

**DESENVOLVIMENTO E SELEÇÃO DE MODELOS DE ALERTA PARA A FERRUGEM DO CAFEIEIRO EM ANOS DE ALTA CARGA PENDENTE DE FRUTOS<sup>1</sup>**

Cesare Di Girolamo Neto<sup>2</sup>, Luiz Henrique Antunes Rodrigues<sup>3</sup>, Thiago Toshiyuki Thamada<sup>4</sup>, Carlos Alberto Alves Meira<sup>5</sup>

<sup>1</sup>Trabalho financiado pelo Consórcio Brasileiro de Pesquisa e Desenvolvimento do Café – Consórcio Pesquisa Café

<sup>2</sup>Bolsista Consórcio Pesquisa Café, Me., cesare.neto@gmail.com.

<sup>3</sup>Professor, PhD, Faculdade de Engenharia Agrícola - Universidade Estadual de Campinas, Campinas-SP, lique@feagri.unicamp.br.

<sup>4</sup>Bolsista Consórcio Pesquisa Café, Bel., thiago.tamada@colaborador.embrapa.br.

<sup>5</sup>Pesquisador, Dr., Embrapa Informática, Campinas-SP, carlos.meira@embrapa.br.

**RESUMO:** Este trabalho teve como objetivo desenvolver e selecionar modelos de alerta para prever o aumento da taxa de progresso mensal da ferrugem do cafeeiro para lavouras em anos de alta carga pendente de frutos. Os modelos foram desenvolvidos por meio de quatro técnicas de mineração de dados: redes neurais artificiais, árvores de decisão, máquinas de vetores suporte e florestas aleatórias. A seleção dos modelos ocorreu de forma gráfica e por meio de suas medidas de desempenho e o resultado mostrou que os modelos desenvolvidos neste trabalho apresentaram desempenho superior a outros previamente desenvolvidos. Estes modelos de alerta fornecem melhores subsídios para o monitoramento da doença da ferrugem do cafeeiro em anos de alta carga pendente de frutos.

**PALAVRAS CHAVE:** modelos de alerta, doenças de plantas, mineração de dados, modelagem.

**DEVELOPMENT AND SELECTION OF WARNING MODELS FOR COFFEE RUST ON HIGH FRUIT LOAD YEARS**

**ABSTRACT:** The aim of this study was to develop and select warning models to predict the development of the monthly progress rate for coffee rust on high fruit load years. The models were developed by four data mining techniques: neural networks, decision trees, support vector machines and random forest. The models were selected by graphs and performance measures. The results showed that the models developed on this study obtained better performance than others previously developed. These warning models provide better insights on dealing with the coffee rust on years of high fruit load.

**KEYWORDS:** warning models, plant disease, data mining, modeling.

**INTRODUÇÃO**

A ferrugem (causada pelo fungo *Hemileia vastatrix* Berk. e Br.) é a principal doença do cafeeiro. De acordo com Zambolim et al. (2002), as perdas de produção devido à esta doença podem chegar a 50%, caso nenhuma medida de controle seja adotada. O controle da ferrugem pode ser feito de forma eficiente com aplicação de fungicidas, entretanto, métodos tradicionais de controle, como a aplicação baseada em calendário fixo, podem levar ao uso desnecessário e em períodos incorretos, gerando gastos excessivos para o produtor, na compra e mão de obra para aplicação dos fungicidas, além de causar impactos ambientais.

Ferramentas como modelos de predição, ou alerta, podem ser utilizadas visando antecipar quando uma doença de planta ocorrerá, sendo que quando uma predição for realizada corretamente pode-se evitar aplicações incorretas de defensivos. A partir da divulgação no meio científico da Mineração de Dados (*Data Mining* - DM) por Fayyad et al. (1996), tem-se notado um aumento no número de modelos gerados por meio de tal metodologia. Neste sentido, modelos de alerta para determinar a taxa de progresso mensal da ferrugem do cafeeiro foram desenvolvidos por Meira et al. (2009) e Cintra et al. (2011), os quais desenvolveram Árvores de Decisão (AD) e AD *fuzzy*, respectivamente. Outras técnicas na área de DM também têm sido utilizadas para a modelagem da ferrugem do cafeeiro, como Redes Neurais Artificiais (RNA) *fuzzy*, desenvolvidas por Alves et al. (2010); Máquinas de Vetores Suporte (*Support Vector Machines* - SVM), induzidas por Luaces et al. (2011) e Florestas Aleatórias (*Random Forest* - RF), que também foram mencionadas no trabalho de Cintra et al. (2011).

A escolha pelo uso de uma ou outra técnica de modelagem requer a análise do problema em questão, sendo que cada uma destas técnicas têm vantagens e desvantagens. Os modelos em AD são fáceis de serem interpretados (Witten et al., 2011), a medida que as RF, além de evitar o sobreajuste do modelo aos dados utilizados, também são pouco sensíveis a ruídos (Breiman, 2001). A principal vantagem das RNA é resolver problemas que apresentam uma solução muito difícil de ser encontrada (Haykin, 2009), enquanto que as SVM também evitam sobreajuste e normalmente produzem classificadores muito precisos (Witten et al., 2011).

Sendo assim, a geração de novos modelos de alerta da ferrugem do cafeeiro por meio de diferentes técnicas de modelagem poderia gerar modelos com poder de predição maior do que os gerados anteriormente. Logo, o objetivo

deste trabalho foi desenvolver, avaliar e selecionar modelos de alerta para determinar a taxa de progresso mensal da ferrugem do cafeeiro (TP) para lavouras em anos de alta carga pendente de frutos.

## MATERIAL E MÉTODOS

Os dados utilizados neste trabalho referem-se ao acompanhamento mensal da incidência da ferrugem do cafeeiro em uma fazenda experimental da Fundação PROCAFÉ. Esta fazenda está localizada na cidade de Varginha/MG, latitude sul de 21° 34' 00'', longitude oeste de 45° 24' 22'' e altitude de 940 m. Foram coletados dados para períodos entre outubro de 1998 a outubro de 2011. As lavouras selecionadas tinham idade entre 6 e 20 anos e estavam em anos de alta carga pendente de frutos (acima de 30 sacas beneficiadas/ha), sendo para os cultivares Catuaí (Vermelho e Amarelo) e Mundo Novo. Havia lavouras em espaçamento largo (por volta de 3,5 m entre linhas e 0,7 m entre plantas – densidade média de 4.000 plantas/ha) e em espaçamento adensado (por volta de 2,5 m entre linhas e 0,5 m entre plantas – densidade média de 8.000 plantas/ha). O processo de amostragem foi realizado ao final de cada mês, conforme recomendação de Chalfoun (1997): coleta de 100 folhas do terço médio das plantas em cada talhão, entre o terceiro e o quarto par de folhas; contagem do número de folhas com lesões de ferrugem; e determinação da incidência (percentual de folhas atacadas). Não houve controle da doença durante o ano agrícola. O período de colheita foi entre junho e agosto.

Além dos dados referentes à doença, foram obtidos dados meteorológicos, como temperatura (média, máxima e mínima), precipitação pluviométrica, umidade relativa do ar (UR) e velocidade do vento. Estes dados foram registrados a cada 30 minutos por uma estação meteorológica automática presente na fazenda experimental.

A variável dependente, ou também chamada de atributo meta, foi definida como sendo a TP, a qual consiste no aumento, diminuição ou manutenção da incidência da doença entre dois meses subsequentes. As diferenças percentuais das incidências foram mapeadas em um atributo de origem binária, sendo que a classe '1' indica que a TP foi maior ou igual a 5 p.p. (pontos percentuais) e a classe '0' indica que TP foi inferior a 5 p.p. Uma segunda opção de atributo meta também foi utilizada, com limites de 10 p.p. Os valores do atributo meta foram baseados em Meira et al. (2008).

Os atributos preditivos, ou variáveis independentes, partiram de um nível de construção horário (forma em que foram coletados da estação meteorológica) e passaram por transformações, levando-os até um nível que possibilitou uma análise em conjunto com o atributo meta (Meira et al., 2008). Diversos atributos foram calculados utilizando-se de médias e somatórios simples das variáveis meteorológicas. Atributos relacionados com molhamento foliar foram derivados da variável de UR, utilizando valores superiores a 90% para considerar o molhamento das folhas. Os atributos preditivos foram completados com o atributo que dava a característica do espaçamento de cada lavoura utilizada. Foram construídos um total de 23 atributos, além da opção de atributo meta.

O conjunto de dados utilizado na modelagem totalizou 738 registros, sendo que alguns registros foram eliminados devido à falha na estação meteorológica. As duas opções de atributo meta foram utilizadas separadamente no conjunto de dados. O conjunto de dados com a opção de atributo meta 10 p.p. continha cerca de 25% dos registros com classe "1" e 75% com classe "0". Esta distribuição prejudicaria o processo de aprendizado do modelo, logo os registros foram equilibrados seguindo métodos de balanceamento de classes (Batista et al., 2004). Estes métodos deixaram cada classe com cerca de 50% dos registros. O conjunto de dados não balanceado também foi utilizado na indução de modelos. Os modelos gerados com arquivos balanceados tiveram seu desempenho avaliado nos conjuntos originais.

A partir destes conjuntos foi realizada a seleção de atributos, que visou escolher quais eram os melhores e mais importantes atributos do conjunto de dados. Esta seleção foi realizada de duas formas: uma delas subjetiva, a qual consistiu na seleção de atributos de acordo com a complexidade e dificuldade de obtenção dos mesmos, gerando três conjuntos de dados (Meira et al., 2008); e a outra por meio de cinco métodos objetivos (algoritmo já implementado) amplamente utilizados na área de DM: Wrapper, CFS, Chi-quadrado, InfoGain e GainRatio (Witten et al., 2011). Os métodos de seleção foram aplicados ao conjunto contendo todos os atributos.

O software utilizado na indução dos modelos foi o WEKA, versão 3.7.9. (Hall et al., 2009). Foram utilizadas quatro técnicas de modelagem para induzir os modelos: AD, RNA, RF e SVM. Todos os parâmetros de modelagem foram calibrados visando otimizar a taxa de acerto dos modelos e evitar o sobreajuste. Para a opção de atributo meta 5 p.p. foram induzidos 32 modelos (8 métodos de seleção de atributos x 4 técnicas de modelagem), já para opção de atributo meta 10 p.p. foram induzidos 64 modelos (arquivos balanceados e não balanceados - 2 x 8 métodos de seleção de atributos x 4 técnicas de modelagem).

A avaliação e seleção dos modelos desenvolvidos ocorreu em duas etapas. Os modelos foram inicialmente dispostos em gráficos do tipo ROC (*Receiver Operating Characteristic*), onde tornou-se possível selecionar os melhores modelos dentre todos os desenvolvidos. Estes modelos encontram-se dispostos em vértices de um "envelope externo convexo", presente nos gráficos. Os modelos que não fizeram parte do envelope convexo puderam ser descartados (Provost & Fawcett, 2001). A partir dos modelos selecionados, eles foram avaliados por meio de suas medidas de desempenho e atributos no conjunto de dados.

As medidas de desempenho foram provenientes de uma matriz contendo os acertos e erros do modelo, chamada matriz de confusão, onde as principais medidas são a taxa de acerto, sensibilidade e especificidade. A taxa de acerto determina o percentual de acertos total para um modelo, já a sensibilidade trata o percentual de exemplos de aumento da TP que foram classificados corretamente, enquanto que a especificidade mostra o percentual de exemplos de não aumento da TP que foram classificados corretamente. Um modelo foi considerado superior a outro por obter medidas mais elevadas

da taxa de acerto, sensibilidade e especificidade, além do fato de que as duas últimas devem estar equilibradas, evitando que o modelo acerte muito de uma classe e erre muito da outra.

O conjunto de atributos utilizado para gerar cada um dos modelos também foi utilizado para sua avaliação. Conjuntos com muitos atributos e, principalmente, atributos complexos e difíceis de serem calculados, podem dificultar a aplicabilidade destes modelos. Dentre dois modelos com medidas de desempenho similares, optou-se pelo mais simples, justamente pelo fato deste ser mais fácil de ser implantado. Em contrapartida, analisou-se qual conjunto de dados continha mais atributos relacionado ao desenvolvimento da ferrugem em campo. Quanto mais condições relacionadas à ferrugem, mais representativo foi considerado o modelo.

## RESULTADOS E DISCUSSÃO

A Figura 1 representa o gráfico ROC para os modelos desenvolvidos com atributo meta 10 p.p. Neste caso, nota-se que seis modelos foram selecionados no envelope convexo (15, 26, 52, 59, 61 e 62/64). As medidas de desempenho destes modelos estão na Tabela 1.

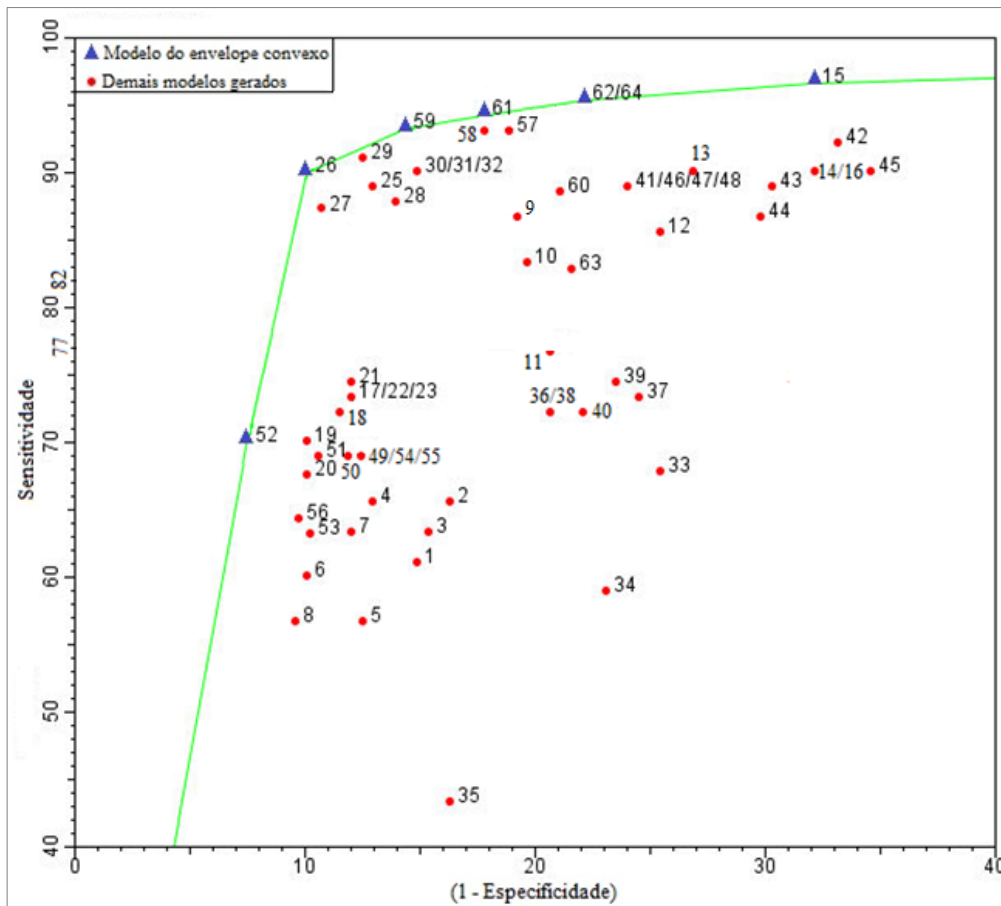


Fig.1. Gráfico ROC para os modelos gerados com opção de atributo meta 10 p.p.

Ao analisar os modelos da Tabela 1, notou-se que o modelo 15 obteve uma taxa de acerto inferior aos demais e foi o primeiro a ser descartado. Os modelos 52, 61 e 62/64 apresentaram grandes variações nos valores de sensibilidade e especificidade, mostrando que pendem muito para acertar, percentualmente, mais exemplos de aumento da TP (61 e 62/64) ou mais exemplos de não aumento da TP (52), e também foram descartados.

Tabela 1: Medidas de desempenho (em %) para os modelos selecionados com opção de atributo meta 10 p.p.

Modelos	15	26	52	59	61	62/64
Técnica de modelagem	AD	RF	SVM	SVM	SVM	SVM
Taxa de acerto	76,5	89,9	85,3	88,6	86,0	83,5
Erro	23,5	10,1	14,7	11,4	14,0	16,5
Sensitividade	96,7	90,0	70,1	93,3	94,3	95,4
Especificidade	67,8	89,9	92,4	86,5	82,2	77,8

Com relação aos dois modelos restantes (26 e 59), seus conjuntos de atributos foram avaliados. O modelo 59 conteve 11 atributos, sendo três de temperaturas médias (mínima, máxima e média), um atributo de UR, dois atributos de precipitação (acumulada e média), um atributo relacionado ao número de dias chuvosos e três atributos de temperatura durante o período de incubação do *H. vastatrix*. Já o modelo 26 conteve todos estes atributos, além de outros dois atributos relacionados à duração do período de molhamento foliar (número de horas diárias e número de horas noturnas com UR acima de 90%) e mais um atributo relacionado à temperatura média durante este período.

Em termos de complexidade, os conjuntos estão equivalentes, com uma ligeira vantagem para o modelo 59, que contém três atributos a menos. Entretanto, o modelo 26 tem atributos de extrema importância para o desenvolvimento da ferrugem do cafeeiro. Para que o fungo se desenvolva são necessários períodos mínimos de seis horas de molhamento foliar contínuo (Kushalappa et al., 1983), sendo que caso esta situação ocorra com uma temperatura durante este período variando de 22 a 24 °C, as condições ótimas de desenvolvimento do fungo são atingidas (Zambolim et al., 2002). Por representar essas condições com seus atributos e também obter medidas de avaliação levemente superiores, o modelo 26 foi considerado o indicado para representar o aumento da ferrugem quando o limite para a TP for de 10 p.p. Tratando-se agora dos modelos gerados com opção de atributo meta de 5 p.p., a Figura 2 representa o gráfico ROC para estes modelos. Neste caso, nota-se que dois modelos foram selecionados no envelope convexo (22 e 28). As medidas de desempenho destes modelos estão na Tabela 2.

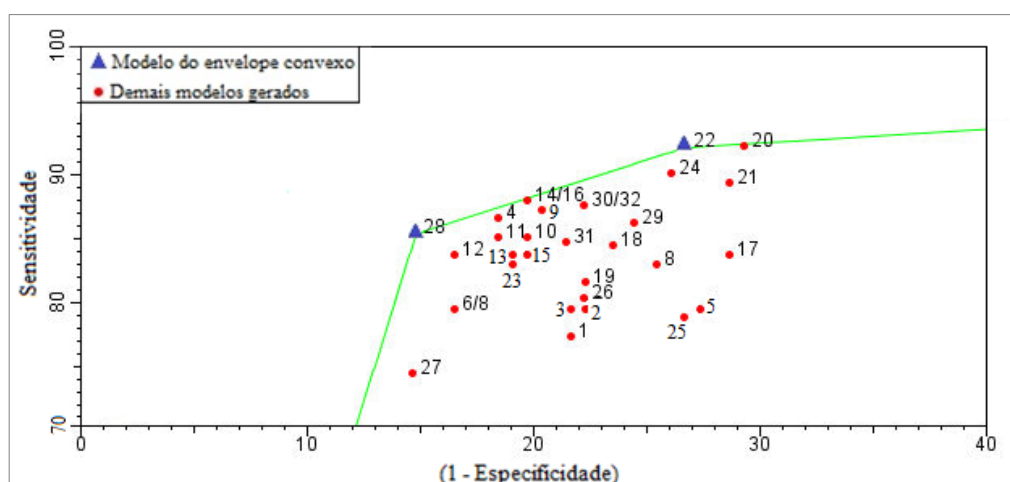


Fig.2. Gráfico ROC para os modelos gerados com opção de atributo meta 5 p.p.

Tabela 2: Medidas de desempenho (em %) para os modelos selecionados com opção de atributo meta 5 p.p.

Modelos	22	28
Técnica de modelagem	RNA	SVM
Taxa de acerto	82,2	85,3
Erro	17,8	14,7
Sensitividade	92,2	85,4
Especificidade	73,3	85,2

Ao analisar estes dois modelos, nota-se que o modelo 28 obteve medidas de avaliação superiores ao modelo 22, com melhores valores de taxa de acerto e medidas de sensibilidade e especificidade mais equilibradas. Isto dá indícios de que o modelo 28 teve melhor desempenho e deveria ser o escolhido para representar o aumento da TP neste caso.

De qualquer forma, os conjuntos de atributos para estes dois modelos foram analisados. O modelo 28 conteve quatro atributos, um relacionado à precipitação (média), um à velocidade do vento, um à temperatura durante o molhamento foliar e um atributo especial (atributo complexo, que reúne condições de molhamento foliar, temperatura e luminosidade relacionados ao desenvolvimento da ferrugem). Já o modelo 22 conteve mais atributos, um total de seis, sendo que dois estavam relacionados às temperaturas médias (mínima e média) além de outros quatro atributos especiais (Meira et al., 2008).

O conjunto de dados para gerar o modelo 22 é mais complexo, pois além de um número maior de atributos, ele contém mais atributos difíceis de serem gerados (atributos especiais). Ainda há o fato de que o modelo 28 contém atributos que representam mais as condições relacionadas ao desenvolvimento da doença, sendo que os atributos de precipitação e velocidade do vento contribuem para fatores relacionados à dispersão dos esporos do *H. vastatrix* na lavoura (Zambolim et al., 2002), sem contar o fato da precipitação média contribuir para que ocorram ou não períodos de molhamento foliar. A temperatura durante este período também é importante, como mencionado para os modelos com opção de

atributo meta 10 p.p. Sendo assim, com um conjunto de dados menos complexo, mais condições relacionadas ao desenvolvimento da ferrugem e melhores medidas de desempenho, o modelo 28 foi considerado o mais indicado. Os modelos gerados neste trabalho foram comparados com modelos gerados com metodologias similares, sendo estes gerados por Meira et al. (2009) e Cintra et al. (2011), conforme a Tabela 3. Todos estes modelos foram desenvolvidos para lavouras em anos de alta carga pendente de frutos.

Tabela 3: Medidas de desempenho (em %) dos modelos de alerta desenvolvidos neste trabalho e por outros autores.

Nome do modelo	Atributo meta	Taxa de acerto	Erro	Sensitividade	Especificidade
28	5 p.p.	85,3	14,7	85,4	85,2
MA-TI5*	5 p.p.	81,3	18,7	79,9	82,6
MIG5**	5 p.p.	84,7	15,3	***	***
26	10 p.p.	89,9	10,1	90,0	89,9
MA-TI10*	10 p.p.	79,2	20,8	70,3	82,8
M3G10**	10 p.p.	83,5	16,5	***	***

\* Modelo gerado por Meira et al., (2009). \*\* Modelo gerado por Cintra et al., (2011). \*\*\* Cintra et al., (2011) não forneceu estas medidas de avaliação.

Para a opção de atributo meta 5 p.p., verificou-se que a taxa de acerto do modelo 28 foi a melhor dentre os três modelos avaliados e seus valores de sensibilidade e especificidade foram superiores ao modelo MA-TI5. As medidas de sensibilidade e especificidade também estão levemente mais próximas entre si no modelo 28 do que no modelo MA-TI5. Assim, o modelo 28 apresentou desempenho superior aos demais, mostrando ser o mais adequado para realizar a predição do aumento da taxa de progresso da ferrugem para tal opção de atributo meta.

Com relação à opção de atributo meta 10 p.p., nota-se que a taxa de acerto do modelo desenvolvido neste trabalho (26) foi a maior dentre os três modelos avaliados, chegando próximo aos 90%. Já seus valores de sensibilidade e especificidade foram relativamente superiores ao modelo MA-TI10, principalmente para a sensibilidade, a qual foi de 90,0% contra 70,3%, mostrando uma melhora significativa do modelo 26 ao classificar corretamente exemplos positivos ou de aumento da TP. Este modelo também se mostrou mais equilibrado do que o modelo MA-TI10, pela diferença de apenas 0,1 p.p. entre as medidas de sensibilidade e especificidade. Consequentemente, o modelo 26 se mostrou superior aos demais para esta opção de atributo meta.

Os modelos desenvolvidos neste trabalho estão sendo avaliados para safras agrícolas dos anos de 2011/2012 e 2012/2013. Após a avaliação de quais modelos, ou conjuntos de modelos, obtiverem melhor desempenho, eles serão incorporados em um sistema de alerta. Este sistema está disponível na internet para uso exclusivo pelos técnicos da fundação PROCAFÉ, para auxiliar na elaboração das recomendações veiculadas nos boletins mensais de avisos fitossanitários. Quanto mais confiáveis forem os avisos de aumento da doença emitidos pelo sistema de alerta, mais confiáveis serão os avisos emitidos pela fundação, auxiliando, ainda mais, o produtor sobre a época correta de aplicar os fungicidas para efetuar o controle da ferrugem do cafeeiro.

## CONCLUSÕES

Os modelos de predição da taxa de progresso mensal da ferrugem do cafeeiro desenvolvidos neste trabalho fornecem melhores subsídios para o monitoramento da doença da ferrugem do cafeeiro em anos de alta carga pendente de frutos do que outros modelos desenvolvidos. Um sistema de monitoramento da ferrugem do cafeeiro, que esteja baseado nestes modelos de alerta, pode trazer ganhos econômicos ao produtor, além de reduzir os impactos ambientais causados por aplicações incorretas de fungicidas.

## AGRADECIMENTOS

À fundação PROCAFÉ por ceder os dados relacionados ao monitoramento de incidência da ferrugem do cafeeiro.

## REFERÊNCIAS BIBLIOGRÁFICAS

- ALVES, M. C.; CARVALHO, L. G.; POZZA, E. A.; ALVES, L. S. A Soft Computing Approach For Epidemiological Studies of Coffee And Soybean Rusts. *International Journal of Digital Content Technology and its Applications* 4:149-154. (2010).
- BATISTA, G. E. A. P. A.; PRATI, R. C.; MONARD, M. C. A study of the behavior of several methods for balancing machine learning training data. *SIGKDD Explorations* 6:20-29. (2004).
- BREIMAN, L. Random forests. *Machine Learning Journal* 45:5-32. (2001).
- CHALFOUN, S. M. Doenças do cafeeiro: importância, identificação e métodos de controle. Lavras: UFLA/FAEPE. (1997).

- CINTRA, M. E.; MEIRA, C. A. A.; MONARD, M. C.; CAMARGO, H. A.; RODRIGUES, L. H. A. The use of fuzzy decision trees for coffee rust warning in Brazilian crops. In: 11<sup>th</sup> International Conference on Intelligent Systems Design and Applications, 2011, Córdoba – ES. 11<sup>th</sup> International Conference on Intelligent Systems Design and Applications. Córdoba - ES:IEEE, 2011.
- FAYYAD, U.; PIATETSKY-SHAPIO, G.; SMYTH, P. From data mining to knowledge discovery in databases. *AI Magazine* 17(3):37-54. (1996).
- HALL, M. A.; FRANK, E.; HOLMES, G.; PFAHRINGER, B.; REUTEMANN, P.; WITTEN, I. H. The WEKA Data Mining Software: An Update; *SIGKDD Explorations* 11:10-18. (2009).
- HAYKIN, S. *Neural Networks and Learning Machines*. 3ed., Englewood Cliffs:Prentice-Hall. (2009).
- KUSHALAPPA, A. C.; AKUTSU, M.; LUDWIG, A. Application of survival ratio for monocyclic process of *Hemileia vastatrix* in predicting coffee rust infection rates. *Phytopathology* 73:96-103. (1983).
- LUACES, O.; RODRIGUES, L. H. A.; MEIRA, C. A. A.; BAHAMONDE, A. Using nondeterministic learners to alert on coffee rust disease. *Expert systems with applications* 38(11):14276-14283. (2011).
- MEIRA, C. A. A.; RODRIGUES, L. H. A.; MORAES, S. A. Análise da epidemia da ferrugem do cafeeiro com árvore de decisão. *Tropical Plant Pathology* 33(2):114-124. (2008).
- MEIRA, C. A. A.; RODRIGUES, L. H. A.; MORAES, S. A. Modelos de alerta para o controle da ferrugem-do-cafeeiro em lavouras com alta carga pendente. *Pesquisa Agropecuária Brasileira* 44(3):233–242. (2009).
- PROVOST, F.; FAWCETT, T.; KOHAVI, J. The case against accuracy estimation for comparing induction algorithms. In: 15<sup>th</sup> International Conference on Machine Learning, 1998, San Francisco – EUA. 15<sup>th</sup> International Conference on Machine Learning. San Francisco - EUA:Morgan Kaufmann, 1998.
- WITTEN, I. H.; FRANK, E.; HALL, M. A. *Data mining: practical machine learning tools and techniques*. 3ed. San Francisco:Morgan Kaufmann. (2011).
- ZAMBOLIM, L.; VALE, F. X. R.; COSTA, H.; PEREIRA, A. A.; CHAVES, G. M. Epidemiologia e controle integrado da ferrugem do cafeeiro. In: ZAMBOLIM, L. *O estado da arte de tecnologias na produção de café*. Viçosa: Suprema Gráfica e Editora, 369-449. (2002).