. A motivação principal em se usar árvores de decisão está na sua boa capacidade de representação de dados. Esta técnica permite mapear dados, mesmo que de alta dimensão, de forma simples e de fácil entendimento, o que é bastante útil em sistemas de reconhecimento de padrões [5].

Desta forma, o presente documento apresenta, o desenvolvimento de duas ferramentas baseadas, em redes neurais artificiais e árvores de decisão, capazes de estimar o potencial hídrico em cafeeiros (*Coffea arabica* L.), bem como duas ferramentas capazes de classificá-lo. Por fim, é apresentado um estudo comparativo entre as metodologias abordadas.

Recentemente, o uso de índices espectrais em aprendizado de máquina vem se tornando um assunto recorrente em pesquisas. [14] desenvolveu uma estrutura de regressão envolvendo vários modelos de aprendizado de máquina para estimar parâmetros de água com base em dados espectrais. Ainda nessa linha, [15] implementou técnicas de aprendizado de máquina para prever e monitorar as condições de seca devido às mudanças climáticas utilizando-se de índices espectrais. Por fim, o uso dos índices espectrais em técnicas de aprendizado de máquinas se mostram úteis como um método eficiente e prático de mapeamento da cobertura de determinada vegetação, como mostra o trabalho de [16]. Assim, o campo de estudo desse trabalho vem se mostrando cada vez mais em evidência.

II. METODOLOGIA

Para a realização da metodologia adotada, foram necessários os seguintes equipamentos:

- 1 Mini-espectrômetro foliar CI-710;
- 1 Bomba de Scholander SEC-3115-P40G4V;
- 1 Licença do Software MatLab.

A. Banco de Dados

Utilizou-se neste projeto um banco de dados espectrais montado pela equipe de pesquisadores da EPAMIG e EM-BRAPA/Café. As amostras foram coletadas no período de 2014-2018, em cafeeiros arábica nos municípios de Santo Antônio do Amparo e Diamantina, Minas Gerais, Brasil.

O conjunto de dados é composto por 1280 eventos com 9 características, sendo elas PRI (*Photochemical Reflectance Index*), PSRI (*Plant Senescence Reflectance Index*), NDVI (*Normalized Difference Vegetation Index*), WBU (*Water Band Index*), ARI1 (*Anthocyanin Reflectance Index*), CRI1 (*Carotenoid Reflectance Index*), SIPI (*Structure Insensitive Pigment Index*), FRI (*Flavonol Reflectance Index*) e data da coleta, obtidas por meio do mini-espectrômetro, e o potencial hídrico (ψ_{am}), medido com uma bomba de Scholander. Os dados foram coletados em onze diferentes datas, visando capturar o efeito de variações climáticas sazonais da região.

A fim de possibilitar o processamento da data de coleta, foi feita sua conversão para um valor numérico, de acordo com a estação do ano, sendo atribuído um valor de 1 a 4 para cada estação: primavera, verão, outono e inverno, respectivamente.

Para a estimativa dos dados, considerou-se os valores puros normalizados do potencial hídrico, entretanto para fins de classificação os alvos foram divididos em classes, de acordo com o valor do potencial hídrico, intervalados a cada -0,9429MPa (e.g. classe1 [0 -0,9429], classe2 [-0,9429 - 1,8857], e assim por diante), totalizando 7 classes. O número de classes foi determinado de acordo com a disposição dos dados da faixa de valores de potencial hídrico disponíveis.

A divisão de classes resultou na seguinte disposição dos dados: 906 dados na Classe 01; 173 dados na Classe 02; 99 dados na Classe 03; 38 dados na Classe 04; 33 dados na Classe 05; 10 dados na Classe 06 e 23 dados na Classe 07.

B. Pré-processamento

Para o início do desenvolvimento foi necessário um pré-processamento dos dados. Para o uso em *MatLab*, primeiramente alocaram-se os dados em forma matricial, onde as linhas são as amostras e as colunas são as características, sendo os alvos, dispostos em uma matriz coluna separada.

Para o uso no projeto dos classificadores, os alvos foram separados em 7 classes, igualmente intervaladas de acordo com o valor máximo e mínimo de potencial hídrico das amostras. Essa configuração de classes foi escrita na forma matricial, na qual cada linha representa uma classe e cada coluna uma amostra, sendo atribuído o valor "1" para a classe à qual a amostra pertence e "0" para as demais, onde as linhas representam as sete classes, respectivamente. A motivação da utilização de classes é a possibilidade de observar qual faixa de dados de potencial hídricos os classificadores implementados obterão uma maior quantidade de acertos.

Em seguida foram eliminados manualmente os *outliers* e normalizadas as amostras, de modo a evitar eventuais priorizações de dados considerando apenas seu valor absoluto mais elevado em relação aos demais [13]. Esses dados eliminados dizem respeito às linhas com dados faltantes ou com

valor discrepante do esperado, e foram retirados de acordo com especialistas da área.

Foi adotada como regra de normalização a escala para o conjunto de [0,1] para todos os dados. A equação (1) foi utilizada para a realização desse procedimento.

$$P_n = \frac{(P - P_{min})}{(P_{max}) - (P_{min})} \quad (1)$$

Onde, P_n é o valor normalizado, P é o valor original, P_{min} é o menor valor e P_{max} é o maior valor.

Uma vez normalizados os dados, foi feito um embaraalhamento, a fim de evitar que amostras muito similares sejam utilizadas no treinamento das redes ou definição das árvores, prejudicando sua capacidade de generalização [4].

Por fim, para a seleção de características, desenvolveu-se um código em *MatLab* a fim de realizar uma análise de correlação das variáveis, visando identificar redundâncias e a relevância de cada característica, por meio do coeficiente de *Pearson* de cada atributo com ele próprio e com o alvo, eliminando aqueles que apresentaram baixa correlação com o alvo (baixa relevância) e/ou alta correlação entre si (alta redundância), no segundo caso um dos atributos de cada par com alta correlação foi eliminado, diminuindo a dimensionalidade do problema [6].

Os dados foram divididos em 70% para treinamento, 20% para validação e 10% para teste, sendo esta divisão embasada pela sugestão de 70% para treinamento e 30% para testes, proposta por [13], finalizando o pré-processamento.

C. Métricas de Avaliação de Desempenho

1) *Índices de desempenho para os classificadores*: Visando mensurar o desempenho dos classificadores, foi determinado o erro quadrático médio em cada interação, sua média e desvio padrão nos testes, bem como seu percentual de acerto, pior classe e percentual de acerto na pior classe.

Visando a avaliação de desempenho dos classificadores, foram montadas suas matrizes de confusão, otimizando a realização dos estudos comparativos dos resultados obtidos com o uso de cada técnica, além de facilitar a identificação dos pontos de maior dificuldade de classificação. Além disso também foi utilizado como parâmetro o índice de *Kappa*, proposto por [7], e utilizado por [8] e por [9] para avaliar seus trabalhos.

O índice *Kappa* é obtido por meio da análise da matriz de confusão e é determinado seguindo a Equação (2) e segundo [8], é uma medida satisfatória para avaliar a precisão de uma classificação, pois considera, todos os elementos da matriz de confusão, diferente da acurácia global que utiliza somente a diagonal principal.

$$K = \frac{(P_o - P_e)}{1 - P_e} \quad (2)$$

onde,

$$P_o = \sum \frac{n_{ii}}{n} \quad (3)$$

e

$$P_e = \sum \frac{(n_i \cdot n_j)}{n^2} \quad (4)$$

Sendo que, K significa o coeficiente de *Kappa*, n_{ii} representa a quantidade de elementos classificados corretamente, n representa o total de amostras, n_i representa o somatório da i -ésima linha e n_j o somatório da j -ésima coluna, considerando $i = j$. Uma vez determinado o coeficiente de *Kappa*, os resultados podem ser classificados segundo a Tabela I, fornecendo uma medida de desempenho do sistema.

Tabela I
CLASSIFICAÇÃO DOS ÍNDICES DE DESEMPENHO

Valor de <i>Kappa</i>	Desempenho
0,81 a 1,00	Excelente
0,61 a 0,80	Muito Bom
0,41 a 0,60	Bom
0,21 a 0,40	Razoável
0,00 a 0,20	Ruim
< 0,0	Péssimo

2) *Índices de desempenho para os estimadores*: A fim de medir o desempenho dos estimadores, inicialmente foram determinados o erro relativo médio, sua média e desvio padrão para as iterações dos testes, além do índice de confiança médio, melhor índice de confiança e desvio padrão.

A análise de desempenho dos estimadores foi feita por meio do cálculo dos índices de correlação “ r ”, que representa a precisão do modelo, concordância “ d ”, que representa a exatidão do modelo e confiança “ id ”, que representa a confiabilidade do modelo.

O índice de correlação foi obtido através do coeficiente de *Pearson*, segundo a equação (5).

$$\rho = \frac{cov(X, Y)}{\sqrt{var(X) \cdot var(Y)}} \quad (5)$$

onde X e Y são os vetores a serem correlacionados e ρ é o coeficiente de correlação de *Pearson*.

O índice de concordância é dado pela Equação (6), onde P_i é o valor estimado, O_i é o valor observado e O é a média dos valores observados [10].

$$d = 1 - \frac{\sum (P_i - O_i)^2}{\sum (|P_i - O| + |O_i - O|)^2} \quad (6)$$

O índice de confiança foi obtido por meio do produto da correlação pela concordância e pode ser usado como um indicador do desempenho de estimadores. [11] propuseram a classificação de desempenho de estimadores com base na Tabela II, utilizada também por [10] e por [12] como metodologia de avaliação dos resultados de seus trabalhos.

III. RESULTADOS E DISCUSSÕES

A. Seleção de Características

Com a realização dos procedimentos propostos, foi montada uma tabela mostrada na Figura 1, que ilustra a matriz de correlação entre cada atributo com os demais e com o alvo

Tabela II
CLASSIFICAÇÃO DOS ÍNDICES DE DESEMPENHO

Valor de "id"	Desempenho
> 0,85	Ótimo
0,76 a 0,85	Muito Bom
0,66 a 0,75	Bom
0,61 a 0,65	Mediano
0,51 a 0,60	Sofrível
0,41 a 0,50	Mau
< 0,40	Péssimo

(ψ_{am}). Primeiramente foram selecionados os índices com maiores valores de correlação com o alvo, ou seja, com alta relevância (células marcadas em azul). Em seguida foram selecionados os maiores valores de correlação desses índices com os demais (células marcadas em vermelho), visando identificar a existência de redundância.

De posse desses dados, foram eliminados os índices PSRI, ARII e FRI, que apresentaram, simultaneamente, alta redundância e baixa relevância. Também foi removido o índice WBU que, apesar de não apresentar redundância com as demais características, apresentou uma baixa correlação com o alvo.

Desta forma, os índices selecionados foram PRI, NDVI, CRII, SIPI e Data.

	Ψ_{am}	PRI	PSRI	NDVI	WBU	ARII	CRII	SIPI	FRI
Ψ_{am}	1								
PRI	0,2735	1							
PSRI	-0,0424	-0,1885	1						
NDVI	0,1988	0,4898	-0,2365	1					
WBU	0,038	-0,0268	0,0248	0,0671	1				
ARII	-0,0401	0,3383	0,0325	0,5653	-0,1125	1			
CRII	0,2045	0,4180	0,0373	0,6091	-0,1763	0,5329	1		
SIPI	0,2721	0,2987	-0,0793	0,6017	0,0757	0,0671	0,4572	1	
FRI	0,1026	0,0546	0,0240	0,2476	-0,1037	0,0952	0,5487	0,5823	1
Data	0,2175	0,0187	-0,0885	0,0199	-0,3762	0,0096	0,1281	0,0864	0,0795

Figura 1. Matriz de Correlação.

B. Classificadores

Nesta seção serão mostrados e discutidos os resultados obtidos para os dois classificadores desenvolvidos. Para a exposição dos mesmos, será utilizada a matriz de confusão, que pode ser interpretada da seguinte forma:

- Os valores em negrito nas células vermelhas e verdes representam a quantidade de amostras, e o percentual em cada uma dessas células representa a parcela que esses dados representam do conjunto de amostras total.
- As células em verde representam os dados classificados corretamente em cada classe;
- As células em vermelho representam os dados classificados de forma errada pelo classificador;
- As células em cinza na última coluna apresentam, em verde, o somatório dos acertos e em vermelho o somatório dos erros, ambos relativos à linha em que se enquadram;
- As células em cinza na última linha apresentam, em verde, o somatório dos acertos e em vermelho o somatório dos erros, ambos relativos à coluna em que se enquadram;

- A célula em azul mostra, em verde, o percentual de acerto total (acurácia global) do classificador, e em vermelho o seu percentual de erro total.

1) *Redes Neurais Artificiais*: A fim de determinar o número de neurônios em cada camada da rede neural artificial, foram realizadas 10 execuções para cada arquitetura, e para determinar as funções de ativação foram feitos testes com 100 iterações.

Dessa forma, montou-se a matriz de confusão (Figura 2), utilizando-se do banco de dados de teste, contendo os resultados do melhor modelo testado. Utilizou-se 2 camadas intermediárias, sendo 13 neurônios na primeira camada e 7 na segunda camada. Em ambas, a função de ativação utilizada foi a Tangente Hiperbólica.

Na análise da matriz de confusão, é possível observar que a maioria dos dados são pertencentes as classes 1, 2, 3 e 7, justificando o melhor desempenho do classificador nesses eventos. Uma solução que pode vir a melhorar o desempenho da RNA para as classes 4, 5 e 6, que possuem uma menor quantidade de ocorrências, é a coleta de mais dados das classes com poucas amostras, almejando obter um maior equilíbrio entre as classes. Além disso, uma outra alternativa seria o reequilíbrio da quantidade de amostras utilizando técnicas como a interpolação de dados.

		Matriz de Confusão							
		1	2	3	4	5	6	7	
Classes Preditas	1	9543 74.6%	521 4.1%	199 1.6%	144 1.1%	95 0.7%	2 0.0%	0 0.0%	90.9% 9.1%
	2	17 0.1%	745 5.8%	3 0.0%	12 0.1%	2 0.0%	1 0.0%	0 0.0%	95.5% 4.5%
	3	28 0.2%	21 0.2%	634 5.0%	61 0.5%	32 0.3%	2 0.0%	4 0.0%	81.1% 18.9%
	4	33 0.3%	9 0.1%	48 0.4%	51 0.4%	67 0.5%	7 0.1%	0 0.0%	23.7% 76.3%
	5	33 0.3%	10 0.1%	11 0.1%	58 0.5%	74 0.6%	30 0.2%	2 0.0%	33.9% 66.1%
	6	1 0.0%	2 0.0%	1 0.0%	7 0.1%	16 0.1%	14 0.1%	11 0.1%	26.9% 73.1%
	7	0 0.0%	0 0.0%	0 0.0%	0 0.0%	14 0.1%	46 0.4%	189 1.5%	75.9% 24.1%
		98.8% 1.2%	57.0% 43.0%	70.8% 29.2%	15.3% 84.7%	24.7% 75.3%	13.7% 86.3%	91.7% 8.3%	87.9% 12.1%
		1	2	3	4	5	6	7	Classes Alvo

Figura 2. Matriz de confusão do melhor resultado com 100 iterações para a RNA. Os valores em negrito representam a quantidade de amostras, e o percentual em cada uma dessas células representa a parcela desses dados no conjunto de amostras total. Fonte: Do Autor (2020)

Considerando a matriz de confusão do melhor resultado (Figura 2), calculou-se o seu índice de *Kappa*, segundo a Equação (2). Este apresentou um valor de 67.21%, descrevendo o classificador como "muito bom" segundo a Tabela I. A média do erro quadrático nessa configuração foi de 0.1203, com desvio padrão de ± 0.0523 .

2) *Árvores de Decisão*: Para a construção da árvore de decisão, foram executadas 100 iterações para cada configuração de poda. Estes valores foram maiores que os da RNA devido ao menor tempo de execução e a menor quantidade de testes necessários, uma vez que o único parâmetro a ser determinado é o nível de poda da árvore. O melhor modelo da técnica de árvores de decisão possui um nível de poda de 85 nós.

Matriz de Confusão

	1	2	3	4	5	6	7	
1	9424 73.6%	496 3.9%	169 1.3%	124 1.0%	78 0.6%	0 0.0%	0 0.0%	91.6% 8.4%
2	45 0.4%	838 6.5%	13 0.1%	0 0.0%	4 0.0%	0 0.0%	0 0.0%	93.1% 6.9%
3	20 0.2%	23 0.2%	648 5.1%	26 0.2%	25 0.2%	0 0.0%	0 0.0%	87.3% 12.7%
4	28 0.2%	2 0.0%	44 0.3%	117 0.9%	98 0.8%	8 0.1%	0 0.0%	39.4% 60.6%
5	31 0.2%	8 0.1%	24 0.2%	79 0.6%	90 0.7%	17 0.1%	3 0.0%	35.7% 64.3%
6	0 0.0%	0 0.0%	0 0.0%	0 0.0%	10 0.1%	34 0.3%	13 0.1%	59.6% 40.4%
7	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	46 0.4%	215 1.7%	82.4% 17.6%
	98.7% 1.3%	61.3% 38.7%	72.2% 27.8%	33.8% 66.2%	29.5% 70.5%	32.4% 67.6%	93.1% 6.9%	88.8% 11.2%
	1	2	3	4	5	6	7	

Classes Preditas (eixo Y) e **Classes Alvo** (eixo X)

Figura 3. Matriz de confusão com 100 iterações para a Árvore de Decisão. Os valores em negrito representam a quantidade de amostras, e o percentual em cada uma dessas células representa a parcela desses dados no conjunto de amostras total. Fonte: Do Autor (2020).

Considerando ainda a matriz de confusão do melhor resultado mostrada na Figura 3, foi calculado o seu índice de *Kappa*, segundo a equação (2), apresentando um valor de 71,07%, descrevendo o classificador como "muito bom" segundo a Tabela I, colocando-o na mesma categoria do classificador baseado em redes neurais artificiais, apesar de seu coeficiente de *Kappa* mais elevado. A média do erro quadrático nessa configuração foi de 0.1113, com desvio padrão de ± 0.0255 .

C. Estimadores

1) *Redes Neurais Artificiais*: A fim de dimensionar o número de neurônios nas camadas da RNA utilizada para regressão, bem como determinar as melhores funções de ativação, foi realizado um procedimento similar ao utilizado para a rede de classificação.

A fim de melhorar a exibição dos resultados obtidos e possibilitar uma melhor análise dos pontos de erro, foram registradas as Figuras 4 e 5. A primeira imagem ilustra a distribuição dos dados (pontos) em relação a reta ideal, sendo que quanto mais distantes da reta, maiores os erros associados

aos pontos em questão. A segunda mostra um comparativo entre as curvas real (pontilhado azul) e estimada (vermelha).

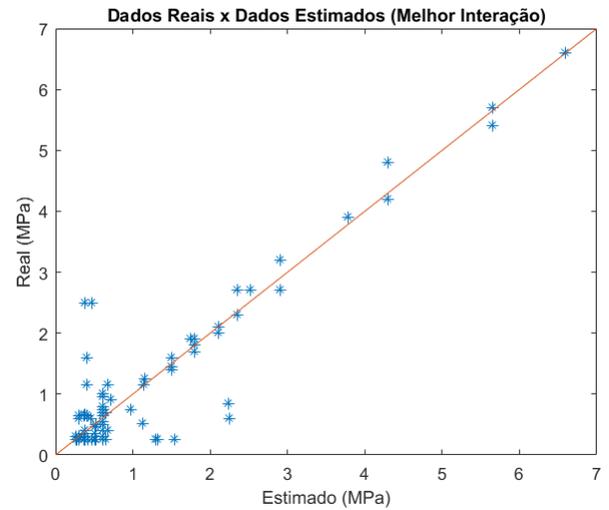


Figura 4. Dispersão dos dados estimados em uma interação da RNA. Fonte: Do Autor (2020).

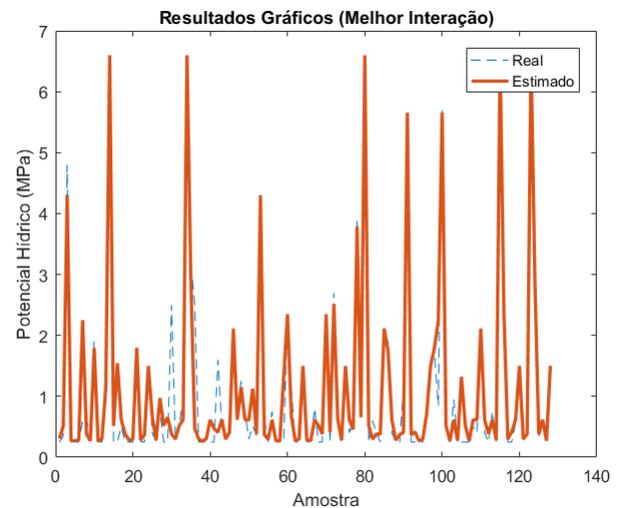


Figura 5. Comparativo da curva real pela curva estimada em uma interação da RNA. Fonte: Do Autor (2020).

Como exposto nas figuras, os maiores erros estão, em geral, associados aos valores de potencial hídrico intermediários, correspondentes às classes 4, 5 e 6, corroborando com os resultados obtidos pelos classificadores, elucidando o impacto do pequeno número de amostras pertencentes essas classes.

No contexto geral o desempenho do estimador pode ser considerado satisfatório, uma vez que foi classificado como "ótimo" pelo critério de [11], com um baixo desvio padrão em 100 iterações, demonstrando a constância dos resultados, além do erro relativo médio de $7,27 \pm 4,44\%$, reforçando a integridade dos resultados.

2) *Árvores de Decisão*: Para a construção da árvore de decisão de regressão, foram executadas 100 iterações para cada

configuração de poda. Novamente, o número de iterações foi maior que o da RNA devido ao menor tempo de execução e a menor quantidade de testes necessários, uma vez que o único parâmetro a ser determinado é o nível de poda da árvore.

A fim de melhorar a exibição dos resultados obtidos e possibilitar uma melhor análise dos pontos de erro, foram registradas as Figuras 6 e 7. A primeira imagem ilustra a distribuição dos dados (pontos) em relação a reta ideal, sendo que quanto mais distantes da reta, maiores os erros associados aos pontos em questão. A segunda mostra um comparativo entre as curvas real (pontilhado azul) e estimada (vermelha).

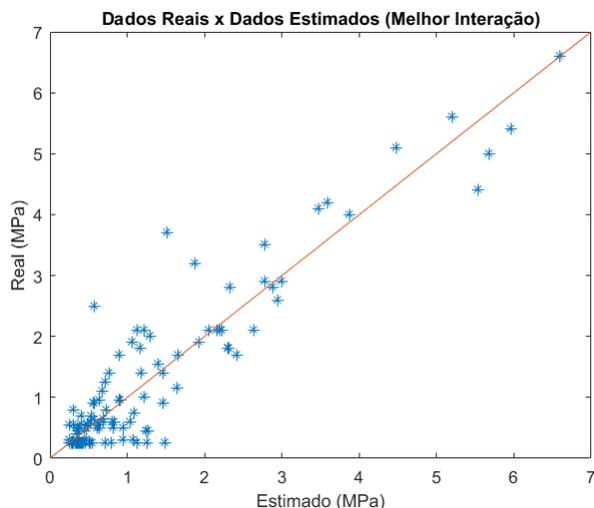


Figura 6. Dispersão dos dados estimados em uma interação da árvore. Fonte: Do Autor (2020).

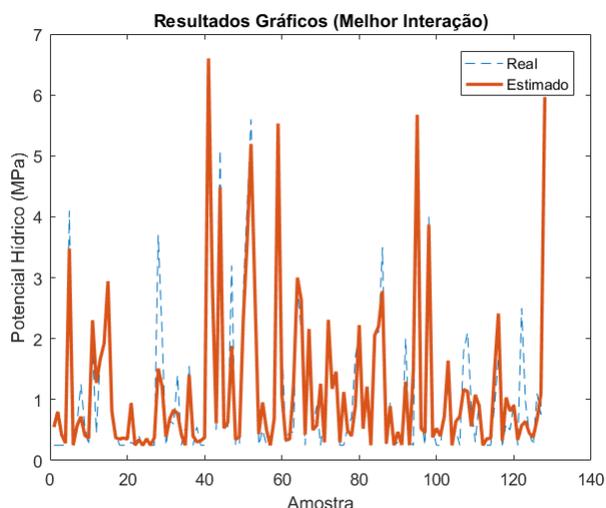


Figura 7. Comparativo da curva real pela curva estimada em uma interação da árvore. Fonte: Do Autor (2020).

Corroborando com os resultados do estimador baseado em RNA, os maiores erros estão, em geral, associados aos valores de potencial hídrico intermediários, correspondentes às classes

4, 5 e 6, elucidando o impacto do pequeno número de amostras pertencentes à essas classes. Entretanto, no estimador baseado em árvore de decisão, os erros associados aos eventos fora dessa faixa são maiores, o que pode ser observado pela maior dispersão desses dados em relação à reta ideal (Figura 6), se comparado ao estimador utilizando redes neurais artificiais (Figura 4).

No contexto geral, o desempenho do estimador utilizando árvores de decisão, assim como o estimador baseado em RNA, pode ser considerado satisfatório, uma vez que foi classificado como “ótimo” pelo critério de [11], com um baixo desvio padrão em 100 iterações (0,0909), demonstrando a constância dos resultados, além do erro relativo médio de $3,85 \pm 4,27$ %, reforçando a integridade dos resultados.

D. Estudo Comparativo dos Resultados

A fim de definir qual a melhor técnica, foram montadas as Tabelas III e IV, que expõem os melhores resultados obtidos na regressão e classificação respectivamente.

Tabela III
COMPARAÇÃO DOS MELHORES RESULTADOS PARA REGRESSÃO.

Técnica Utilizada	ERM	σ_{ER}	σ_{id}	id	id _{max}
Rede Neural Artificial	7,27	4,44	0,0739	0,8550	0,9544
Árvore de decisão	3,85	4,27	0,0909	0,8299	0,9443

Por meio da análise da Tabela III, é possível notar que a árvore de decisão obteve menores valores de erro relativo médio e desvio do erro relativo, e pode ser justificado pelo fato de que elas retornam apenas valores fixos, sua estimativa é feita por meio de uma classificação em muitas classes, evitando que estimem valores pontuais muito discrepantes dos reais.

Entretanto, a rede neural artificial apresentou um índice de confiança médio superior à árvore de decisão, bem como um menor desvio deste valor, elucidando sua maior constância. Por fim a rede obteve um índice de confiança maior em sua melhor interação se comparado ao obtido nas mesmas circunstâncias pela árvore.

Estes resultados demonstram uma maior capacidade de generalização da RNA, em se tratando do problema de regressão em pauta, elencando-a como a melhor ferramenta para este cenário, quando comparada à árvore de decisão, uma vez que o principal critério de avaliação de desempenho utilizado foi o índice de confiança.

Tabela IV
COMPARAÇÃO DOS MELHORES RESULTADOS PARA CLASSIFICAÇÃO.

Técnica	EQM	σ_{EQM}	↓EQM	Acerto	P.C.	P.C.
RNA	0,1203	0,0523	0,0055	87,9%	4	23,7%
A.D.	0,1113	0,0255	0,0547	88,8%	5	35,7%

Com base no exposto na Tabela IV, é possível inferir que a árvore de decisão foi superior à rede neural em quase todos os aspectos, exceto no menor erro quadrático médio de uma

única interação, o qual para a RNA foi de 0,0055, ao passo que para a árvore foi de 0,0547. Isto possivelmente é justificado por singularidades das amostras utilizadas nos testes desta execução em específico, não refletindo necessariamente a superioridade da rede, uma vez que, conforme já mencionado, seu desempenho foi inferior nos quesitos de valor médio e no desvio padrão do EQM.

Outro ponto considerado ao analisar estes resultados é o percentual de acerto na pior classe, como pode ser visto, este valor foi substancialmente maior para as árvores, evidenciando sua maior capacidade de generalização para este caso, corroborando com os demais resultados expostos na tabela.

Considerando o percentual de acerto, a média do EQM nas iterações e o desvio padrão do EQM, as diferenças não foram significativas, porém a árvore continuou ligeiramente superior também nestes quesitos. Por fim, ao observar o seu coeficiente de *Kappa*, nota-se que a árvore de decisão se manteve superior, justificando sua escolha como a melhor metodologia para classificação dentre as abordadas neste documento.

Desta forma, fica evidente a superioridade da rede neural desenvolvida para estimação bem como a superioridade da árvore de decisão desenvolvida para a classificação, permitindo a aplicação de ambas para o problema discutido neste texto, de acordo com a abordagem à ser utilizada (classificação ou regressão).

É válido considerar ainda, que as árvores de decisão facilitam o entendimento das operações que levam aos resultados, uma vez que são compostas de regras “Se, Então”, ao passo que as redes neurais artificiais mesmo podendo ser facilmente replicadas utilizando uma sequência de somas e multiplicações após serem treinadas, são consideradas, como um sistema em caixa preta, no qual não se tem conhecimento sobre as razões que levaram aos resultados. Essas características devem ser consideradas ao escolher a metodologia à ser utilizada, uma vez que os resultados, em geral, não apresentaram grandes diferenças de desempenho.

IV. CONCLUSÃO

Os resultados obtidos com a realização das metodologias apontadas, de forma geral, foram positivos, uma vez que ambas as técnicas abordadas conseguiram realizar as atividades propostas com desempenho satisfatório.

Contudo, observou-se que para valores de potencial hídrico na faixa intermediária das amostras coletadas (aproximadamente entre -2,75MPa e -5,4MPa), tanto os classificadores quanto os regressores, utilizando ambas as técnicas, apresentaram resultados abaixo do esperado. Isso provavelmente se deve a pequena quantidade de eventos com o valor do ψ_{am} nesta faixa em comparação ao número de eventos com o valor do potencial hídrico nas demais faixas.

Dentre as ferramentas de aprendizado de máquina estudadas neste texto, as redes neurais artificiais foram melhores para a estimação do potencial hídrico foliar em cafeeiros das regiões abordadas, e as árvores de decisão foram superiores quando o objetivo era a classificação das amostras.

Entretanto, em ambas as aplicações o desempenho não foi substancialmente distinto, sendo que para o problema de regressão a RNA obteve um índice de confiança de 0,8550, contra 0,8290 da árvore de decisão, ambos classificados como “ótimo” segundo o critério proposto por [11]. Já para o problema de classificação o percentual de acerto foi de 88,8% para a árvore de decisão e de 87,9% para a RNA, com índices *Kappa* de 70,07% e 67,21%, respectivamente, ambos classificados como “muito bom” segundo exposto por [8].

Observou-se que a separação de dados em classes ocasiona em um desbalanceamento da base de dados. Esse problema pode ser tratado fazendo-se uso de algoritmos de *oversampling* ou *undersampling*, que pretende-se aplicar em trabalhos futuros, bem como a técnica de *k-fold*, na separação de dados de treino e teste.

Desta forma é válido que sejam considerados outros aspectos característicos de cada metodologia ao decidir sobre a sua aplicação, como por exemplo, facilidade de implementação, custo computacional e entendibilidade, que não foram objetos de estudo deste trabalho.

AGRADECIMENTOS

Agradeço principalmente aos pesquisadores do Programa de Pós-Graduação da Universidade Federal de Lavras e aos pesquisadores da EPAMIG, onde todos, em equipe, foram essenciais para a realização desse projeto, com o apoio em conhecimento e fomento.

REFERÊNCIAS

- [1] J.D. Barnes and L. Balaguer and E. Manrique and S. Elvira and A.W. Davison, "A reappraisal of the use of DMSO for the extraction and determination of chlorophylls a and b in lichens and higher plants". *Environmental and Experimental Botany*. Volume 32, number 2. 85 - 100 pag. 1992.
- [2] ZHANG, C. and Pattey, E. and Liu, J. and Cai, H. and Shang, J. and T. Dong. "Retrieving Leaf and Canopy Water Content of Winter Wheat Using Vegetation Water Indices". *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*. 2018.
- [3] GENC, L. AND INALPULAT, M. AND KIZIL, U. AND MERIK, M. AND SMITH, S. E. AND MENDES, M. "Determination of water stress with spectral reflectance on sweet corn (*Zea mays* L.) using classification tree (CT) analysis". *Photochemistry and Photobiology*. Volume 100. 81-90 pag. 2013.
- [4] HAYKIN, S. "Redes Neurais: Princípios e Prática". Bookman. Porto Alegre. 2001.
- [5] SIMÕES, A. C. A. "Mineração de Dados Baseada em Árvore de Decisão para Análise do Perfil de Contribuintes". Universidade Federal de Pernambuco. 140 pag. 2008.
- [6] SILVA, I. N. da and SPATTI, D. H. and FLAUZINO, R. A. "Redes Neurais Artificiais para engenharia e ciências aplicadas". Artliber. São Paulo. 2010.
- [7] COHEN, J. "A COEFFICIENT OF AGREEMENT FOR NOMINAL SCALES". *Educational and Psychological Measurement*. Volume 20. 37-46 pag. Nova York. 1960.
- [8] GONCALVES, W. G. AND RIBEIRO, H. M. C. AND SÁ, J. A. S. AND MORALES, G. P. AND FILHO, H. R. F. AND ALMEIDA, A. C. "Classificação de estratos florestais utilizando redes neurais artificiais e dados de sensoriamento remoto". *Revista Ambiente Água*. Taubaté. 2016.
- [9] SANTOS, C. J. "Avaliação do uso de classificadores para verificação de atendimento a critérios de seleção em programas sociais". Universidade Federal de Juiz de Fora. 88 pag. Juiz de Fora. 2017.

- [10] SOARES, F. C. AND ROBAINA, A. D. AND PEITER, M. X. AND RUSSI, J. L. AND VIVIAN, G. A. "Redes neurais artificiais na estimativa da retenção de água do solo". Revista Ciência Rural. Volume 44. 293-300 pag. 2014.
- [11] CAMARGO, A.P. and SENTELHAS, P.C. "Avaliação do desempenho de diferentes métodos de estimativa da evapotranspiração potencial no Estado de São Paulo". Revista Brasileira de Agrometeorologia. Volume 5. 89-97 pag. 1997.
- [12] BATISTA, L. A. and GUIMARÃES, R. J. and PEREIRA, F. J. and CARVALHO, G. R. and CASTRO, E. M. de. "Anatomia foliar e potencial hídrico na tolerância de cultivares de café ao estresse hídrico". Revista Cia Agronômica. Volume 41. 475-481 pag. 2010.
- [13] BRAGA, A. P. de. and LEON, A. P. and LUDEMIR, T. B. "Redes Neurais Artificiais: Teoria e Aplicações". LTC. 2ª edição. Belo Horizonte. 2007.
- [14] Maier, Philipp M. e Keller, Sina. "Machine Learning Regression on Hyperspectral Data to Estimate Multiple Water Parameters". 2018 9º Workshop sobre Imagem Hiperespectral e Processamento de Sinal: Evolução em Sensoriamento Remoto (Sussurros). 2018.
- [15] Perera, Sachi e Li, Wenzhao e Linstead, Erik e El-Askary, Hesham. "Forecasting Vegetation Health in the MENA Region by Predicting Vegetation Indicators with Machine Learning Models". IGARSS 2020 - 2020 IEEE International Geoscience and Remote Sensing Symposium. 2020.
- [16] Chafik, Hassan e Berrada, Mohamed e Legdou, Anass e Amine, Aouatif e Lahssini, Said. Exploração de índices espectrais NDVI, NDWI amp; SAVI no modelo de classificador Random Forest para mapeamento de cobertura fraca de alecrim: aplicação na região de Gourrama, Marrocos. 2020 IEEE International conference of Moroccan Geomatics (Morgeo). 2020.