

Avaliação de uma Abordagem de Aprendizado Supervisionado para Operações no Mercado de Ações

Bárbara B. C. Silva¹, Marcelo S. Lauretto¹,
Pablo M. Andrade², Luciano V. Araújo¹

¹Escola de Artes, Ciências e Humanidades – Universidade de São Paulo
(EACH-USP) – São Paulo – SP – Brasil

²Instituto de Matemática e Estatística – Universidade de São Paulo
(IME-USP) – São Paulo – SP – Brasil

{marcelolauretto, barbara.barbosa.silva, pablo.andrade, lvaraujo}@usp.br

Resumo. *Técnicas de mineração de dados têm sido extensivamente aplicadas a mercados financeiros, para previsão de tendências, recomendação de operações de compra/venda e negociação automática. Neste trabalho, apresentamos uma proposta de aplicação de aprendizado supervisionado, utilizando Random Forests, para a recomendação de operações no mercado de ações. Dentre as 68 ações que compõem o Índice Bovespa, em 40 ações o método obteve uma taxa de aplicações bem sucedidas superior a 75%; em 38 ações a taxa de aproveitamento de oportunidades foi superior a 50%; e em 30 ações, o retorno líquido médio por operação foi superior a 4%. Esses resultados preliminares são bastante promissores, motivando estudos adicionais e extensões futuras.*

Abstract. *Data mining methods have been widely applied in financial markets, for trend forecasting, buy/sell decision support and automatic trading. This paper describes the application of a supervised learning approach, using Random Forests, for decision support in stock markets. Among the 68 stocks that compose the Bovespa Index, in 40 stocks our approach achieved a rate of successful operations above 75%, in 38 stocks, the rate of yielded opportunities was higher than 50%; and in 30 stocks the average net return per operation was higher than 4%. These preliminary results provide a strong motivation for further studies and extensions.*

1. Introdução

Mercados de ações exercem um papel de fundamental importância na economia dos países. De um lado, permitem que as empresas possam captar recursos para investimentos em tecnologia, infraestrutura ou outras formas de expansão, a partir da oferta pública de ações. Ao mesmo tempo, as ações representam, para os compradores individuais ou coletivos, uma importante alternativa de investimento para preservar ou aumentar sua poupança para uso futuro. De acordo com a NYSE Euronext (2012), mais de 90 milhões de cidadãos norte-americanos possuem quotas de ações utilizando investimentos individuais ou de fundos mútuos (fundos coletivos administrados por bancos ou outros agentes financeiros). Ademais, muitos outros participam no mercado de ações por meio de investimentos de fundos de pensão, companhias de seguros, entre outros.

Por outro lado, o preço das ações negociadas é volátil e depende de vários fatores, tais como o desempenho da empresa, o nível de atividade econômica de seu setor, indicadores econômicos, entre outros. Assim, investidores e gestores de fundos necessitam monitorar constantemente o comportamento dos preços das ações, a fim de tomar decisões corretas de investimento e evitar uma exposição excessiva a ativos de alto risco.

Técnicas de mineração de dados têm sido amplamente propostas para análise do mercado de ações, a fim de identificar padrões de comportamento das séries temporais. Essas técnicas partem da premissa de que tais padrões podem ser adequadamente explorados na previsão de preços, na recomendação de estratégias de operações ou até mesmo na negociação (*trading*) automática. Nessas abordagens, normalmente o vetor de atributos consiste de indicadores técnicos tradicionais, calculados a partir das séries temporais de preços e volumes negociados.

Entre as técnicas de mineração de dados, árvores de classificação têm algumas vantagens sobre outros métodos de aprendizado supervisionado. O primeiro é a facilidade de interpretação, já que árvores de classificação podem ser interpretadas pelos especialistas sem a necessidade de conhecimento prévio sobre o processo de treinamento. A segunda vantagem é que os algoritmos de criação da árvore normalmente escolhem os atributos mais relevantes, eliminando a necessidade de procedimentos externos de seleção de atributos. Tais características tornam os algoritmos de árvores de classificação métodos atraentes de aprendizado supervisionado em diversos domínios, particularmente no mercado de ações. De fato, os tomadores de decisão nesse mercado buscam sempre analisar e validar os classificadores construídos, a fim de confirmar se as regras de classificação geradas e os indicadores técnicos selecionados são consistentes.

Miró-Julià et al (2010) propõem a aplicação de árvores de decisão para detectar oportunidades de compra e venda, utilizando atributos binários derivados de indicadores técnicos. Lauretto (1996) e Stern et al (1998) apresentaram um algoritmo para indução de árvores de classificação denominado REAL (Real-Valued Attribute Learning), juntamente com um estudo de caso na aplicação de estratégias de operação de compra e venda.

Random forests são estruturas compostas por *ensembles* de árvores de classificação, em que a predição da classe para novas instâncias é baseada em um sistema de votação: a nova instância é submetida a cada árvore na floresta, e a classe designada é (em geral) aquela mais votada. As *Random forests* são usualmente mais estáveis que árvores de classificação individuais, no sentido de serem menos sensíveis à perturbação nos dados de treinamento [Breiman, 2001].

O objetivo desse trabalho é apresentar uma avaliação empírica da utilização de *Random Forests* para recomendação de estratégias de operação no mercado de ações. Propomos uma abordagem de aprendizado supervisionado na qual os atributos são constituídos por indicadores técnicos tradicionais, e as classes se baseiam em duas estratégias de operação:

- *Compra-Venda*: o investidor compra uma quota da ação e revende-a se o preço variar acima ou abaixo de limites determinados, ou após um certo período (em dias); e
- *Venda-Compra*: o investidor vende uma quota da ação (caso não possua, deve alugar essa quota) e recompra uma quota equivalente se o preço variar acima ou

abaixo dos limites ou após um certo período.

Por meio da simulação dessas operações sobre séries históricas de preços, cada dia de negociação recebe o rótulo de uma dentre três classes, representando a melhor decisão a ser tomada naquele dia: *Compra-Venda*, *Venda-Compra* (se alguma dessas operações terminou com lucro) ou *Nenhuma operação* (se nenhuma dessas estratégias foi lucrativa).

Esta abordagem é inspirada na proposta de Lauretto (1996) e Stern et al (1998), sendo que as principais diferenças neste trabalho são: a incorporação da estratégia *Venda-Compra*, a adoção de um método de validação cruzada adaptada para séries temporais [Hyndman e Athanasopoulos, 2012] e a incorporação de um procedimento de seleção ótima dos parâmetros das estratégias para cada papel. Os três principais índices de desempenho analisados neste trabalho foram: porcentagem de oportunidades aproveitadas pelo classificador, porcentagem de operações que obtiveram sucesso e média de retorno por aplicação.

2. Materiais e Métodos

Esse estudo é baseado em séries históricas fornecidas pela BM&F Bovespa, contendo as principais informações sobre os papéis negociados: código da empresa e da ação, preços diários (de abertura, de fechamento, médio, mínimo e máximo), quantidade de transações e volume negociados¹. Para a análise de desempenho, foram selecionadas as 68 ações que integram o Índice Bovespa [BM&F BOVESPA, 2012], devido à sua alta liquidez e seu alto volume de negociações.

Silva et al (2010) analisaram o comportamento do Índice Bovespa no período de 2005 a 2009, verificando uma forte influência da recessão econômica de 2008 sobre os preços das ações da Bovespa. Em particular, foi observada uma acentuada queda de preços no último quadrimestre de 2008 e posterior recuperação ao longo de 2009. Para evitar que o comportamento atípico dos preços nesses dois anos interferisse nos resultados das análises, foi selecionado o período de janeiro de 2010 a outubro de 2012 para avaliação da metodologia proposta.

As rotinas de processamento das séries históricas e dos testes descritos nas próximas seções foram implementadas em R, um ambiente integrado para processamento estatístico e análise de dados [R Core Team, 2013]. Sua escolha neste estudo se deu em razão de algumas de suas qualidades: distribuição livre e gratuita; disponibilidade de implementações de algoritmos de aprendizado, incluindo *Random Forest* [Liaw & Wiener, 2002]; linguagem de programação de alto nível e de rápida aprendizagem.

2.1. Random forests

Random Forests, introduzidas por Breiman (2001), são classificadores agregados compostos de *ensembles* de árvores de classificação independentes. A classificação de uma nova instância é feita com base em um sistema de votação, onde cada instância é classificada por cada árvore individualmente e os votos das classes são contados. Embora na maioria dos casos o critério de maioria dos votos seja utilizado (a classe com maior número de votos é escolhida), é possível estabelecer um limiar mínimo tal que uma classe

¹Uma descrição detalhada dos dados utilizados por ser encontrada em:
<http://www.bmfbovespa.com.br/pt-br/cotacoes-historicas/FormSeriesHistoricas.asp>

é selecionada apenas se atingir um percentual mínimo de votos entre as árvores. Como discutiremos adiante, essa funcionalidade é importante nesse trabalho, já que nosso interesse é somente recomendar a aplicação de uma estratégia se houver uma alta convicção de seu sucesso (medida pelo percentual de árvores que recomendam a aplicação dessa estratégia).

Apresentamos um breve sumário da construção de cada árvore em uma *Random Forest*. Denotamos por N a quantidade de exemplos e por M a quantidade de atributos no conjunto de treinamento original.

1. Uma amostragem de *bootstrap* de tamanho N é retirada do dado original, e é usada para a indução de uma nova árvore.
2. A cada divisão de um nó, $m \ll M$ atributos são selecionados aleatoriamente dentre os M atributos originais, e a melhor separação desses m atributos é utilizada para dividir o nó.

O valor de m é fixado durante a construção da floresta, e deve ser calibrado pelo usuário. O pacote *randomForest* [Liaw and Wiener, 2002], usado nesse trabalho, usa o valor $m = \sqrt{M}$ como default.

3. Cada árvore é expandida o máximo possível, sendo que nenhum procedimento de poda é adotado.

Breiman (2001) demonstra que a taxa de erro de classificação da *Random Forest* aumenta com a correlação entre as árvores e diminui com a força individual de cada árvore da floresta. A seleção aleatória dos exemplos e dos atributos tem como objetivo diminuir a correlação entre as árvores.

Na Seção 2.3 é apresentada uma ilustração da aplicação de *Random Forests* dentro do contexto deste trabalho.

2.2. Indicadores técnicos

Para a construção dos vetores de características, foram adotados seis indicadores técnicos dentre os mais utilizados [Puga et al., 2010; ADVFN, 2013]. As implementações desses indicadores estão disponíveis no Pacote *TTR* [Ulrich, 2012]. Apresentamos a seguir um breve sumário desses indicadores, denotando por $P(t)$ preço de fechamento do ativo no dia t .

Média móvel simples (MMS): É uma média aritmética simples, calculada sobre os últimos n preços de fechamento do ativo desejado:

$$\text{MMS}_n(t) = \sum_{i=0}^{n-1} P(t-i) / n. \quad (1)$$

Média móvel exponencial: Trata-se de uma média móvel com um fator de suavização que atribui maior peso aos valores mais recentes, sendo portanto capaz de acompanhar as oscilações de preço mais rapidamente do que a média móvel simples. Seu cálculo é dado por

$$\text{MME}_n(t) = \sum_{i=0}^{n-1} q(1-q)^i P(t-i) / \sum_{i=0}^{n-1} q(1-q)^i, \quad q = 2/(n+1). \quad (2)$$

Taxa de variação (ROC - rate of change): É considerado um indicador de movimento do preço, pois mede a variação percentual (positiva ou negativa) do preço no dia t em relação a $t - n$:

$$ROC_n(t) = 100 \left[\frac{P(t) - P(t-n)}{P(t-n)} \right]. \quad (3)$$

Para alguns ativos, o ROC apresenta padrões cíclicos, possibilitando prever movimentações futuras dos preços.

Oscilador Estocástico: É composto por três elementos. O primeiro, $\%K_n$, realiza uma mudança de escala do preço de fechamento no dia t em relação à amplitude de preços nos últimos n dias:

$$\%K_n = 100 \left[\frac{P(t) - Pm_n(t)}{PM_n(t) - Pm_n(t)} \right], \quad (4)$$

onde $PM_n(t)$ e $Pm_n(t)$ denotam, respectivamente, o preço máximo e o preço mínimo do ativo no período $t, t-1, \dots, t-(n-1)$. A motivação para esse indicador é que, em um período de tendência de alta, o preço de fechamento $P(t)$ tende a ficar próximo do preço máximo $PM_n(t)$, e portanto $\%K$ tende a ficar próximo de 100. Analogamente, em um período de tendência de baixa, $P(t)$ tende a ficar próximo do preço mínimo $Pm_n(t)$ e $\%K$ se aproxima de 0.

O segundo elemento, $\%D_n$, é uma média móvel simples de 3 dias de $\%K_n$. Quando $\%K_n$ cruza $\%D_n$ para cima, considera-se o início de uma tendência de alta; quando $\%K_n$ cruza $\%D_n$ para baixo, considera-se o início de uma tendência de queda.

O terceiro elemento, denotado por $Slow\%D_n$, é uma média móvel simples de 3 dias de $\%D_n$. De forma similar a $\%K_n$ com $\%D_n$, os cruzamentos de $\%D_n$ com $Slow\%D_n$ são sinais de inícios de tendências de alta ou de queda.

Índice de força relativa (IFR): Esse indicador realiza uma comparação entre os movimentos diários de alta (induzidos pela força dos compradores) e de baixa (induzidos pela força dos vendedores) nos últimos n dias. A *força relativa* de um papel é definida como:

$$FR_n(t) = \frac{\text{média de altas}}{\text{média de baixas}} = \frac{\sum_{i=0}^{n-1} \max[0, P(t-i) - P(t-i-1)]}{\sum_{i=0}^{n-1} \max[0, P(t-i-1) - P(t-i)]}. \quad (5)$$

O índice de força relativa é uma bijeção de FR_n no intervalo $[0, 100]^2$:

$$IFR_n(t) = 100 - \left[100 / (1 + FR_n(t)) \right]. \quad (6)$$

De forma análoga a $\%K_n$, valores de IFR_n próximos de 100 indicam tendências de alta, e valores de IFR_n próximos de 0 indicam tendências de baixa.

Convergência/divergência de médias móveis (MACD): Sua sigla vem do nome em inglês *Moving average convergence/divergence*, e é composto por três elementos. O primeiro, denominado linha MACD, é formado pela diferença entre uma média móvel exponencial de 12 dias e uma média móvel exponencial de 26 dias:

$$MACD(t) = MME_{12}(t) - MME_{26}(t). \quad (7)$$

$MACD(t) > 0$ indica que as expectativas mais recentes são mais favoráveis para alta do que as anteriores, e $MACD(t) < 0$ indica o comportamento inverso. Valores de MACD próximos de zero indicam que a oferta e a demanda estão em equilíbrio.

²Convenciona-se que, quando não ocorrem altas no período, $IFR_n(t) = 0$, e quando não ocorrem baixas no período, $IFR_n(t) = 100$.

O segundo elemento, denominado linha de sinal e denotado por $S(t)$, é calculado pela média móvel exponencial de 9 dias do MACD. $S \gg 0$ indica alta demanda e $S \ll 0$ indica alta oferta do papel.

O terceiro elemento, denominado histograma MACD, é obtido pela diferença entre a linha MACD e a linha de sinal:

$$\text{HMACD}(t) = \text{MACD}(t) - S(t). \quad (8)$$

Quando HMACD está acima da linha zero e começa a se aproximar da mesma, este pode ser uma indicação de que a tendência de alta está perdendo força. Em sentido oposto, quando HMACD está abaixo da linha zero e começa a se aproximar da mesma, este pode ser um sinal de que a tendência de queda está perdendo força.

A construção dos vetores de atributos foi feita com 23 elementos, compreendendo:

- MMS_n de 3, 13 e 21 dias;
- MME_n de 5, 13 e 21 dias;
- ROC_n de 13 e 21 dias;
- $\%K_n$, $\%D_n$, $\text{Slow}\%D_n$ de 7, 14 e 21 dias;
- IFR_n de 9, 14 e 21 dias;
- linhas MACD, S e HMACD.

2.3. Estratégias de operação e validação dos dados

Uma estratégia de operação de mercado é um conjunto de regras predefinidas que determinam as ações de um operador no mercado. Consideramos neste trabalho dois tipos de estratégias de operações. A seguinte notação é utilizada:

- t denota o dia de início da operação
- g é o máximo ganho esperado (*stop-gain*)
- l é a máxima perda tolerada (*stop-loss*)
- d é a duração máxima de uma operação (em dias)

Compra-Venda(t, g, l, d): Comprar a ação no dia t e vendê-la quando uma das seguintes condições é satisfeita:

1. O preço de fechamento da ação sobe acima de $g\%$ em relação ao preço do dia t ;
2. O preço de fechamento da ação cai abaixo de $l\%$ em relação ao preço do dia t ;
3. Após d dias, se nenhuma das condições acima foi satisfeita no período $[t+1, t+d]$.

Venda-Compra(t, g, l, d): Alugar uma quota da ação para vendê-la no dia t , e recomprar uma quota equivalente da mesma quando a primeira das condições abaixo é satisfeita:

1. O preço de fechamento da ação cai abaixo de $g\%$ em relação ao preço do dia t ;
2. O preço de fechamento da ação sobe acima de $l\%$ em relação ao preço do dia t ;
3. Após d dias, se nenhuma das condições acima foi satisfeita no período $[t+1, t+d]$.

Note que nas estratégias de *Compra-Venda* e *Venda-Compra*, o retorno líquido é calculado pela diferença entre os preços de compra e venda, descontados o custo da negociação (ex. taxas de corretagem). Na estratégia de *Venda-Compra*, há um custo adicional de aluguel que deve ser considerado. Nas simulações realizadas, assumiu-se um custo de operação (compra e venda) de $c = 1\%$ do valor da ação, e para o aluguel

da ação foi assumido o valor fixo de 0.05% por dia. Uma operação é considerada *bem-sucedida* se seu retorno líquido é positivo, e *mal-sucedida* caso contrário.

A Figura 1 mostra dois exemplos hipotéticos de aplicação da estratégia de Compra-Venda. No primeiro caso (a), a variação de preço (linha vermelha) atinge o ganho esperado (g) e a operação termina em sucesso (com retorno líquido positivo) antes do dia $t + d$. No segundo caso (b) a variação do preço oscila entre $-l$ e g até o dia $t + d$, quando a operação termina. Nesse caso, uma vez que o retorno líquido é negativo, ou seja, a variação final do preço fica abaixo do custo (c), a operação é considerada mal-sucedida.

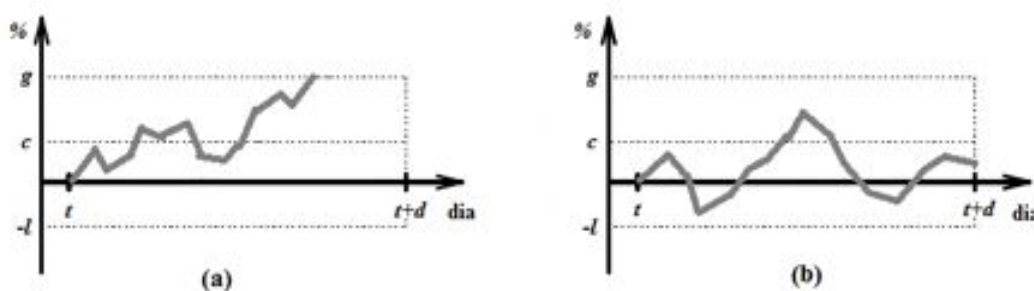


Figura 1. Exemplos de aplicação bem sucedida (a) e mal sucedida (b) da estratégia Compra-Venda (adaptado de Stern et al, 2008)

A classificação dos exemplos é feita da seguinte forma: fixados os parâmetros g, l e d , verificamos o sucesso/fracasso da estratégia $Compra-Venda(t, g, l, d)$ e $Venda-Compra(t, g, l, d)$ para cada dia t da série histórica. Se alguma das estratégias é bem-sucedida, rotula-se o dia t com a classe correspondente ($C=Compra-Venda$, $V=Venda-Compra$). Se nenhuma das estratégias é bem-sucedida, atribuímos a classe $N=Nenhuma$ operação.

Esse critério de classificação assume que uma operação do tipo $Compra-Venda$ e uma do tipo $Venda-Compra$ iniciadas no mesmo dia t não poderiam ser simultaneamente bem-sucedidas. Em geral, isso será verdadeiro se ambas as estratégias foram parametrizadas com o mesmos valores de g, l e d (máximo ganho esperado, máxima perda tolerada e duração máxima da operação) e se os valores desses parâmetros forem definidos sob premissas razoáveis (por exemplo, $l < g$).

Além do rótulo da classe, também é armazenado, para cada dia t , o retorno líquido obtido para cada operação (se alguma foi executada). Essa informação é usada para medir os retornos acumulados resultantes das sugestões dadas pelo classificador em um certo período (por exemplo, um ano).

Note-se que não há, *a priori*, valores ótimos para os parâmetros g, l e d – já que a definição de cada estratégia depende, por exemplo, da variabilidade do preço de cada ação. Por essa razão, foi implementado um procedimento automatizado para fixar esses parâmetros, como descrito na Subseção 2.4.

A Figura 2 ilustra a aplicação de uma *Random Forest* sobre um vetor de atributos $\mathbf{x} = (x_1, x_2, \dots, x_{23})$ composto pelos indicadores técnicos de um ativo calculados no dia t . A escolha da classe final (operação recomendada: C, V ou N) é feita a partir da classificação de \mathbf{x} por cada árvore da floresta e da consolidação dos votos nas classes. Os nós e

arestas destacados representam as regras de decisão aplicadas sobre x .

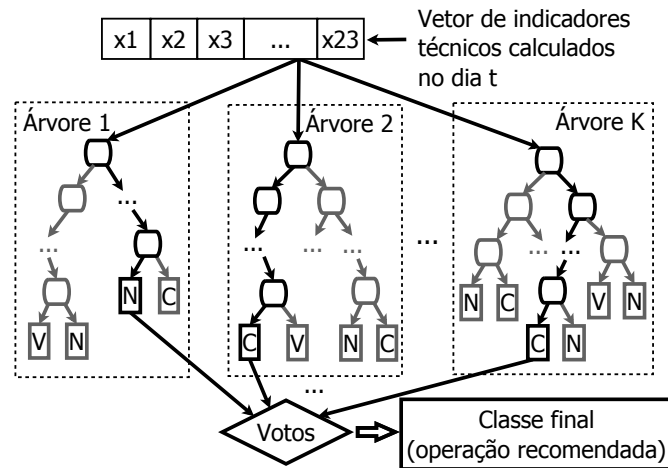


Figura 2. Exemplo de aplicação de uma *Random Forest* para a recomendação de operações sobre um vetor de indicadores técnicos

2.4. Validação cruzada

Uma das premissas assumidas nos esquemas usuais de validação cruzada k -fold e *leave-one-out* é que amostras dos conjuntos de treinamento e de validação são independentes. Contudo, nas séries temporais usualmente se observa uma forte dependência entre as observações, especialmente em janelas de tempo pequenas [Arlot e Celisse, 2010]. Para esse estudo foi aplicado o procedimento proposto por Hyndman e Athanasopoulos (2012), similar ao procedimento *leave-one-out*, exceto que o conjunto de treinamento consiste apenas de observações que ocorrem anteriormente à observação do conjunto de testes (evitando que observações futuras sejam usadas na construção do classificador). Essa abordagem requer que as primeiras observações sejam usadas apenas para treinamento e nunca consideradas para o conjunto de teste.

Denotando por T o tamanho total do *dataset*, e supondo que k observações sejam necessárias para produzir um conjunto de treinamento confiável, o processo funciona da seguinte maneira:

1. Repita os seguintes passos para $i = 1, 2, \dots, T - k$:
2. Construa a *Random Forest* usando as observações nos instantes $i, i + 1, i + 2, \dots, i + k - 1$, e teste-a na observação do instante $k + i$. Compute o acerto/erro (comparando a classe real e a prevista) e o retorno correspondente, caso alguma estratégia de operação tenha sido sugerida pela floresta.
3. Compute a acurácia total e os retornos líquido obtido para as $T - k$ amostras de teste.

Sabe-se que, no mercado financeiro, aplicar erroneamente uma estratégia de operação é mais crítico do que deixar de aplicá-la quando esta poderia ser bem sucedida. Enquanto o segundo tipo de erro implica em perda de oportunidade (mas sem perda de

capital), o primeiro erro resulta em perda direta de capital. Por essa razão, para evitar perdas elevadas causadas por operações mal sucedidas, adotamos uma abordagem mais conservadora do que a simples classificação por maioria dos votos: para cada exemplo de testes, somente são consideradas sugestões de operações Compra/Venda ou Venda/Compra (Classes C e V, respectivamente) se o percentual de votos na respectiva classe é superior a 70%. Se essa condição não é satisfeita, assume-se a classe N (nenhuma estratégia é aplicada).

A partir da validação cruzada, obtém-se uma matriz de confusão 3×3 no formato apresentado pela Tabela 1, onde as linhas representam as classes reais, e as colunas representam as classes preditas. A célula n_{ij} representa o número de exemplos de teste da classe i que foram classificados pela *Random Forest* como classe j , para $i, j \in \{C, N, V\}$.

Tabela 1. Formato da matriz de confusão obtida pela validação cruzada.

		Classe Predita		
		C	N	V
Classe Real	C	$n_{C,C}$	$n_{C,N}$	$n_{C,V}$
	N	$n_{N,C}$	$n_{N,N}$	$n_{N,V}$
	V	$n_{V,C}$	$n_{V,N}$	$n_{V,V}$

Três indicadores de desempenho foram considerados nesse trabalho:

- *AprovOport*: Taxa de aproveitamento de oportunidades: razão entre o número de operações bem-sucedidas e número de oportunidades (total de operações que seriam bem sucedidas, caso todas fossem aproveitadas):

$$AprovOport = \frac{(n_{V,V} + n_{C,C})}{(n_{V,V} + n_{V,N} + n_{V,C} + n_{C,V} + n_{C,N} + n_{C,C})}. \quad (9)$$

- *OperBemSuc*: Taxa de operações bem-sucedidas: razão entre o número de operações bem-sucedidas e o número total de operações sugeridas:

$$OperBemSuc = \frac{(n_{V,V} + n_{C,C})}{(n_{V,V} + n_{N,V} + n_{C,V} + n_{V,C} + n_{N,C} + n_{C,C})}. \quad (10)$$

- *RetMedOper*: Retorno médio por operação: razão entre a soma dos retornos líquidos produzidos pelas operações sugeridas (bem ou mal sucedidas) e o número total de operações sugeridas:

$$RetMedOper = \frac{ret(n_{V,V} + n_{N,V} + n_{C,V} + n_{V,C} + n_{N,C} + n_{C,C})}{(n_{V,V} + n_{N,V} + n_{C,V} + n_{V,C} + n_{N,C} + n_{C,C})}, \quad (11)$$

onde $ret(S)$ denota a soma dos retornos líquidos produzidos pelas S operações.

Esses indicadores de desempenho são combinados em uma função de *Score* definida pela seguinte combinação convexa:

$$Score = \alpha_{AO} * AprovOport + \alpha_{OBS} * OperBemSuc + \alpha_{RMO} * RetMedOper \quad (12)$$

Neste trabalho, os pesos α_{AO} , α_{OBS} , α_{RMO} foram escolhidos empiricamente de maneira que a função *Score* representasse uma posição conservadora, no sentido de favorecer estratégias com altas taxas de operações bem-sucedidas (alto α_{OBS}), mesmo com menores

taxas de aproveitamento de oportunidades ou retornos médios por operação. Os resultados apresentados nesse estudo foram obtidos utilizando-se $\alpha_{AO} = 0.10$, $\alpha_{OBS} = 0.85$, $\alpha_{RMO} = 0.05$.

A seleção dos valores de g , l e d é feita individualmente para cada papel, como segue. Primeiramente, são estabelecidos, para cada parâmetro, um conjunto de valores candidatos. Para esse estudo, foram definidos intervalos que abrangem os ganhos e perdas mínimos e máximos por período, bem como os números mínimo e máximo de dias esperado por operação. Esses conjuntos são:

- $g \in \{10\%, 15\%, 20\% \dots 35\%\}$
- $l \in \{3\%, 6\%, 9\%, \dots 15\%\}$
- $d \in \{10, 15, 20, \dots 35\}$

Para cada combinação de valores dos parâmetros g , l e d , as estratégias de operações são simuladas na série histórica, e os exemplos são rotulados com as classes correspondentes. Executa-se a validação cruzada e calculam-se os indicadores *AprovOport* (Eq. 9), *OperBemSuc* (Eq. 10), *RetMedOper* (Eq. 11) e *Score* (Eq. 12). Para cada papel, é escolhida a combinação de valores de g , l e d que maximiza a função *Score*. Para esse procedimento foram utilizadas as séries históricas dos anos de 2010 e 2011, adotando-se uma janela móvel de um ano para o conjunto de treinamento.

Após a escolha dos parâmetros ótimos, uma nova validação cruzada foi executada sobre os anos de 2011 e 2012, para avaliar o desempenho das *Random Forests*. Como nessa etapa também se adotou uma janela móvel de um ano para o conjunto de treinamento, os indicadores *AprovOport*, *OperBemSuc*, *RetMedOper* e *Score* foram calculados exclusivamente sobre o ano de 2012, garantindo-se assim que a avaliação final fosse realizada sobre um período disjunto daquele adotado para escolha dos parâmetros.

3. Resultados

A Figura 3 apresenta os indicadores de desempenho *AprovOport* (Eq. 9), *OperBemSuc* (Eq. 10), *RetMedOper* (Eq. 11) e *Score* (Eq. 12) para 40 ações – dentre as 68 originais – com maiores valores de *Score*.

Conforme descrito na seção anterior (Seção. 2.4), os pesos da função *Score* foram escolhidos a fim de maximizar a taxa de operações bem-sucedidas (*OperBemSuc*). O gráfico indica que esse objetivo foi alcançado, já que, para todas as 40 ações apresentadas, a taxa *OperBemSuc* foi superior a **75%**. Adicionalmente, em 29 ações essa taxa foi superior a **80%**. Esse resultado mostra que a abordagem foi conservadora, evitando operações que apresentavam alta probabilidade de insucesso.

O segundo indicador de desempenho importante é a taxa de aproveitamento de oportunidades (*AprovOport*). Ele indica o percentual das oportunidades (operações que, se executadas, resultariam em sucesso) que foram aproveitadas (foram de fato sugeridas pelo modelo). Uma baixa taxa de aproveitamento indicaria um modelo excessivamente conservador, que deixaria de executar a maioria das operações com retorno líquido positivo. Observa-se que, para 38 das 40 ações, a taxa *AprovOport* foi superior a **50%**, e que em 32 ações essa taxa foi superior a **70%**. Isso indica que, apesar de conservador, o modelo proposto é eficaz na indicação de operações com potencial para gerar lucro.

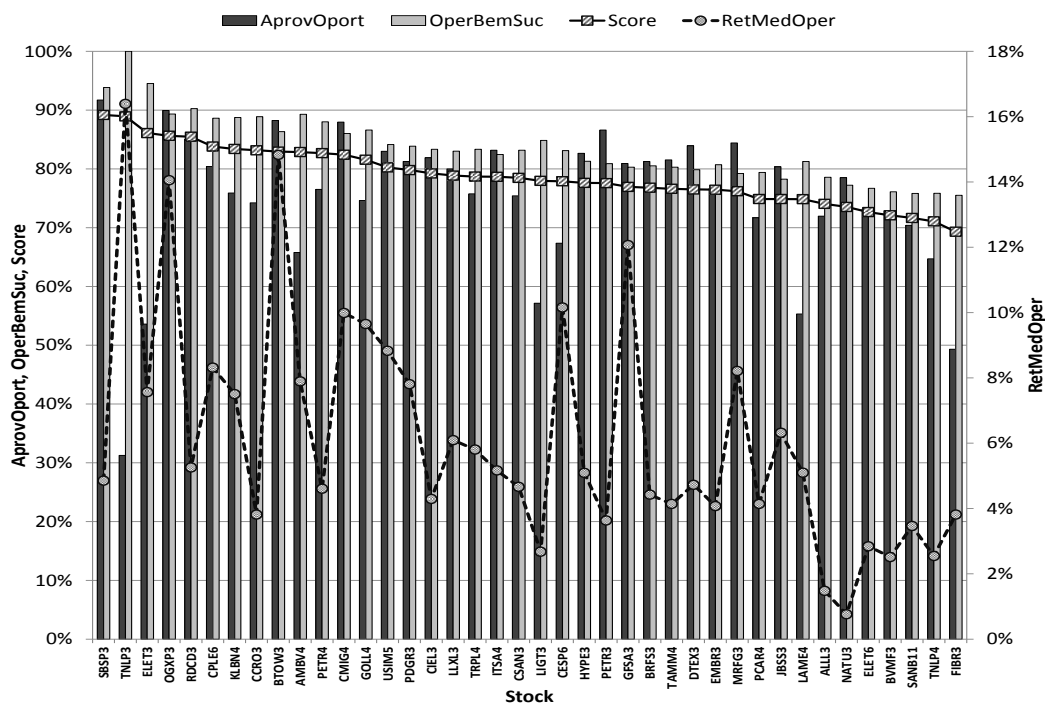


Figura 3. Indicadores de desempenho *AprovoOport*, *OperBemSuc*, *RetMedOper* e *Score* para as 40 ações com maiores valores de *Score*

O terceiro e último indicador, retorno médio por operação (*RetMedOper*), confirma o bom resultado obtido pelos outros dois indicadores. Para 30 das 40 ações, o retorno médio por operação foi superior a 4%, tendo alcançado mais de 5% para 25 papéis. Conforme descrito na Seção 2.4, a duração máxima de uma operação não ultrapassa 35 dias. Logo, retornos médios entre 4% e 5% por operação podem ser considerados resultados bastante expressivos, quando se leva em consideração a curta duração das operações.

4. Conclusões

Este trabalho discutiu uma abordagem para utilização de métodos de aprendizado supervisionado para a recomendação de estratégia de operações no mercado de ações, e apresentou uma análise empírica do desempenho das *Random Forests* para esse problema.

Os resultados apresentados são bastante promissores, motivando estudos adicionais e possíveis extensões. Uma análise da influência dos pesos da função de *Score* sobre os índices de desempenho é desejável, a fim de se obter uma calibração mais apropriada dos parâmetros das estratégias de acordo com o perfil do investidor. Outros índices de desempenho adotados no mercado financeiro serão também incorporados, para realização de comparações de nossos resultados com outros estudos.

Sob o aspecto de mineração de dados, ajustes mais refinados dos parâmetros do algoritmo de *Random Forests*, a adoção de outros métodos de aprendizado supervisionado e a incorporação de indicadores técnicos adicionais devem ser também considerados, pois poderão incrementar a acurácia do método.

Os autores são gratos pelo apoio e financiamento recebidos da EACH-USP, do IME-USP, da Coordenação de Aperfeiçoamento de Pessoal de Nível Superior (CAPES), do Conselho Nacional de Desenvolvimento Científico e Tecnológico (CNPq) e da Fundação de Apoio à Pesquisa do Estado de São Paulo (FAPESP). Parte deste trabalho foi realizado no âmbito do Programa de Educação Tutorial (MEC/SESu) na EACH-USP.

Referências

- ADVFN (2013). *Análise Técnica (análise gráfica)*. Disponível em: <http://br.advfn.com/educacional/analise-tecnica>. Acesso em: jan. 2013.
- Arlot, S. & Celisse, A. (2010). A survey of cross-validation procedures for model selection. *Statistical Surveys* 4, p. 40–79.
- Breiman, L. & Cutler, A. (2012). *Random Forests*. Disponível em http://www.stat.berkeley.edu/~breiman/RandomForests/cc_home.htm. Acesso em: set. 2012.
- BM&F BOVESPA (2012). *Índice Bovespa - Ibovespa*. Disponível em <http://www.bmfbovespa.com.br/indices/ResumoCarteiraTeorica.aspx?Indice=IBOVESPA&idioma=pt-br>. Acesso em: nov. 2012.
- Hyndman, R. J. & Athanasopoulos, G. (2012). *Forecasting: principles and practice*. Livro digital disponível em <http://otexts.com/fpp/>. Acesso em: out. 2012.
- Lauretto, M. S. (1996). *Árvores de Classificação para Escolha de Estratégias de Operação em Mercados de Capitais*. Dissertação de Mestrado, Instituto de Matemática e Estatística, Universidade de São Paulo.
- Liaw, A. & Wiener, M. (2002). Classification and Regression by randomForest. *R News* 2(3), p. 18–22.
- Miró-Julà, M., Fil-roig, G. & Isern-deyà, A. P. (2010). Decision Trees in Stock Market Analysis: Construction and Validation. In: N. García-Pedrajas et al. (Eds.): *IEA/AIE 2010*, Part I, LNAI 6096, p. 185–194.
- Mueller, W. & Wysotzki, F. (1994). Automatic construction of decision trees for classification. *Annals of Operations Research* 52, p. 231–247.
- NYSE Euronext (2012). *Why We Invest*. Disponível em <https://nyse.nyx.com/financial-literacy/all-about-investing/investing-basics/why-we-invest>. Acesso em: dez. 2012.
- Puga, R., Rodrigues, M. & Cerbassi, G. (coord) (2010). *Formação de traders: faça dinheiro na bolsa com a análise técnica*. Rio de Janeiro: Campus.
- Silva, R.; Zembruski, M.; Correa, F. C.; Lamb, L. C. (2010). Stock markets and criticality in the current economic crisis. *Physica A* 389, p. 5460–5467.
- Stern, J. M., Nakano, F., Lauretto, M. S. & Ribeiro, C. O. (1998). Algoritmo de Aprendizagem para Atributos Reais e Estratégias de Operação em Mercado. In: *Sixth Iberoamerican Conference on Artificial Intelligence - IBERAMIA'98*, Lisboa.
- R Core Team (2013). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Viena, Áustria. Disponível em <http://www.R-project.org>.
- Ulrich, J. (2012). *The TTR Package Reference Manual*. Disponível em <http://cran.r-project.org/web/packages/TTR/index.html>.