

**UNIVERSIDADE ESTADUAL PAULISTA "JÚLIO DE MESQUITA FILHO"**

**FACULDADE DE CIÊNCIAS - CAMPUS BAURU**

**DEPARTAMENTO DE COMPUTAÇÃO**

**BACHARELADO EM CIÊNCIA DA COMPUTAÇÃO**

**GUSTAVO TRIELLI AVILA**

**REGRESSÃO EM SÉRIES TEMPORAIS FINANCEIRAS COM RNN:  
UM ESTUDO COM MILHO FUTURO**

**BAURU-SP**

**2019**

GUSTAVO TRIELLI AVILA

**REGRESSÃO EM SÉRIES TEMPORAIS FINANCEIRAS COM RNN:  
UM ESTUDO COM MILHO FUTURO**

Trabalho de Conclusão de Curso do Curso de Bacharelado em Ciência da Computação da Universidade Estadual Paulista “Júlio de Mesquita Filho”, Faculdade de Ciências, Campus Bauru.

Orientador: Prof. Dr. Clayton Reginaldo Pereira

BAURU-SP

2019

Gustavo Trielli Avila    Regressão em séries temporais financeiras com RNN:  
um estudo com milho futuro/ Gustavo Trielli Avila. – Bauru-SP, 2019-    39  
p. : il. (algumas color.) ; 30 cm.  
Orientador: Prof. Dr. Clayton Reginaldo Pereira  
Trabalho de Conclusão de Curso – Universidade Estadual Paulista “Júlio de  
Mesquita Filho”  
Faculdade de Ciências  
Bacharelado em Ciência da Computação, 2019.  
1. Redes Neurais Artificiais 2. Previsão séries temporais financeiras 3. Long  
short-term memory (LSTM)

Gustavo Trielli Avila

## **Regressão em séries temporais financeiras com RNN: um estudo com milho futuro**

Trabalho de Conclusão de Curso do Curso de Bacharelado em Ciência da Computação da Universidade Estadual Paulista "Júlio de Mesquita Filho", Faculdade de Ciências, Campus Bauru.

Banca Examinadora

---

**Prof. Dr. Clayton Reginaldo Pereira**

Orientador

Universidade Estadual Paulista "Júlio de Mesquita Filho"

Faculdade de Ciências

Departamento de Computação

---

**Prof. Dra. Simone das Graças Domingues Prado**

Universidade Estadual Paulista "Júlio de Mesquita Filho"

Faculdade de Ciências

Departamento Computação

---

**Dr. Leandro Aparecido Passos Júnior**

Universidade Estadual Paulista "Júlio de Mesquita Filho"

Faculdade de Ciências

Departamento Computação

Bauru, \_\_\_\_\_ de \_\_\_\_\_ de \_\_\_\_\_.

# Resumo

Investimentos no mercado financeiro representam uma parcela significativa na economia de todo país. No Brasil o interesse da população nessa área cresceu muito nos últimos anos, se tornando cada vez mais um foco de pesquisas para aplicações nesse setor. O presente trabalho busca utilizar redes neurais LSTM na previsão de preços de uma série temporal financeira. Serão apresentados conceitos do mercado financeiro e de aprendizagem de máquina para fundamentar o trabalho. Foram propostos modelos para avaliar a dependência temporal para a previsão dos preços bem como para avaliar o impacto de indicadores técnicos e séries exógenas como variáveis independentes para compor a predição. O estudo foi feito utilizando a série temporal do milho futuro, um derivativo agrícola negociado na bolsa de valores.

**Palavras-chave:** série temporal financeira, redes neurais, LSTM, previsão de valores.

# Abstract

The financial market is an important part of every country's economy. Brazilians interest in this sort of investment increased a lot in the last years, so as the number of researches and applications related to this task. This work investigates the use of LSTM neural networks to forecast the prices of a financial time series. A theoretical background regarding the financial market and machine learning concepts will be presented to support the work. Models were proposed to evaluate the time dependence for price forecasting as well as to evaluate the impact of technical indicators and exogenous series as independent variables to compose the prediction. The study was done using the future corn time series, a stock-traded agricultural derivative.

**Keywords:** financial time series, neural networks, LSTM, price forecasting.

# Lista de figuras

Figura 1 – Especificações do contrato de Milho Futuro. . . . .	16
Figura 2 – Representação de um neurônio artificial . . . . .	20
Figura 3 – Representação de uma arquitetura LSTM . . . . .	21
Figura 4 – Dados de negociação dos contratos de milho futuro . . . . .	24
Figura 5 – Série temporal do preço de fechamento fechamento . . . . .	27
Figura 6 – Box plot do fechamento por ano e por trimestre . . . . .	28
Figura 7 – Box plot do volume por trimestre . . . . .	28
Figura 8 – Distribuição dos preços e plot de probabilidade. . . . .	29
Figura 9 – Conjunto de dados após manipulações . . . . .	29
Figura 10 – Previsão dos preços para os dados de treino . . . . .	34
Figura 11 – Previsão dos preços para os dados de validação . . . . .	35
Figura 12 – Previsão dos preços para os dados de teste . . . . .	36
Figura 13 – Convergência do modelo . . . . .	36

# Lista de quadros

Quadro 1 – Treino . . . . .	31
Quadro 2 – Validação . . . . .	32
Quadro 3 – Teste . . . . .	33



# Lista de abreviaturas e siglas

AM	Aprendizado de Máquina
API	Application Programming Interface
B3	Brasil, Bolsa, Balcão
IA	Inteligência Artificial
LSTM	Long Short-Term Memory
MACD	Moving Average Convergence/Divergence
MM	Média Móvel
MME	Média Móvel Exponencial
PLN	Processamento de Linguagem Natural
RMSE	Rooten Mean Squared Error
RN	Redes Neurais
RNA	Redes Neurais Artificiais
RNN	Redes Neurais Recorrentes
RSI	Relative Strength Index

# Sumário

<b>1</b>	<b>INTRODUÇÃO</b>	<b>11</b>
<b>2</b>	<b>FUNDAMENTAÇÃO TEÓRICA</b>	<b>14</b>
<b>2.1</b>	<b>Mercado Financeiro</b>	<b>14</b>
2.1.1	Bolsa de Valores	14
2.1.2	Derivativos e mercado futuro	14
2.1.3	Milho futuro	15
2.1.4	Série temporal financeira	16
2.1.5	Análise técnica	17
2.1.5.1	Média móvel simples	17
2.1.5.2	Média móvel exponencial	18
2.1.5.3	Histograma Moving Average Convergence/Divergence (MACD)	18
2.1.5.4	Índice de força relativa	18
<b>2.2</b>	<b>Aprendizado de máquina</b>	<b>19</b>
2.2.1	Aprendizagem supervisionada	19
2.2.2	Redes Neurais Artificiais	19
2.2.3	LSTM	20
<b>2.3</b>	<b>Métricas de desempenho</b>	<b>21</b>
2.3.1	Erro percentual absoluto médio	22
2.3.2	Erro quadrático médio	22
2.3.3	Raiz quadrada do erro quadrático médio	22
2.3.4	Acurácia	22
<b>3</b>	<b>METODOLOGIA DE PESQUISA</b>	<b>24</b>
<b>3.1</b>	<b>Base de dados</b>	<b>24</b>
<b>3.2</b>	<b>Ferramentas utilizadas</b>	<b>25</b>
3.2.1	Python	25
3.2.2	Scikit-learn	25
3.2.3	Pandas	25
3.2.4	Matplotlib	25
3.2.5	Keras	25
3.2.6	MetaTrader 5	26
<b>3.3</b>	<b>Modelos Propostos</b>	<b>26</b>
<b>4</b>	<b>DESENVOLVIMENTO</b>	<b>27</b>
<b>4.1</b>	<b>Análise dos dados</b>	<b>27</b>

<b>4.2</b>	<b>Tratamento dos dados</b>	<b>28</b>
4.2.1	Adição de novas variáveis	28
4.2.2	Pré-processamento	29
4.2.2.1	Normalização dos dados	29
4.2.2.2	Montagem das sequências	30
4.2.2.3	Treino, validação e teste	30
4.2.3	Treinamento dos Modelos	30
<b>4.3</b>	<b>Análise dos Resultados</b>	<b>31</b>
<b>5</b>	<b>CONCLUSÃO</b>	<b>37</b>
	<b>REFERÊNCIAS</b>	<b>38</b>

# 1 Introdução

O mercado financeiro é um importante elemento da economia destinado ao fluxo de recursos financeiros entre poupadores e tomadores, o que viabiliza os investimentos de modo geral, possibilitando uma otimização da utilização dos recursos financeiros da economia, que resulta no crescimento econômico (ANDREZO; LIMA, 2007). No Brasil o interesse da população no mercado financeiro tem crescido muito nos últimos anos. Segundo dados de 2018 do relatório da ANBIMA (ANBIMA, 2019) apenas 42% dos brasileiros investem em algum produto financeiro, mas o mercado financeiro vive uma transformação digital com forte processo de digitalização dos canais de relacionamento, o que contribui para a democratização do acesso aos produtos. Além disso, dados mais recentes do relatório disponibilizado pela B3 (B3, 2019 A), bolsa de valores oficial do Brasil, mostra que o número de investidores cresceu consideravelmente nos últimos dois anos. O número de CPFs cadastrados na bolsa cresceu 133% de 2017 até Setembro de 2019, enquanto no período de 2011 até 2017 cresceu apenas 6%.

Isso mostra que o interesse das pessoas sobre esse assunto tende a crescer ainda mais. Como este mercado possui um papel importante na economia de um país existe uma motivação por parte dos investidores em utilizar técnicas e estratégias que consigam prever o seu comportamento, para que possam fazer investimentos que gerem maiores retornos e menos risco. Existem diferentes hipóteses sobre a previsão do preço no mercado de ações, alguns autores dizem que não é possível realizar previsões do valor de um ativo com base em informações históricas sobre ele, entretanto outros dizem por meio de experimentos que é possível prever o comportamento de um ativo até certo nível (DAMETTO, 2018). Nesse sentido a utilização de RNA para previsão de valores futuros em séries temporais está entre as técnicas mais populares atualmente (GIACOMEL, 2016).

Existem dois tipos de análises predominantes sobre um ativo na hora de um investidor tomar uma decisão: a análise fundamentalista e a análise técnica. Na análise fundamentalista o investidor procura analisar fatores econômicos e o desempenho da companhia, como relatórios contábeis e demonstrações de resultados do exercício. Já na análise técnica o investidor analisa gráficos do movimento de um ativo e tenta identificar padrões para projetar movimentos futuros (ABE, 2017). Os indicadores técnicos são uma ferramenta auxiliar da análise técnica, eles são resultados de operações matemáticas sobre a série financeira de um ativo, e ajudam a prever o movimento do mesmo.

Grande parte dos trabalhos de pesquisa nessa área utiliza RNA *perceptron* multicamadas na previsão de valores sobre os papéis mais conhecidos, como índice Ibovespa no Brasil e ações de grandes companhias como a *Apple*. Um possível problema encontrado em utilizar esse tipo

de ativo como objeto de estudo é o viés que o modelo pode ter por conta da tendência em que o ativo se encontra. O índice Ibovespa por exemplo é o principal indicador da economia brasileira e é composto por ações das empresas mais importantes do mercado (B3, 2019 B), trata-se de um cenário muito amplo que reflete toda a economia de um país, podendo tornar complexo devido ao número de variáveis que podem influenciar no valor do ativo. Ações de grandes empresas também podem ser problemáticas pela maioria já estar consolidada em algum setor, sofrendo grandes variações apenas quando decorrentes de fatores externos como notícias ou avanços na área. Por esses motivos foi escolhido para esse trabalho utilizar um derivativo como objeto de estudo, mais especificamente o contrato de milho futuro.

O mercado de derivativos é um segmento do mercado financeiro composto por conjuntos de operações cujos valores derivam de outros ativos (ANDREZO; LIMA, 2007). O contrato de milho futuro é um derivativo agrícola no qual o seu preço deriva do preço da saca de milho, o ativo-objeto. A escolha de um derivativo agrícola deve-se ao fato de estar em um segmento de mercado mais bem definido e com menos fatores externos à esse mercado influenciando no preço, como notícias por exemplo. Por ser um derivativo agrícola seu preço pode ter um comportamento sazonal por conta do clima, que pode influenciar na oferta do milho. O contrato de milho futuro é um dos derivativos agrícolas mais negociados, e apesar de possuir uma alta volatilidade existe a possibilidade, por conta de fatores comentados acima, de que existam padrões nos dados que possam facilitar a previsão.

Além disso um derivativo do mercado futuro possui um atrativo em relação ao investimento com ações comuns que é a possibilidade de operar alavancado, o que atrai especuladores e traz um retorno muito maior se for possível prever o comportamento do papel. O conceito de mercado futuro e operação alavancada serão explicados nas seções seguintes. Um especulador é um tipo de investidor que tem o objetivo de lucrar com a movimentação dos preços em um curto período de tempo. É fato que um derivativo não tem a mera finalidade de especulação, mas os especuladores possuem o importante papel de trazer liquidez para esse mercado (HULL, 2016). Para fins de estudo a especulação se encaixa melhor no contexto de previsão de série financeira de preços.

Por se tratar de uma série temporal, este trabalho pretende utilizar um modelo *Long Short-Term Memory* (LSTM), que é um tipo de rede neural recorrente (RNN), para realizar a previsão dos preços, e avaliar se a dependência temporal da série de preços têm influência na previsão do preço futuro. Uma LSTM é um modelo de RNA que consegue absorver a dependência temporal das entradas do modelo (HOCHREITER; SCHMIDHUBER, 1997). Como foi mostrado em Giacomel (2016), vários trabalhos nessa área mostraram bons resultados utilizando RNAs *perceptron* multicamadas. Este trabalho pretende testar se a dependência temporal de uma série financeira pode influenciar na previsão do preço futuro. Em Giacomel (2016) não houve ganhos em inserir indicadores técnicos como variáveis de entrada e levanta-se a hipótese de que isso ocorreu pelo fato de RNAs aprenderem relações entre suas variáveis de

entrada, que é de certa forma a mesma informação que um indicador técnico traz, tornando-se uma informação redundante. Nesse trabalho também será aplicado indicadores técnicos como variáveis de entrada a fim de testar a hipótese levantada.

## 2 Fundamentação Teórica

Neste capítulo serão apresentados fundamentos e conceitos sobre o mercado financeiro, séries temporais financeiras e RNAs. Esses fundamentos são necessários para entendimento do completo do trabalho.

### 2.1 Mercado Financeiro

#### 2.1.1 Bolsa de Valores

A bolsa de valores é um elemento do mercado financeiro que viabiliza a negociação valores mobiliários em âmbito público (ANDREZO; LIMA, 2007). De acordo com (FLEURIET; GALVAO; MENDES, 2005) os valores mobiliários foram definidos em lei como sendo:

- As ações, debêntures e bônus de subscrição.
- Os cupons, direitos, recibos de subscrição e certificados de desdobramento relativos aos valores mobiliários.
- Os certificados de depósitos de valores mobiliários.
- As cédulas de debêntures.
- As cotas de fundos em valores mobiliários ou clubes de investimentos em quaisquer ativos.
- As notas comerciais.
- Os contratos futuros, de opções e outros derivativos, cujos ativos subjacentes sejam valores mobiliários.

As bolsas de valores são associações civis, sem fins lucrativos e com funções de interesse público. Entre suas principais atribuições estão oferecer um mercado para a cotação de títulos e valores mobiliários e proporcionar liquidez às aplicações de curto e longo prazo através de um mercado contínuo (FLEURIET; GALVAO; MENDES, 2005). O Brasil possui atualmente uma única bolsa de valores em operação, a B3 S.A (Brasil, Bolsa, Balcão).

#### 2.1.2 Derivativos e mercado futuro

Os derivativos são instrumentos financeiros cujo preço deriva de algum outro ativo subjacente. O ativo subjacente pode ser físico como *commodities* agrícolas (milho, café, boi

gordo) ou financeiro (ações, índices, taxas de juros). No Brasil os derivativos são negociados em bolsa de valores (HULL, 2016).

No mercado futuro são negociados contratos padronizados de *commodities* agrícolas, pecuárias, avícolas, metálicas e financeiras. Esses contratos firmam a compra ou venda de um ativo em uma data futura por um valor determinado, esse valor determinado é o preço de futuro do ativo. Um contrato é formado por uma parte compradora e outra vendedora, no qual a parte compradora se compromete a pagar em uma data futura a quantia e o ativo determinado no contrato. Da mesma forma a parte contrária se compromete a vender em uma data futura a quantia e o ativo determinado no contrato. No Brasil a B3 determina as especificações de cada contrato e também garante a compensação das partes para garantir a solidez do mercado.

Os derivativos do mercado futuro possuem duas principais características, a proteção para investidores contra a oscilação de preços e a liquidez do mercado (HULL, 2016). Logo, existem três tipos principais de investidores nesse mercado, os *hedgers*, os especuladores e os arbitradores. *Hedge* é o nome dado a operações com estratégia de proteção para os riscos de investimento. Os *hedgers* são justamente investidores que buscam eliminar os riscos de perdas devido a oscilação dos preços no mercado à vista. Os especuladores não possuem interesse comercial na *commodity*, eles entram e saem rapidamente buscando lucrar com a oscilação dos preços. Os arbitradores buscam garantir um lucro sem risco pela diferença de preços entre dois ativos.

### 2.1.3 Milho futuro

Segundo dados da B3 (2019 B) o contrato de milho futuro é atualmente o derivativo agrícola mais negociado em número de contratos. O principal indicador do preço do milho é o Indicador de Preços do Milho Esalq/BM&FBOVESPA, este indicador representa a média dos preços da saca de milho de 60Kg praticados no mercado à vista na região de Campinas, no estado de São Paulo. Esses contratos possuem somente liquidação financeira, isso significa que não existe a entrega física do produto, portanto quando um investidor assume uma posição em um contrato futuro, ele pode encerrar sua posição até a data de vencimento e obrigatoriamente no vencimento, liquidando a variação do preço futuro do milho enquanto a posição estava aberta. Os contratos são definidos pela B3 com as especificações mostradas na Figura 1

Um exemplo de operação para proteção é um produtor de milho que quer se proteger do risco de possível queda do preço da saca realiza uma operação de *hedge*, tomando uma posição de venda em um contrato futuro para garantir seu lucro em um determinado valor. Dessa forma, se o preço do milho à vista cair, o lucro do produtor será o preço de venda no mercado à vista mais a liquidação do contrato, e se o preço do milho subir, o lucro será preço de venda no mercado à vista menos a liquidação do contrato. Assim o produtor garante o lucro no futuro independente da oscilação de preço da saca de milho.



Figura 1 – Especificações do contrato de Milho Futuro.

Objeto de negociação	Milho em grão a granel, com odor e aspectos normais, duro ou semiduro e amarelo.
Código de negociação	CCM
Tamanho do contrato	450 sacas de 60kg líquidos (equivalentes a 27 toneladas métricas).
Cotação	Reais por saca, com duas casas decimais.
Variação mínima de apregoação	R\$0,01.
Lote padrão	1 contrato.
Último dia de negociação	Dia 15 do mês de vencimento.
Data de vencimento	Dia 15 do mês de vencimento. Caso não haja sessão de negociação, a data de vencimento será a próxima sessão de negociação.
Meses de vencimento	Janeiro, março, maio, julho, agosto, setembro e novembro.
Liquidação no vencimento	Financeira.

Fonte: Disponível em: <[http://www.b3.com.br/pt\\_br/produtos-e-servicos/negociacao/commodities/ficha-do-produto-8AE490CA6D41D4C7016D45F3CB0A38F0.htm](http://www.b3.com.br/pt_br/produtos-e-servicos/negociacao/commodities/ficha-do-produto-8AE490CA6D41D4C7016D45F3CB0A38F0.htm)>. Acesso em 3 de novembro de 2019

O milho futuro possui contratos de 450 sacas, e a cotação em reais por saca, o que significa que se um contrato está cotado em R\$30,00 reais o valor total dele é R\$13.500,00 reais. Entretanto isso não significa que um investidor precisa de R\$13.500,00 reais para operar um contrato, a bolsa permite operar um contrato com uma parcela do valor total. Normalmente esse valor gira em torno de 5,49% do valor total do contrato, então no caso de um contrato cotado a R\$30,00 reais seria necessário R\$741,15 reais para realizar uma operação. Este tipo de operação é denominada operação alavancada, quando se opera uma quantia maior do que a que se tem. A operação alavancada é permitida nesse caso para atrair especuladores, uma vez que se pode ganhar mais dinheiro investindo menos, e assim trazer liquidez ao mercado. Os contratos de milho futuro possuem vencimento nos meses de Janeiro, Março, Maio, Julho, Agosto, Setembro e Novembro.

#### 2.1.4 Série temporal financeira

Uma série temporal é uma coleção de dados sequenciais ao longo do tempo. Uma série temporal financeira possui informações sobre a negociação de um ativo ao longo do tempo. Usualmente essas informações são o preço de abertura, o preço máximo atingido no período, o preço mínimo atingido no período, o preço de fechamento no período, e o volume de negociações no período. Dentro desse cenário existe uma linha de pesquisa chamada Econometria, que estuda a modelagem de séries temporais financeiras através de ferramentas estatísticas e matemáticas.

Séries temporais financeiras possuem algumas características particulares:

- **Pontos influentes** – são valores incomuns, que fogem do padrão do restante da série temporal. Em séries financeiras, podem ser considerados pontos influentes, momentos de alta volatilidade no mercado, gerando grandes altas ou quedas nos preços, que logo depois voltam ao seu patamar normal
- **Heteroscedasticidade condicional** – a variância para os valores de entrada e saída da série temporal não é constante com o passar do tempo, tornando o comportamento da série temporal aleatório
- **Não-linearidade** - devido à sua complexidade e ao seu comportamento estocástico, não é possível modelar este tipo de série temporal com uma função linear.
- **Sazonalidade** - os valores alcançados por uma série temporal financeira podem variar de acordo com a época do ano e repetir o padrão de variação nos anos seguintes.

### 2.1.5 Análise técnica

A análise técnica possibilita ao investidor fazer uma leitura ,através dos gráficos, dos movimentos da massa de investidores do mercado e tentar acompanhá-los (ABE, 2017). Os indicadores técnicos são resultados de manipulações matemáticas sobre a série de preços. Eles tem a finalidade de auxiliar na análise tentando indicar de alguma forma o movimento do mercado. Existem três grupos de indicadores, os rastreadores de tendência, osciladores e mistos. Rastreadores de tendência indicam a tendência que o mercado se encontra. Osciladores captam pontos de inflexão na série de preços, e os mistos tentam indicar a psicologia de massa do mercado, isto é, as expectativas dos investidores em relação ao mercado (ELDER, 1993). Para este trabalho foram utilizados 4 indicadores técnicos entre os mais populares, sendo 3 rastreadores de tendência (média móvel simples, média móvel exponencial, histograma MACD) e 1 oscilador (índice de força relativa).

#### 2.1.5.1 Média móvel simples

A média móvel (MM) simples mostra o valor médio dos dados em um determinado período. Uma média móvel de 10 dias mostra o valor médio dos dados nos últimos 10 dias.

$$MM_{simples} = \frac{P_1 + P_2 + \dots + P_N}{N} \quad (2.1)$$

Onde:

$P$  = preço no instante

$N$  = número de dias da média móvel

### 2.1.5.2 Média móvel exponencial

A média móvel exponencial (MME) é semelhante a MM, a diferença é que ela atribui maior peso a dados mais recentes e reage com mais rapidez às mudanças.

$$MME = P_t * K + MME_{t-1} * (1 - K) \quad (2.2)$$

Onde:

$$K = \frac{2}{N + 1}$$

$N$  = número de dias da MME

$P_t$  = preço no instante  $t$

$MME_{t-1}$  = MME no instante  $t - 1$

### 2.1.5.3 Histograma Moving Average Convergence/Divergence (MACD)

O MACD original é um dos indicadores mais utilizados na análise técnica. Ele é um rastreador de tendência composto por três MMEs. Os períodos das MMEs ficam a critério do investidor, mas as mais utilizadas são de 26 dias, de 12 dias, e de 9 dias. O MACD é representado por duas linhas, a linha MACD e a linha de sinal. A linha MACD é composta pela subtração da MME de 12 dias sobre os preços de fechamento pela MME de 26 dias sobre os preços de fechamento. A linha de sinal é a MME de 9 dias sobre a linha MACD. Finalmente o histograma MACD é calculado a partir da subtração da linha MACD pela linha de sinal. O histograma MACD proporciona uma visão não só da tendência mas também da sua força (ELDER, 1993).

- a) Linha MACD = MME de 12 dias - MME de 26 dias
- b) Linha Sinal = MME de 9 dias sobre a linha MACD
- c) Histograma MACD = Linha MACD - Linha Sinal

O cálculo da MME é fornecido pela equação 2.2.

### 2.1.5.4 Índice de força relativa

O índice de força relativa (Relative Strength Index - RSI) é um indicador técnico criado por Welles Wilder Junior em 1978. Ele é um indicador da classe dos osciladores e mede a mudança recente do valor dos ativos, para avaliar se situações de sobre compra (super valorizado) ou sobre venda (sub valorizado). O RSI é demonstrado em um intervalo entre 0 e 100, onde valores acima de 70 indicam que o preço do ativo pode estar super valorizado, enquanto valores abaixo de 30 indicam uma desvalorização (ELDER, 1993).

$$RSI = 100 - \frac{100}{1 + RS} \quad (2.3)$$

Onde:

$$RS = \frac{\text{média dos fechamentos em ALTA para } n \text{ dias}}{\text{média dos fechamentos em BAIXA para } n \text{ dias}}$$

## 2.2 Aprendizado de máquina

O aprendizado de máquina AM é um campo da Inteligência Artificial (IA) e pode ser definido como um processo automático de captura relações entre os dados com o objetivo de reconhecimento e detecção de padrões ([KELLEHER; NAMEE; D'ARCY, 2015](#)).

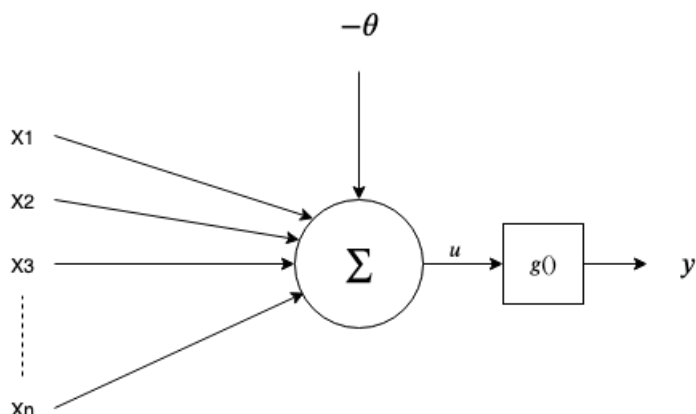
### 2.2.1 Aprendizagem supervisionada

O aprendizado supervisionado é a técnica utilizada para construir modelos preditivos. Essa técnica constrói um modelo que aprende relações entre um conjunto de dados de entrada e de saída, através da análise de uma base histórica, com a finalidade de reproduzir dados de saída para dados de entrada que não estavam presentes na base ([KELLEHER; NAMEE; D'ARCY, 2015](#)).

### 2.2.2 Redes Neurais Artificiais

Uma RNA é um modelo matemático baseado no funcionamento dos neurônios do cérebro humano, simulando conexões sinápticas ([DAMETTO, 2018](#)). Uma RNA é representada por um grafo orientado onde cada nó do grafo representa um neurônio e cada aresta representa uma sinapse. A figura 2 representa um modelo simples onde um neurônio recebe uma lista de entradas  $x$  e cada conexão com o neurônio possui um peso. Então o neurônio realiza a soma ponderada  $\sum$  das entradas pelos pesos, do resultado é subtraído um valor  $\theta$  chamado limiar de ativação, então finalmente esse resultado passa por uma função de ativação  $g()$  que limita a saída o neurônio em um certo intervalo, gerando uma saída  $y$  ([DAMETTO, 2018](#)). Vários neurônios podem ser combinados formando uma rede complexa capaz de resolver problemas maiores.

Figura 2 – Representação de um neurônio artificial



Fonte: Elaborada pelo autor

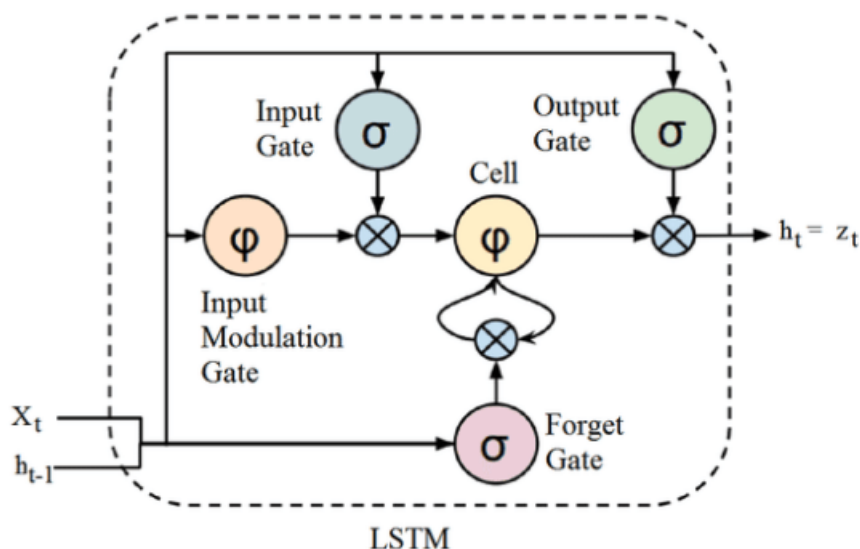
### 2.2.3 LSTM

Uma RNN é um tipo de rede neural na qual os neurônios possuem conexões de retroalimentação que permite representar uma entrada recente em forma de ativação. Isso é útil em muitos casos para resolver tarefas em que existe uma dependência sequencial entre os dados, como em problemas de Processamento de Linguagem Natural (PLN) ou em previsões de séries temporais (HOCHREITER; SCHMIDHUBER, 1997). Entretanto uma RNN comum não é capaz de guardar a informação de entradas anteriores por muito tempo. Uma rede LSTM é um modelo proposto por (HOCHREITER; SCHMIDHUBER, 1997) para contornar o problema de dissipação do gradiente das RNNs convencionais por não conseguirem armazenar longas dependências temporais. A figura 3 representa uma célula LSTM, onde  $\sigma$  representa a função *sigmoid* e  $\varphi$  representa a função *tanh*.

As fundamentações sobre as células LSTM que serão declaradas a seguir foram feitas com base no material de Academy (2019). Uma célula LSTM possui uma estrutura em cadeia que contém quatro RNAs, os portões, e diferentes blocos de memória chamados de células. As informações de memória são retidas nas células e as manipulações são feitas pelos portões.

- **Forget Gate:** Aqui é decidido quais informações devem ser esquecidas. A célula possui duas entradas,  $x_t$  que é a entrada no momento  $t$  e  $h_{t-1}$  que é a saída da célula anterior. As entradas são alimentadas ao *gate* e passam por uma função de ativação, gerando uma saída binária. Se para um determinado estado de célula a saída for 0, a informação é esquecida, para saída 1 a informação é retida.
- **Input Gate:** No *input gate* novas informações são adicionadas ao estado da célula. De maneira análoga ao *forget gate*, uma função *sigmoid* filtra os valores de entrada a serem

Figura 3 – Representação de uma arquitetura LSTM



Fonte: Disponível em: <http://deeplearningbook.com.br> Acesso em 3 de novembro de 2019

lembrados. Uma função *tanh* cria um vetor contendo todos os valores possíveis de  $h_{t-1}$  e  $x_t$ , os vetores são multiplicados a fim de obter novas informações úteis.

- **Output Gate:** No *output gate* é realizada a tarefa de retirar informações úteis do estado da célula para ser apresentado como saída. Primeiro um vetor é gerado aplicando uma função *tanh* no estado da célula, então a informação é regulada usando a função *sigmoid* que filtra os valores a serem lembrados usando as duas entradas  $x_t$  e  $h_{t-1}$ . Os valores regulados e o vetor da célula são multiplicados gerando a saída e entrada para a próxima célula.

Este trabalho utilizou, nos seus modelos, redes LSTM bidirecionais. Redes LSTM bidirecionais possuem a arquitetura da célula análoga a uma LSTM convencional, a diferença é que redes bidirecionais possuem uma segunda camada que passa informações de um instante  $t$  para um instante  $t - 1$ . Portanto, em redes bidirecionais o fluxo de informação recorrente acontece em dois sentidos, saindo do estado  $t$  para o estado  $t + 1$  e também para o estado  $t - 1$ .

## 2.3 Métricas de desempenho

Para analisar o desempenho da RNA na previsão dos valores são utilizadas métricas de que avaliam a diferença entre os valores previstos e os valores reais. O objetivo desta seção é listar as métricas de desempenhos que foram utilizadas no trabalho.

### 2.3.1 Erro percentual absoluto médio

O Erro Médio Percentual Absoluto (ou MAPE – mean absolute percentage error) obtém as diferenças percentuais entre todos os valores reais e previstos obtidos e faz uma média simples destes valores. Como todos os elementos da série temporal têm igual peso no resultado final, resultados isolados muito diferentes dos demais não fazem tanta diferença. Essa métrica é útil para ter uma visão geral do erro médio gerado pelo algoritmo de previsão escolhido. O MAPE é calculado conforme equação 2.4, na qual  $P_{real,i}$  é o valor real da série no dia  $i$  e  $P_{previsto,i}$  é o valor previsto para série, também no dia  $i$ .

$$MAPE = \frac{1}{N} \sum_{i=1}^N \left( \frac{|P_{previsto,i} - P_{real,i}|}{P_{real,i}} \right) * 100 \quad (2.4)$$

### 2.3.2 Erro quadrático médio

Essa é uma das métricas mais utilizadas para calcular o desempenho de modelos de previsão. É calculado a partir da soma da variância e dos quadrados das diferenças obtidas entre os valores reais e previstos, e é dado pela equação 2.5 onde  $N$  é o número de dias da série temporal sendo analisada,  $previsto_i$  é o valor previsto para a série no dia  $i$  e  $real_i$  é o valor real da série também no dia  $i$ .

$$MSE = \frac{1}{N} \sum_{i=1}^N (previsto_i - real_i)^2 \quad (2.5)$$

### 2.3.3 Raiz quadrada do erro quadrático médio

Essa métrica, também chamada de RMSE (Root Mean Squared Error). O adicional da raiz quadrada faz com que os erros estejam na mesma dimensão da variável analisada. O RMS é calculado a partir da equação 2.6, onde  $N$  é o número de dias da série temporal sendo analisada,  $previsto_i$  é o valor previsto para a série no dia  $i$  e  $real_i$  é o valor real da série também no dia  $i$ .

$$RMS = \sqrt{\frac{1}{N} \sum_{i=1}^N (previsto_i - real_i)^2} \quad (2.6)$$

### 2.3.4 Acurácia

Para avaliar a performance do modelo em acertar a direção do movimento foram utilizados duas métricas de acurácia. A acurácia em acertar a direção do movimento em relação ao dia anterior, e outra em relação a predição anterior. As acurácias são definidas em 2.7 e 2.8.

$$Acc_1 = \frac{1}{n} \sum_{i=1}^n P_{1i} \quad (2.7)$$

$$Acc_2 = \frac{1}{n} \sum_{i=1}^n P_{2i} \quad (2.8)$$

Onde  $P_{1i}$  e  $P_{2i}$  indicam a predição do movimento para o dia  $i$ .  $P_{1i}$ , definido em 2.9, indica a acurácia da previsão dos movimentos em relação ao dia anterior.  $P_{2i}$ , definido em 2.10, indica a acurácia da previsão dos movimentos em relação à última previsão, onde  $previsto_t$  é a previsão no instante  $t$  e  $real_t$  é o valor real no instante  $t$ .

$$P_i = \begin{cases} 1, (real_{t+1} - real_t)(previsto_{t+1} - real_t) \geq 0 \\ 0, (real_{t+1} - real_t)(previsto_{t+1} - real_t) < 0 \end{cases} \quad (2.9)$$

$$P_i = \begin{cases} 1, (real_{t+1} - real_t)(previsto_{t+1} - previsto_t) \geq 0 \\ 0, (real_{t+1} - real_t)(previsto_{t+1} - previsto_t) < 0 \end{cases} \quad (2.10)$$



### 3 Metodologia de pesquisa

Quanto a metodologia utilizada, primeiramente foi feita uma pesquisa exploratória a fim de compreender a área que o problema está inserido. Entender o problema de previsão em séries financeiras e o domínio do objeto de estudo foram o objetivo nessa etapa para entender o que pode influenciar na movimentação dos preços para uma melhor modelagem dos dados.

Posteriormente foi feita uma análise dos dados da série temporal para entender o comportamento e tentar modelar da melhor forma para realizar a predição dos valores. Por fim, foi realizada uma análise sobre os resultados dos modelos.

#### 3.1 Base de dados

A base de dados utilizada foi a série temporal das negociações de contratos de milho futuro, com dados de Maio de 2014 até Outubro de 2019. Como o milho futuro é negociado em contratos com vencimentos em determinados meses, cada mês de vencimento possui uma série. A série utilizada é uma junção das séries dos contratos de cada mês de vencimento. Os dados foram obtidos através do software *MetaTrader 5*. Cada registro da base de dados possui informações sobre a negociação diária do ativo, sendo o preço de abertura, o preço máximo do dia, o preço mínimo do dia, o preço de fechamento e o volume de negociação. A figura 4 mostra como os dados estão representados.

Figura 4 – Dados de negociação dos contratos de milho futuro

	<b>Abertura</b>	<b>Máxima</b>	<b>Mínima</b>	<b>Fechamento</b>	<b>Volume</b>
<b>Date</b>					
<b>2014-05-15</b>	28.6600	28.8600	28.4200	28.5700	3152
<b>2014-05-16</b>	28.6300	28.7600	28.5200	28.5500	1915
<b>2014-05-19</b>	28.4200	28.4400	28.2100	28.2100	2867
<b>2014-05-20</b>	28.2800	28.3000	27.9600	27.9700	2111
<b>2014-05-21</b>	27.8900	27.9000	27.7100	27.8100	2515

Fonte: Elaborada pelo autor.

## 3.2 Ferramentas utilizadas

Esta seção tem o objetivo de apresentar as ferramentas que foram utilizadas para o desenvolvimento do trabalho.

### 3.2.1 Python

Python é uma linguagem de programação de alto nível, orientada a objeto e *script*. É uma das linguagens mais utilizadas para ciência de dados e aprendizado de máquina pois possui diversas bibliotecas que auxiliam nessa área.

### 3.2.2 Scikit-learn

Scikit-learn é um biblioteca voltada para AM em código aberto para a linguagem Python. Ela possui diversos modelos e recursos para o ML em estado da arte já implementados através de classes e métodos. Informações sobre a estrutura da biblioteca podem ser encontradas em [Pedregosa et al. \(2011\)](#).

### 3.2.3 Pandas

Pandas é uma biblioteca em código aberto para a linguagem Python. Ela provê recursos para trabalhar em alta performance com estruturas de dados e com análise de dados. A figura 4 foi gerada a partir de um *dataframe* utilizando o pandas.

A biblioteca possui diversas funcionalidades que facilitam o trabalho com séries temporais. Toda a parte de modelagem e manipulação dos dados foi feita utilizando essa biblioteca. O projeto pandas é patrocinado pela organização NumFOCUS e sua documentação pode ser encontrada em [Pandas \(2019\)](#).

### 3.2.4 Matplotlib

Matplotlib é uma biblioteca para visualização de dados em 2D. A biblioteca possui funcionalidades para plotar dados em diversos tipos de gráficos de maneira simples. Matplotlib é uma biblioteca em código aberto para a linguagem Python.

### 3.2.5 Keras

Keras é uma *Application Programming Interface* API em alto nível, escrita em Python, para criação de modelos de RNAs. É capaz de rodar sobre *TensorFlow*, *Microsoft Cognitive*

*Toolkit* ou *Theano*, foi desenvolvida com o foco em agilizar a experimentação de modelos partindo da idealização para os resultados de maneira rápida, facilitando a pesquisa. Sua documentação pode ser encontrada em [Chollet et al. \(2015\)](#).

### 3.2.6 MetaTrader 5

A MetaTrader 5 é uma plataforma institucional multimercado para *trading*, análise técnica, uso de sistemas automáticos de negociação (robôs de negociação). Esse software foi utilizado para o levantamento dos dados, pois ele permite o *download* da série de preços de um ativo em arquivo no formato *csv* para estudos e criação de estratégias de operação.

## 3.3 Modelos Propostos

Para este trabalho foram propostas 3 arquiteturas de redes LSTM diferentes com três conjuntos de entradas diferentes, totalizando um total de 9 modelos para serem analisados. As arquiteturas diferem no tamanho da janela temporal que é passada como entrada para a rede.

Foram definidas arquiteturas com entradas de 2 janelas temporais, com 5 janelas e com 15 janelas. Isso para testar se o modelo performa melhor quando recebe dados de mais dias de negociação.

Os conjuntos de entrada foram separados em 3 grupos. Um formado somente pelos dados diários de negociação, outro formado pelos dados diários de negociação mais os indicadores técnicos apresentados em [2.1.5](#), e o terceiro formado pelos dados de negociação, os indicadores técnicos, e mais duas séries temporais exógenas. As séries exógenas são o indicador de preços do Milho Esalq/BM&FBOVESPA, apresentado em [2.1.3](#), e a série de histórica do índice DI diário.

O índice DI é um indicador financeiro que define a taxa de juros que será paga para investidores em aplicações feitas em instituições privadas. A pesquisa exploratória feita sobre objeto de estudo mostrou que a taxa de juros é um fator importante na determinação de preços no mercado futuro, por isso a escolha dessa série como variável de entrada.

As série do indicador Esalq/BM&FBOVESPA é disponibilizada pelo site do Centro de Estudos Avançados em Economia Aplicada ([CEPEA, 2019](#)), e os dados da série histórica do índice DI são disponibilizados pelo site da B3 ([B3, 2019 B](#)).

## 4 Desenvolvimento

A primeira etapa do desenvolvimento foi realizar uma análise sobre os dados para entender como eles se comportam.

### 4.1 Análise dos dados

Figura 5 – Série temporal do preço de fechamento fechamento



Fonte: Elaborada pelo autor.

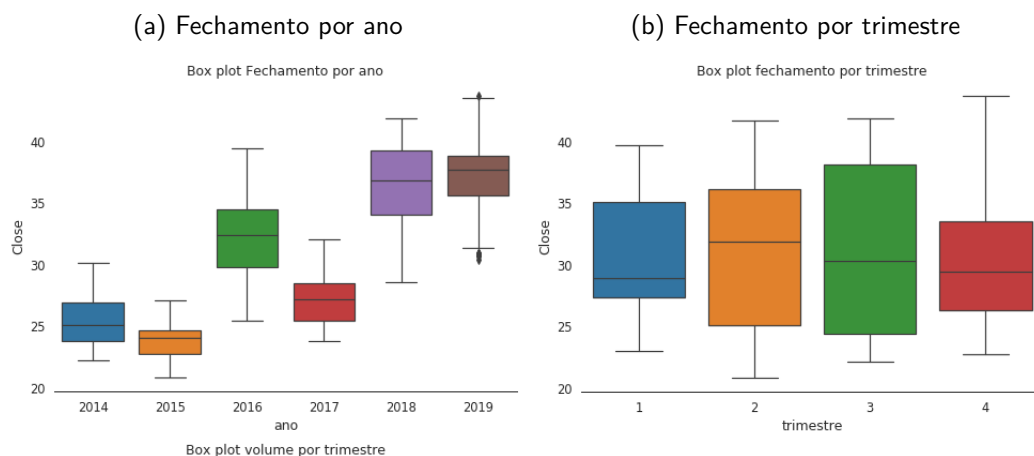
É possível notar que a série não é estacionária. Nota-se a presença de tendências nas movimentações do preço. A figura 6 mostra os preços de fechamento por ano e por trimestre em gráficos *box plot*. Fica claro ao olhar para os preços por ano que existe uma tendência longa de alta nos preços, mas com tendências anuais que oscilam entre alta e queda. No gráfico por trimestre é possível notar que o segundo e o terceiro trimestre do ano são os períodos que possuem maior amplitude de preços, isso indica que é nesse período que ocorre a inversão da tendência anual.

Na figura 7 mostra que o volume de negociações é maior no segundo e terceiro trimestre. O que faz sentido já que são os trimestres onde inverte a tendência anual.

A figura 8 mostra a distribuição dos preços. É possível perceber que a distribuição possui uma certa simetria, indicando que os preços oscilam dentro de uma certa faixa. O valor médio do preço de fechamento é 30.50 com desvio padrão de 5.83.

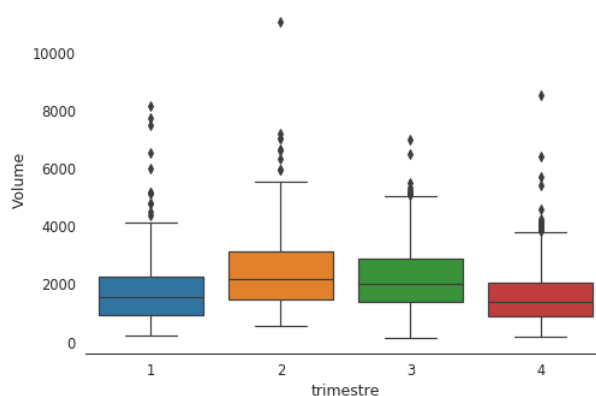
A análise mostra que os preços dos contratos de milho futuro possuem uma certa sazonalidade, principalmente em relação a tendências anuais. Entretanto não foi possível retirar *insights* a partir dos dados para a previsão do período de 1 dia.

Figura 6 – Box plot do fechamento por ano e por trimestre



Fonte: Elaborada pelo autor.

Figura 7 – Box plot do volume por trimestre



Fonte: Elaborada pelo autor.

## 4.2 Tratamento dos dados

### 4.2.1 Adição de novas variáveis

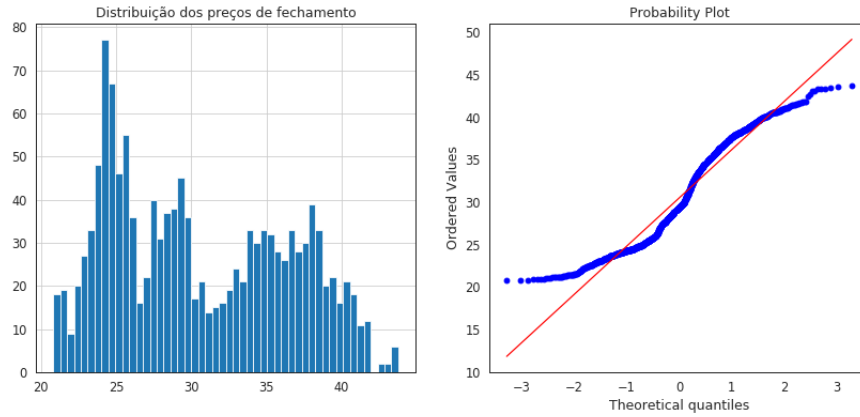
Para adicionar os indicadores técnicos ao conjunto de dados foram implementadas funções em Python que calculam os indicadores. Os indicadores foram calculados conforme apresentado em 2.1.5. As séries exógenas também foram adicionadas ao conjunto de dados. As manipulações nos dados foram feitas utilizando a biblioteca Pandas.

As série do indicador Esalq/BM&FBOVESPA possui duas variáveis, a do preço em reais (R\$) e do preço em dólares (US\$). Ao final dessa etapa tem-se um conjunto de dados com 1317 registros. A figura 9 mostra o estado final dos dados.

Figura 8 – Distribuição dos preços e plot de probabilidade.

(a) Distribuição dos preços

(b) Probabilidade de distribuição normal



Fonte: Elaborada pelo autor.

Figura 9 – Conjunto de dados após manipulações

	Abertura	Máxima	Mínima	Fechamento	Volume	Indicador Esalq R\$	Indicador Esalq U\$	Índice DI	MMS	MME	RSI	MACD
Date												
2014-07-03	24.1300	24.4600	24.0900	24.1900	2217	24.6600	11.1400	1.0004	24.8770	24.3624	34.7166	-0.0257
2014-07-04	23.9900	23.9900	23.8700	23.9200	1003	24.7300	11.1600	1.0004	24.7000	24.2149	36.1963	-0.0231
2014-07-07	23.8800	24.0900	23.7900	23.8500	2421	24.3900	10.9600	1.0004	24.4970	24.0933	41.0732	-0.0157
2014-07-08	23.7800	23.8700	23.6900	23.8000	1066	24.4800	11.0600	1.0004	24.3490	23.9955	43.1925	-0.0044
2014-07-10	23.7800	24.5100	23.6300	24.3000	5354	24.3500	10.9700	1.0004	24.2840	24.0970	41.3302	0.0439

Fonte: Elaborada pelo autor.

## 4.2.2 Pré-processamento

Esta seção apresentará quais foram as técnicas utilizadas no pré-processamento dos dados.

### 4.2.2.1 Normalização dos dados

Para que a rede possa aprender de maneira correta é necessário realizar uma normalização dos dados de entrada, fixando-os dentro de um intervalo. Durante os testes houve melhores resultados quando os dados foram normalizados no intervalo de -1 a 1. Cada coluna do conjunto de dados foi normalizado de -1 a 1 utilizando as equações 4.1 e 4.2

$$x_{std} = (x - x_{min}) / (x_{max} - x_{min}) \quad (4.1)$$

$$x_{norm} = (x_{std} * 2) - 1 \quad (4.2)$$

Onde  $x_{std}$  é a padronização entre 0 e 1,  $x_{min}$  é o valor mínimo da variável,  $x_{max}$  é o valor máximo da variável, e  $x_{norm}$  é o valor normalizado entre 1 e -1.

#### 4.2.2.2 Montagem das sequências

Uma LSTM recebe como entrada uma sequência de dados. Para cada modelo os dados foram ajustados formando uma matriz de 3 dimensões, sendo essas o número de amostras, o tamanho da sequência, e o número de variáveis de entrada.

#### 4.2.2.3 Treino, validação e teste

Para que seja possível avaliar o resultado do modelo, os dados foram separados em 3 conjuntos. O conjunto de treino contém 70% dos registros, e é utilizado para treinar a rede. O conjunto de validação compreende 20% dos dados, ele é utilizado para validar as parametrizações do modelo. Os 10% dos dados restantes forma o conjunto de teste, o qual é utilizado para verificar se o modelo não está superajustado (*overfitting*) para os dados de treino e validação.

Os dados foram separados nesses conjuntos mantendo a sua sequência temporal devido a necessidade do problema.

### 4.2.3 Treinamento dos Modelos

Os modelos propostos possuem 3 camadas, a camada de entrada, a camada LSTM bidirecional, e a camada de saída. A camada de entrada possui  $n$  nós de  $t$  dimensões, com  $n$  sendo o tamanho da sequência e  $t$  o número de variáveis de entrada. A camada LSTM, possui como dimensão de saída duas vezes o valor de  $t$ , e a camada de saída possui 1 dimensão que é a predição do preço de fechamento. O otimizador utilizado foi o RMSprop, e o loss foi o erro quadrático médio. Os treinamentos foram feitos para 2000 épocas com *stop* para quando o modelo para de convergir.

Em relação as funções de ativação, a que teve melhor resultado durante os testes foi a função *tanh* em todas as saídas. Não houve melhoras ao aumentar o espaço dimensional da saída da LSTM, por isso esse valor foi mantido como o mesmo número das variáveis de entrada, mas por ser uma rede bidirecional o valor é o dobro de  $t$ . Também não houve melhores resultados em adicionar mais camadas ao modelo.

### 4.3 Análise dos Resultados

A seguir tem-se os quadros com os resultados da aplicação dos modelos para cada conjunto de dados, treino, validação e teste. A legenda **conjunto** *n* indica o conjunto de variáveis de entrada, sendo que **conjunto 1** são apenas os dados de negociação (abertura, máxima, mínima, fechamento e volume), o **conjunto 2** representa os dados de negociação mais os indicadores técnicos adicionados (MM de 10dias, MME de 5 dias, histograma MACD e RSI). O **conjunto 3** é formado pelos mesmos dados do **conjunto 2** mais duas séries temporais (Indicador Esalq/BM&FBOVESPA, série DI diário).

A primeira linha dos quadros mostram o tamanho da sequência de cada modelo. A última coluna em cada quadro indica a métrica de avaliação. **Acc1** corresponde a acurácia da previsão das direções calculadas de acordo com a equação 2.10 e **Acc2** é a acurácia da previsão das direções calculadas de acordo com a equação 2.9.

A análise do quadro 1 mostra que para os dados de treino, modelo que teve melhor performance foi com uma janela de 5 dias com as variáveis de entrada do conjunto 3, com MAPE de 1,1406%. Para a previsão da direção do movimento baseado na última previsão, o melhor resultado foi obtido com as variáveis de entrada do conjunto 2, porém com uma diferença de apenas 0,1% para o modelo com as variáveis do conjunto 3.

Quadro 1 – Treino

	2 dias	5 dias	15 dias	
<b>conjunto 1</b>	1,2910%	1,6649%	1,5852%	<b>MAPE</b>
<b>conjunto 2</b>	1,2985%	1,2070%	1,5056%	<b>MAPE</b>
<b>conjunto 3</b>	1,7956%	<b>1,1406%</b>	1,4470%	<b>MAPE</b>
<b>conjunto 1</b>	0,4973	0,6090	0,5864	<b>RMSE</b>
<b>conjunto 2</b>	0,4892	0,4547	0,5515	<b>RMSE</b>
<b>conjunto 3</b>	0,6661	<b>0,4304</b>	0,5340	<b>RMSE</b>
<b>conjunto 1</b>	51,1401%	48,5342%	49,1857%	<b>Acc1</b>
<b>conjunto 2</b>	52,1173%	55,5917%	49,2942%	<b>Acc1</b>
<b>conjunto 3</b>	51,0315%	<b>59,6091%</b>	53,6374%	<b>Acc1</b>
<b>conjunto 1</b>	56,3043%	52,6087%	55,6522%	<b>Acc2</b>
<b>conjunto 2</b>	55,8696%	<b>56,9565%</b>	55,9783%	<b>Acc2</b>
<b>conjunto 3</b>	54,2391%	56,8478%	56,4130%	<b>Acc2</b>

Fonte – Autor.

O quadro 2 mostra que quanto ao erro na previsão, os modelos que obtiveram melhor performance foram com janelas de 5 e 15 dias com as variáveis de entrada do conjunto 1. Entretanto a previsão da direção dos preços comparando com o fechamento anterior se mostra ineficiente obtendo um valor significativamente acima de 50% somente com o conjunto de entrada 3. Para a previsão do movimento em relação a última previsão o melhor resultado



foi com uma janela de 15 dias com o conjunto de entradas 2, obtendo uma acurácia de 58%, mas com uma diferença de aproximadamente apenas 1,5% para o modelo com o conjunto de entradas 1.

Quadro 2 – Validação

	2 dias	5 dias	15 dias	
<b>conjunto 1</b>	1,4362%	1,2150%	<b>1,1988%</b>	<b>MAPE</b>
<b>conjunto 2</b>	1,4495%	1,3669%	1,8553%	<b>MAPE</b>
<b>conjunto 3</b>	1,7582%	2,2063%	2,0011%	<b>MAPE</b>
<b>conjunto 1</b>	0,6952	<b>0,5705</b>	0,5927	<b>RMSE</b>
<b>conjunto 2</b>	0,7284	0,6753	0,8916	<b>RMSE</b>
<b>conjunto 3</b>	0,8370	1,0364	0,9796	<b>RMSE</b>
<b>conjunto 1</b>	49,8099%	49,0494%	51,7110%	<b>Acc1</b>
<b>conjunto 2</b>	50,1901%	48,2890%	47,5285%	<b>Acc1</b>
<b>conjunto 3</b>	48,2890%	<b>53,9924%</b>	50,9506%	<b>Acc1</b>
<b>conjunto 1</b>	53,8168%	56,1069%	57,6336%	<b>Acc2</b>
<b>conjunto 2</b>	55,8696%	54,5802%	<b>58,0153%</b>	<b>Acc2</b>
<b>conjunto 3</b>	50,3817%	55,7252%	54,9618%	<b>Acc2</b>

Fonte – Autor.

O quadro 3 mostra as métricas aplicadas às previsões do conjunto de teste. Os resultados mostram os modelos que obtiveram a melhor performance foram utilizando as variáveis de entrada do conjunto 1 e com uma janela de 15 dias. Assim como nos dados de validação, o erro foi menor utilizando menos variáveis de entrada. Quanto a previsão da direção do movimento, a previsão em relação ao fechamento anterior mais uma vez se tornou ineficiente, obtendo um resultado significativamente acima de 50% em apenas 2 modelos. A melhor previsão em relação a movimentação dos preços foi obtida pelo modelo com variáveis de entrada conjunto 1 e uma janela de 5 dias, com uma acurácia de 60%. Isso indica que o modelo consegue acompanhar a direção da movimentação independente do erro no valor da previsão do fechamento.

A análise geral do resultado mostra que um modelo com mais variáveis de entrada se ajustam melhor aos dados de treino, mas não conseguem generalizar bem para outros dados. Os melhores resultados para dados que não foram utilizados no treinamento foram utilizando como variáveis de entrada apenas os dados de negociação dos contratos.

Quanto ao tamanho da sequência, nenhum dos resultados obtidos com uma janela de 2 dias foi superior a janelas maiores, indicando que existe uma dependência temporal em relação ao preço. O modelo com janela de 5 dias se ajustou melhor para os dados de treinamento, mas conforme os dados da série ficam mais distantes dos dados usados para treino o modelo com 15 dias passa a performar melhor. Outro ponto interessante é que os melhores resultados obtidos em relação a previsão da movimentação foram obtidos no conjunto de teste, onde o erro do valor da previsão é maior.

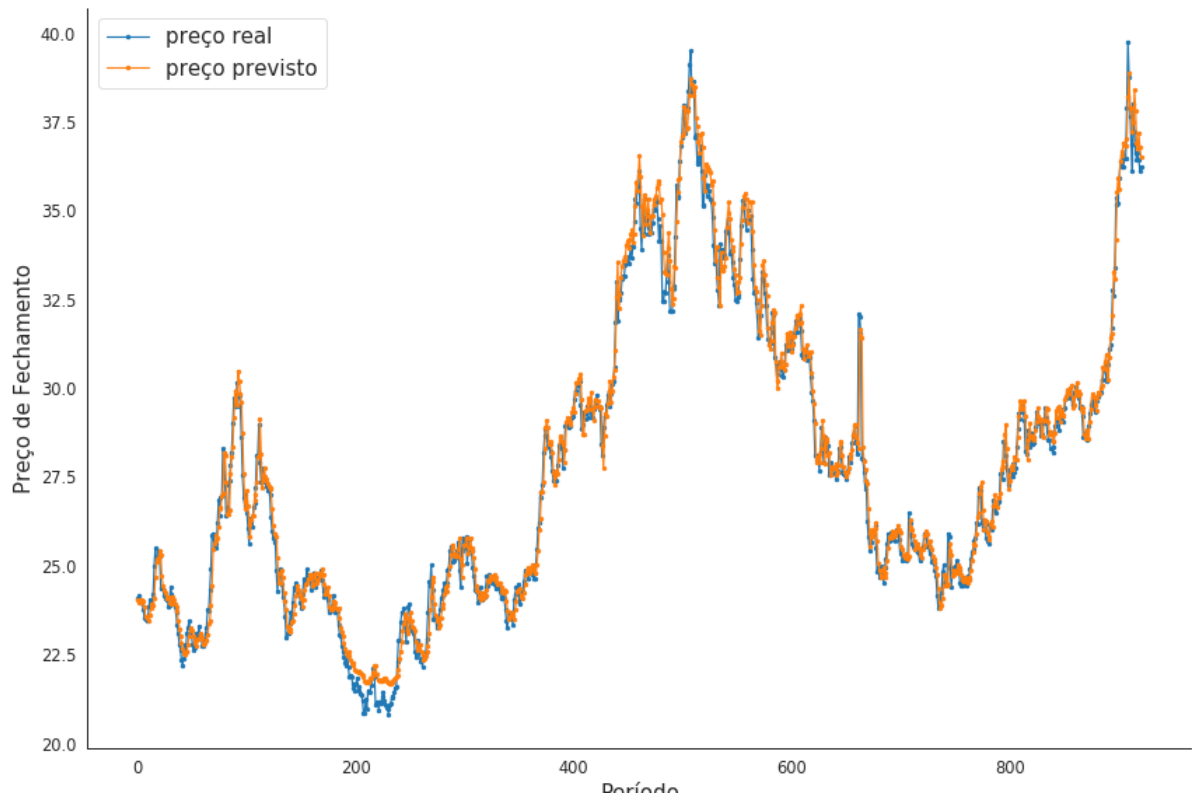
Quadro 3 – Teste

	<b>2 dias</b>	<b>5 dias</b>	<b>15 dias</b>	
<b>conjunto 1</b>	1,6879%	1,6439%	<b>1,4800%</b>	<b>MAPE</b>
<b>conjunto 2</b>	1,8488%	1,7318%	2,1861%	<b>MAPE</b>
<b>conjunto 3</b>	2,1197%	2,0139%	2,1028%	<b>MAPE</b>
<b>conjunto 1</b>	0,9641	0,8559	<b>0,8035</b>	<b>RMSE</b>
<b>conjunto 2</b>	0,8711	0,8829	1,0112	<b>RMSE</b>
<b>conjunto 3</b>	1,0693	1,0046	1,0213	<b>RMSE</b>
<b>conjunto 1</b>	54,5455%	48,0620%	<b>54,6218%</b>	<b>Acc1</b>
<b>conjunto 2</b>	44,6970%	50,3876%	45,3782%	<b>Acc1</b>
<b>conjunto 3</b>	52,2727%	44,1860%	49,5798%	<b>Acc1</b>
<b>conjunto 1</b>	55,7252%	<b>60,1562%</b>	57,6271%	<b>Acc2</b>
<b>conjunto 2</b>	56,4885%	57,8125%	55,9322%	<b>Acc2</b>
<b>conjunto 3</b>	54,9618%	51,5625%	56,7727%	<b>Acc2</b>

Fonte – Autor.

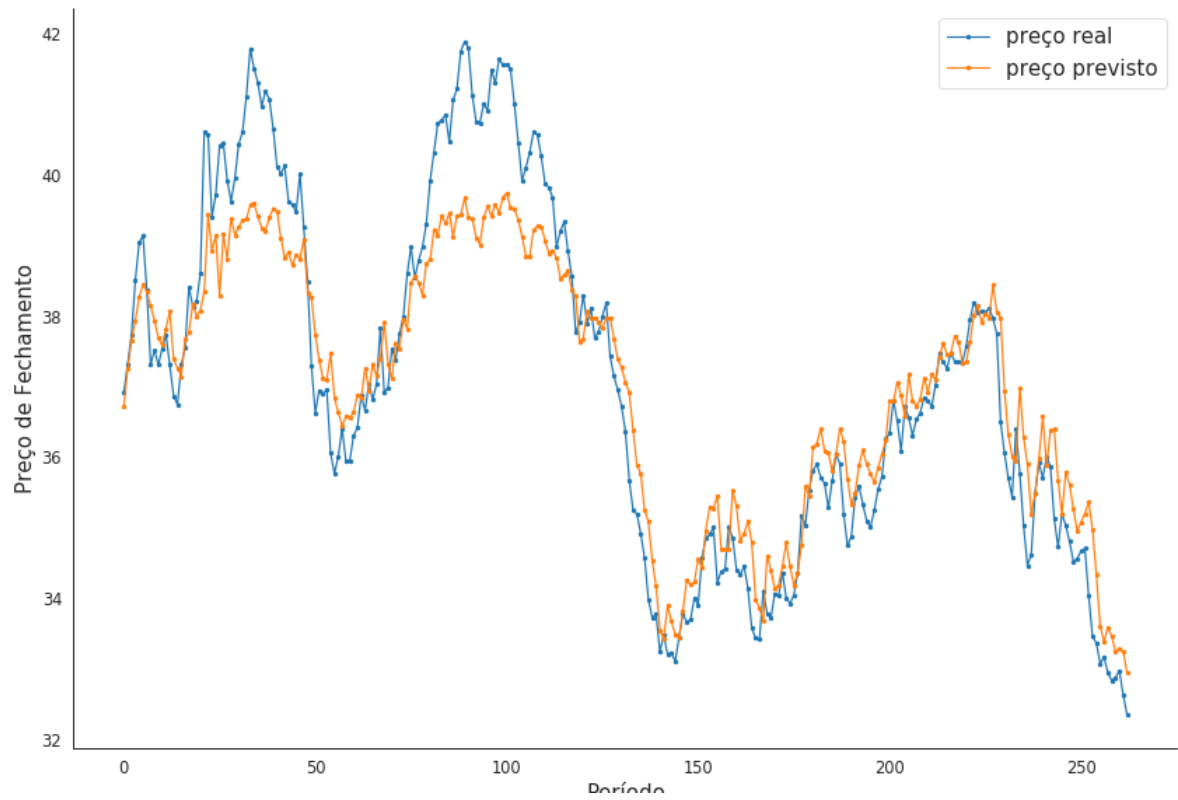
No geral o modelo que melhor performou para os conjuntos de validação e teste foi o modelo com janela de 15 dias e utilizando apenas os dados de entrada do conjunto 1. A figura 10 mostra a previsão deste modelo para os dados de treinamento, a figura 11 a previsão para os dados de validação, e a figura 12 a previsão para os dados de teste. A figura 13 mostra a convergência do modelo durante o treinamento através do decaimento da função de *loss*.

Figura 10 – Previsão dos preços para os dados de treino



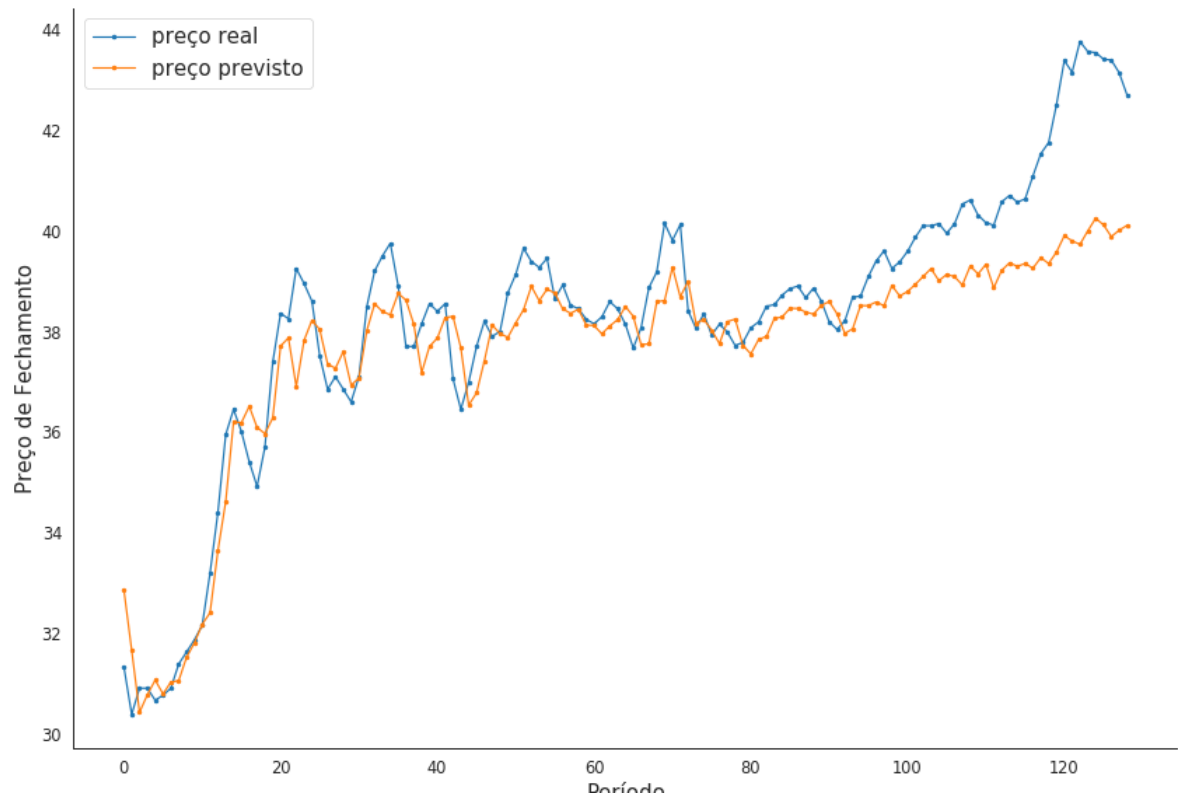
Fonte: Elaborada pelo autor.

Figura 11 – Previsão dos preços para os dados de validação



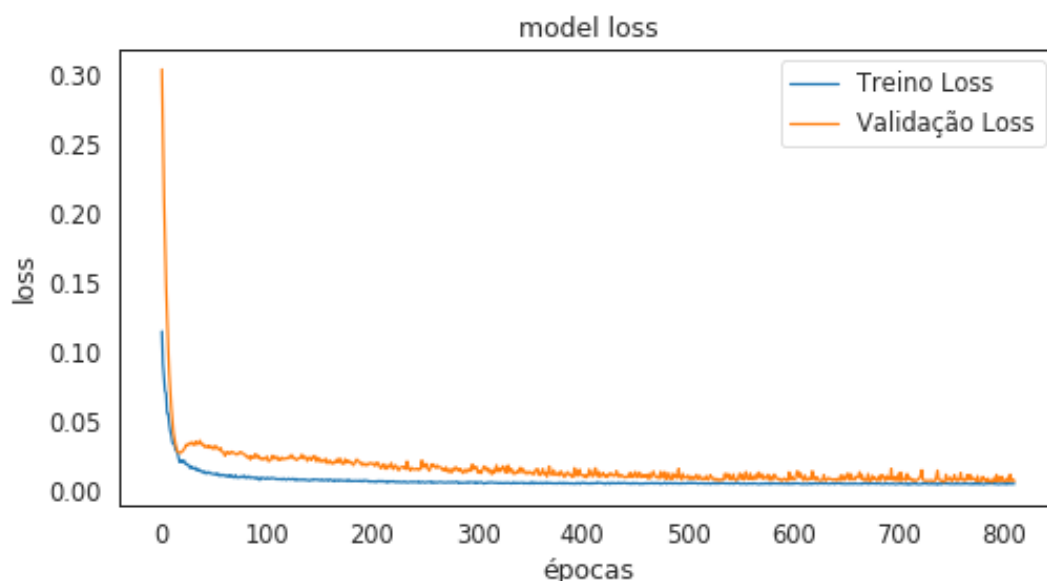
Fonte: Elaborada pelo autor.

Figura 12 – Previsão dos preços para os dados de teste



Fonte: Elaborada pelo autor.

Figura 13 – Convergência do modelo



Fonte: Elaborada pelo autor.

## 5 Conclusão

O trabalho proposto buscou realizar um estudo sobre o uso de redes neurais LSTM em previsões de valores em séries temporais financeiras e a relação temporal entre os preços nessas séries. Primeiro foi feito um estudo sobre o tema do mercado financeiro e sobre os contratos futuros para compreender o domínio do problema. Foram levantadas outras pesquisas realizadas na área de previsão desse tipo de série temporal para serem utilizadas como base teórica.

Foi realizado um estudo sobre o contrato de milho futuro, que foi o objeto de estudo para as previsões, a fim de entender o comportamento dos dados e tirar *insights* para a construção do modelo. Finalmente foi proposto um estudo utilizando 9 modelos de previsão com diferentes tamanhos de sequências de entrada, e diferentes variáveis de entrada, para avaliar o impacto dos preços anteriores no preço futuro e a influência das variáveis de entrada para a previsão.

Os resultados mostraram que de fato existe uma dependência temporal em relação a série de preços. Modelos com maiores janelas de tempo performaram melhor nos conjuntos de validação e teste, e o modelo com menor a menor janela de tempo não obteve a melhor performance em nenhum dos casos. Quanto as variáveis ficou nítido que inserir indicadores técnicos ajudam o modelo a se ajustar no conjunto de treinamento, mas resulta em um modelo menos generalista que realiza previsões piores para dados desconhecidos.

Para especulação no mercado financeiro o interessante é acertar a movimentação do preço. Os resultados mostraram que um erro menor na previsão do preço não garante o acerto na previsão da direção do movimento. Pode-se concluir que a abordagem de previsão do valor do preço não é a melhor no cenário de prever o comportamento do ativo, pois prever corretamente a direção do movimento mesmo que com um erro maior em relação ao preço, traz muito mais retorno para o investidor. Sugere-se outras abordagens como trabalhar com a série de retornos ou tratar como um problema de classificação de alta ou queda do preço.

# Referências

- ABE, M. *Manual de Análise Técnica*. São Paulo, SP: Novatec, 2017.
- ACADEMY, D. S. *Deep Learning Book*. 2019. Acesso em 3 novembro 2019. Disponível em: <http://www.deeplearningbook.com.br/>.
- ANBIMA. *Raio X do investidor brasileiro*. [S.l.], 2019.
- ANDREZO, A. F.; LIMA, I. S. *Mercado financeiro: aspectos conceituais e históricos*. [S.l.]: Atlas, 2007.
- B3. *Histórico de pessoas físicas*. [S.l.], 2019 A. Acesso em 29 Outubro 2019. Disponível em: [http://www.b3.com.br/pt\\_br/market-data-e-indices/servicos-de-dados/market-data/consultas/mercado-a-vista/historico-pessoas-fisicas/](http://www.b3.com.br/pt_br/market-data-e-indices/servicos-de-dados/market-data/consultas/mercado-a-vista/historico-pessoas-fisicas/).
- B3. *Site B3*. 2019 B. [http://www.b3.com.br/pt\\_br/](http://www.b3.com.br/pt_br/). Acessado: 31 out 2019.
- CEPEA. *Centro de Estudos Avançados em Economia Aplicada*. 2019. Acessado em 3 Novembro 2019. Disponível em: <https://www.cepea.esalq.usp.br/>.
- CHOLLET, F. et al. *Keras*. 2015. <https://keras.io>.
- DAMETTO, R. C. *Estudo da aplicação de redes neurais artificiais para predição de séries temporais financeiras*. 2018. Acesso em 29 Outubro 2019. Disponível em: <http://hdl.handle.net/11449/157058>.
- ELDER, A. *Trading for a Living: Psychology, Trading Tactics, Money Management*. 1. ed. Wiley, 1993. ISBN 9780471592242,0471592242. Disponível em: <http://gen.lib.rus.ec/book/index.php?md5=A4340174EA4725C1F9A7AC343F258B2B>.
- FLEURIET, M.; GALVAO, A.; MENDES, L. *Mercado Financeiro. Uma Abordagem Prático Dos Principais Produtos E Serviços*. [s.n.], 2005. ISBN 978-85-352-1336-2. Disponível em: <http://gen.lib.rus.ec/book/index.php?md5=f99ad81dee581a73ddfa8f3bb85d5025>.
- GIACOMEL, F. d. S. Um método algorítmico para operações na bolsa de valores baseado em ensembles de redes neurais para modelar e prever os movimentos dos mercados de ações. 2016.
- HOCHREITER, S.; SCHMIDHUBER, J. Long short-term memory. *Neural computation*, MIT Press, v. 9, n. 8, p. 1735–1780, 1997.
- HULL, J. C. *Opções, Futuros e outros Derivativos*. 9. ed.. ed. Porto Alegre: Bookman, 2016. P.1-2.
- KELLEHER, J. D.; NAMEE, B. M.; D'ARCY, A. *Fundamentals of Machine Learning for Predictive Data Analytics*. The MIT Press, 2015. ISBN 9780262331746. Disponível em: <http://gen.lib.rus.ec/book/index.php?md5=1b9fcc3d26bf8ba9d1e905a2ba4fe21d>.
- PANDAS. *Pandas*. 2019. Acessado em 3 de Novembro de 2019. Disponível em: <https://pandas.pydata.org/index.html>.

PEDREGOSA, F.; VAROQUAUX, G.; GRAMFORT, A.; MICHEL, V.; THIRION, B.; GRISEL, O.; BLONDEL, M.; PRETTENHOFER, P.; WEISS, R.; DUBOURG, V.; VANDERPLAS, J.; PASSOS, A.; COURNAPEAU, D.; BRUCHER, M.; PERROT, M.; DUCHESNAY, E. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, v. 12, p. 2825–2830, 2011.