

ARTIGO ORIGINAL

## Implementação e comparação de técnicas de machine learning aplicadas à predição do desenvolvimento de populações de afídeos

## Implementation and comparison of machine learning techniques applied to predict the development of aphid populations

Alexandre Tagliari Lazzaretti<sup>1</sup>, Vinicius Rafael Schneider<sup>1</sup>, Roberto Wiest<sup>1</sup>, Douglas Lau<sup>2</sup>, José Maurício C. Fernandes<sup>2</sup>, Clyde W. Fraisse<sup>3</sup>, Vinícius Andrei Cerbaro<sup>3</sup>, Maurício Z. Karrei<sup>3</sup>

<sup>1</sup>Instituto Federal Sul-Riograndense - Passo Fundo - RS - Brasil, <sup>2</sup>Embrapa Trigo - Passo Fundo - RS - Brasil,

<sup>3</sup>Universidade da Flórida - Gainesville - FL - EUA

\*[alexandrelazzaretti@ifsul.edu.br](mailto:alexandrelazzaretti@ifsul.edu.br); [vinirafaelsch@gmail.com](mailto:vinirafaelsch@gmail.com); [robertowiest@ifsul.edu.br](mailto:robertowiest@ifsul.edu.br); [douglas.lau@embrapa.br](mailto:douglas.lau@embrapa.br); [mauricio.fernandes@embrapa.br](mailto:mauricio.fernandes@embrapa.br); [cfraisse@ufl.edu](mailto:cfraisse@ufl.edu); [cerbaro@ufl.edu](mailto:cerbaro@ufl.edu); [mauricio.zientar@ufl.edu](mailto:mauricio.zientar@ufl.edu)

Recebido: 25/04/2022. Revisado: 04/04/2023. Aceito: 17/10/2023.

### Resumo

Os insetos possuem importante grau de colaboração para a manutenção do ecossistema no planeta. Porém ao atingir um determinado nível populacional e causar danos às plantas, alguns insetos, passam a ser considerados como insetos-pragas e representam uma ameaça para a agricultura. O afídeo ou pulgão, é um inseto que contém características para atingir este estado pois apresenta um alto potencial biótico e pode causar diferentes tipos de dano às plantas. Fatores meteorológicos como precipitações, ventos e temperaturas interferem no crescimento populacional destes insetos. Portanto, este trabalho se propõe a aplicar a aplicar diferentes técnicas de machine learning com o objetivo de verificar a correlação existente entre variáveis climáticas e a dinâmica populacional dos afídeos. Pode-se concluir que variáveis como precipitação, temperatura, quantidade de dias com chuva na semana e fenômenos climáticos como El niño e La niña possuem influência na população de afídeos. Durante o trabalho, foram implementados 4 (quatro) modelos e aplicados aos dados existentes de população de afídeos com objetivo de avaliar a melhor acurácia. Como resultado obteve-se as acurácias: 11,4% para Regressão Linear; 26,4% para o modelo de Rede Neural Artificial; 29,3% para Árvore de decisão e 41,4% para random forest.

**Palavras-Chave:** Análise de Dados; Árvore de decisão; Extração de conhecimento; Random forest; Redes Neurais Artificiais; Regressão Linear.

### Abstract

Insects have an important degree of collaboration for the maintenance of the ecosystem on the planet. However, after reaching a certain population level and causing damage to plants, some insects are considered as pests and represent a threat to agriculture. Aphids insects that has characteristics to reach this state as it has a high biotic potential and can cause different types of damage to plants. Climatic data as precipitation, winds and temperatures affect the population quantity of these insects. Therefore, this work proposes to apply different machine learning techniques with the objective to verify the existing correlation between climatic variables and the population dynamics of aphids. It can be concluded that variables such as precipitation, temperature, number of days when it rains in the week and climatic phenomena such as El niño and La niña have an influence on the aphid population. During the work, four models were developed in order to predict the population of these insects. The accuracy of the prediction model developed were 11.4% for Linear Regression; 26.4% for the Artificial Neural Network model; 29.3% for Decision Tree and 41.4% for Random Forest.

**Keywords:** Artificial neural networks; Decision tree; Exploratory Data; Knowledge extraction; Linear Regression; Random Forest.

## 1 Introduction

A maioria das espécies de insetos é benéfica ou útil ao homem. Por executarem funções como polinização de plantas, decomposição de matéria orgânica, ter participação ativa no equilíbrio biológico, entre outros, os insetos acabam se tornando importantes na manutenção dos ecossistemas (Finkler, 2013).

Porém, segundo Gassen (1984) alguns insetos são considerados pragas quando alcançam níveis populacionais que possam causar danos às plantas, gerando assim redução no rendimento de grãos e, conseqüentemente, diminuição na lucratividade. Nesse sentido, o autor cita ainda que, ao atingir o estado de praga, compensa ao agricultor o uso de métodos de controle para estes insetos.

O afídeo é um tipo de inseto que apresenta potencial para atingir o estado de praga quando alojado em determinada área. Wiest et al. (2021) expõe a existência de aproximadamente 4.700 espécies classificadas na família Aphididae, que é a família que agrupa os afídeos.

Em um estudo envolvendo diversas variáveis climáticas como precipitação, temperatura, ventos, Oliveira (1971) chegou à conclusão de que de uma forma ou de outra, com maior ou menor intensidade, os fatores climáticos tiveram o seu grau de atuação na diminuição ou aumento de pulgões.

Outro fator que possui correlação com a dinâmica populacional dos afídeos é o fenômeno ENOS (El Niño-Oscilação do Sul), pois conforme citado em Cunha et al. (2001), a ocorrência deste fenômeno acaba alterando o padrão de circulação geral da atmosfera, o que acaba tendo influência no comportamento das variáveis climáticas.

Kamilaris et al. (2017) apontam que a análise de dados possibilita às empresas e agricultores melhora na produtividade de suas lavouras através da extração de conhecimento ou detecção de padrões. A ciência de dados, tem como objetivo gerar valor através de estudos voltados para uma grande quantidade de dados. Fazendo uso de disciplinas como mineração de dados, estatística, banco de dados, entre outras. Esta área torna possível realizar inferências, gerar gráficos e detectar padrões sobre um determinado assunto (Dhar, 2012).

Sendo assim, seria interessante o uso de ciência de dados sobre base de dados coletados em lavouras e base de dados climáticos, no intuito de extrair conhecimento e gerar inferências baseando-se na correlação existente entre variáveis climáticas e dinâmica da população de afídeos, de forma a auxiliar entidades responsáveis na realização do controle populacional destes insetos.

Para isso, em um primeiro momento, foi organizada uma base de dados contendo dados climáticos para as regiões de monitoramento de afídeos. O que tornou possível verificar a correlação entre variáveis climáticas e a quantidade populacional dos afídeos. Foram avaliados os impactos causados pelos fenômenos El Niño e La Niña sob a população destes insetos e detectados períodos do ano em que ocorrem aumento ou diminuição na quantidade de afídeos coletados pelas armadilhas. Por fim, o desenvolvimento de modelos de predição para população de afídeo, com base nas variáveis e fenômenos climáticos abordados no trabalho.

O artigo está estruturado da seguinte maneira: a seção 2

exibe o referencial bibliográfico. Na seção 3 são apresentados o desenvolvimento do trabalho e os resultados obtidos. A seção 4 contém as considerações finais.

## 2 Referencial bibliográfico

### 2.1 Afídeos

Os Afídeos pertencem à família Aphididae, superfamília Aphidoidea, subordem Sternorrhyncha, ordem Hemiptera. Esse inseto possui formato ovalado, coloração variável e o tamanho máximo em que pode alcançar é de 5 milímetros de comprimento (Toebe, 2014).

Segundo Stern (2008), os afídeos, também conhecidos como pulgões, possuem um corpo mole e pequeno. Se alimentam por meio da inserção de parte do seu intestino delgado na área responsável por conduzir o alimento para demais partes da planta, com o objetivo de usufruir dos nutrientes dispostos pela seiva da planta.

Toebe (2014) cita que os danos causados por estes insetos nas plantas podem ser classificados em dois tipos, direto e indireto. O dano direto acontece quando o afídeo se alimenta da seiva da planta e também através da inserção de toxinas presentes em sua saliva. O dano indireto acontece pela transmissão de vírus, como por exemplo o Vírus do Nanismo Amarelo da Cevada (VNAC), que adoece a planta, causando amarelamento em suas folhas e diminuição em seu tamanho. Segundo Auad et al. (2002) as plantas que são infestadas por esses homópteros apresentam folhas enroladas, encarquilhadas e raquíticas, sendo que a secreção açucarada excretada pelos mesmos reduz consideravelmente o valor comercial do produto.

Os afídeos são classificados em três categorias durante sua vida: a fase inicial do afídeo intitula-se ninfa; ao entrar na fase adulta, podem variar entre duas classes, desta forma, são classificados como ápteros, que são os adultos que não possuem asas ou alado, adultos que possuem asas (Stern, 2008).

A dinâmica populacional destes insetos pode ser afetada por fatores bióticos causados por parasitóides, predadores e patógenos e também por fatores abióticos, que ocorrem através de variáveis como temperatura, umidade, luminosidade, entre outras. Sendo que a temperatura é um fator que apresenta interferência direta com o desenvolvimento populacional dos afídeos por afetar a taxa de desenvolvimento, reprodução e sobrevivência, refletindo assim, sua densidade populacional em uma planta (Wiest et al., 2021).

### 2.2 Dados Climáticos

Segundo Torres et al. (2015), a medição de variáveis climáticas apresenta grande importância devido ao fato de que, diversas atividades humanas são direta e indiretamente afetadas por essas variáveis. Sendo assim, para suprir tal necessidade, criaram-se as estações meteorológicas de alta precisão e amostragem.

Na agricultura, o monitoramento automático dos elementos meteorológicos tem contribuído não somente para o aumento da produtividade como, também, para a melhoria da qualidade dos produtos e para a preservação dos recursos naturais (Fernandes et al., 2013).

### 2.3 Enos (El Niño–Oscilação do Sul)

Segundo [INPE \(2021\)](#), os fenômenos El niño e La niña fazem parte do fenômeno atmosférico–oceânico que ocorre no oceano Pacífico Equatorial, denominado El Niño Oscilação Sul (ENOS). Quando a temperatura das águas oceano Pacífico Equatorial estão mais quentes que média normal histórica, este evento é conhecido como El niño e quando o oceano Pacífico Equatorial está mais frio do que a média normal histórica, ocorre o La niña.

Se tratando da região sul do Brasil, em anos onde o fenômeno El niño ocorre existe um aumento na quantidade e ocorrência de precipitações, enquanto em anos de La niña, existe uma diminuição nos valores normais de chuva ([Cunha et al., 2001](#)).

[Rosenzweig et al. \(2001\)](#) citam que as condições climáticas durante os anos de El niño e La niña possuem correlação com danos causados por pragas em algumas regiões.

### 2.4 Banco de Dados

Um banco de dados é definido como um conjunto de dados interrelacionados, livre de redundância, que atende às necessidades de um determinado domínio de aplicação ou de um conjunto de usuários.

Enquanto o SGBD (Sistema de Gerência de Banco de Dados) é uma ferramenta que torna possível a manipulação desses dados armazenados no banco de dados ([Navathe, 2018](#)).

### 2.5 Matriz de Correlação

A Matriz de Correlação é utilizada para iniciar a análise estatística dos dados, através dela é possível identificar, de forma visual, as variáveis estudadas que se relacionam entre si ([Lidiane et al., 2018](#)). O cálculo utilizado para matriz de correlação é o Coeficiente de correlação de Spearman, onde são avaliados a intensidade e o sentido da relação entre duas variáveis.

Este cálculo resulta em um valor entre  $-1$  e  $1$  responsável por expressar o grau de dependência entre duas variáveis. Quando negativas, significa que uma variável diminui com o aumento da outra, caso positiva, significa que uma variável aumenta com o aumento da outra.

### 2.6 Regressão Linear

Segundo [Rodrigues \(2013\)](#), a análise de regressão busca analisar dados com o objetivo de saber se, e como duas ou mais variáveis relacionam-se. Esta análise tem como resultado uma equação matemática que descreve o relacionamento entre essas variáveis.

[Rodrigues \(2013\)](#) expõem ainda que esta análise “pode ser usada para estimar ou prever, valores futuros de uma variável quando se conhecem ou se supõem conhecidos valores da outra variável”.

### 2.7 Redes Neurais Artificiais

Conforme [Teixeira Martins et al. \(2019\)](#), é um sistema computacional paralelo constituído de uma ou mais unidades simples de processamento, também conhecidos como neurônios, que conectados entre si são encarregados de calcular determinadas funções matemáticas que, normalmente, são não-lineares.

Os neurônios são dispostos em uma ou mais camadas interligadas por um grande número de conexões, geralmente unidirecionais. Na maioria dos modelos as conexões são associadas a pesos, que armazenam o conhecimento representado no modelo e servem para ponderar a entrada recebida por cada neurônio da rede ([Teixeira Martins et al., 2019](#)).

### 2.8 Árvore de Decisão

A árvore de decisão pode ser utilizada como classificador ou preditor e possui duas etapas: a de aprendizado e a de classificação/predição. Na fase de aprendizado é gerada uma árvore de decisão com base no conjunto de amostras utilizado para o treinamento. Já a etapa de classificação ou regressão é quando a árvore faz uso da estrutura gerada na fase de aprendizado para classificar ou prever um determinado dado ([Lan et al., 2020](#)).

Segundo [Lan et al. \(2020\)](#), sua estrutura é parecida com a de um fluxograma, que pode ser dividido em três partes: os nós internos, as ramificações e o nós folhas. Cada um dos nós internos representam um teste gerado na fase de aprendizado da árvore, as ramificações correspondem a uma saída dos nós internos, gerando assim um caminho para um nó folha, que por sua vez representa a saída de dados, ou seja, uma classificação ou predição.

### 2.9 Random Forest

A Random Forest aproveita a técnica de árvore de decisões para realizar classificação de dados. ([Lan et al., 2020](#)) cita que a Random Forest combina várias árvores de decisões, intituladas como classificadores fracos, para assim formar um classificador mais forte. Cada árvore utiliza um determinado número de amostras aleatórias para realizar o treinamento, garantindo assim a aleatoriedade dos dados.

Para gerar a classificação ou predição dos dados, a Random Forest faz com que as árvores de decisão elejam, através de votação, a classe mais apropriada para uma entrada de dados ([Lan et al., 2020](#)).

### 2.10 Python

Python é uma linguagem de programação interpretada, orientada a objetos. Permite a construção de código em alto nível, é uma linguagem simples de aprender, suporta módulos e o uso de pacotes, que encoraja a utilização de módulos em sua programação, bem como o reuso de código ([Python.org, 2023](#)).

Esta linguagem de programação foi escolhida para o desenvolvimento do trabalho por apresentar natureza interativa de alto nível e um amadurecimento em seu ecossistema de bibliotecas científicas, facilitando a utilização



de algorítmicos de ciência de dados e a análise exploratória.

### 2.11 Scikit-learn

O Scikit-learn ou Sklearn, é uma biblioteca de aprendizado de máquina de código aberto projetada e escrita em Python (Scikit Learning Development Group, 2023). Essa biblioteca fornece a um conjunto de algoritmos que são utilizados no enfoque de aprendizado de máquina, de forma que até mesmo especialistas de outras áreas consigam fazer uso da mesma.

No scikit-learn, todos os objetos e algoritmos aceitam dados de entrada de matrizes bidimensionais na forma de registros e colunas. Os objetos Scikit-learn compartilham um conjunto uniforme de métodos que depende de sua finalidade: estimadores podem ajustar modelos de dados, preditores podem fazer previsões sobre novos dados e transformadores convertem dados de uma representação para outra.

A biblioteca inclui algoritmos de aprendizado de máquina, sendo eles do tipo supervisionado e não supervisionado. Alguns dos algoritmos presentes na biblioteca são: k-means, que faz uso da metodologia de clusters, e classifica objetos de acordo com a proximidade de suas características com as características de determinado cluster, random forest que cria uma árvore aleatória, onde cada nó representa uma característica de um objeto e os nós folhas representam um objeto cadastrado na base de dados, ou seja, uma saída/classificação.

### 2.12 Pandas

Pandas é uma biblioteca Python voltada para manipulação e análise de dados. Pandas Development Group (2023) cita que a biblioteca fornece estruturas e funções de dados ricas, projetadas para tornar o trabalho com dados estruturados rápido, fácil e expressivo.

A biblioteca fornece rotinas intuitivas integradas para a execução de manipulações e análises de dados comuns em tais conjuntos de dados". Através dela é possível indexar, separar, selecionar e agrupar subconjuntos de dados para tornar mais fácil a manipulação e análise dos mesmos.

## 3 Resultados

### 3.1 Estudo de Caso

Para o desenvolvimento deste trabalho, foi utilizada a base de dados proveniente do Trapsystem, que segundo Lazzaretti et al. (2016), é uma aplicação que "permite o gerenciamento de dados obtidos pela captura/coleta de insetos em armadilhas". Os dados desta aplicação foram cedidos pela Embrapa (Empresa Brasileira de Pesquisa Agropecuária) e possui a quantidade de afídeos coletados em armadilhas amarelas expostas a campo. A espécie destes pulgões que foi disponibilizada e será usada no estudo é *Rhopalosiphum padi* (Linnaeus, 1758) que segundo (CABI, 2022) é uma praga que apresenta ampla distribuição geográfica e abundância.

A base de dados de afídeos foi exportada em formato csv e, conforme Fig. 1, têm seus dados armazenados de

```
station; date; semana; id; description; desinsect; total
435;2017-01-05;1;3;Armadilha 0;Soma Agrupada de Afídeos;4
435;2017-01-12;2;3;Armadilha 0;Soma Agrupada de Afídeos;2
435;2017-01-19;3;3;Armadilha 0;Soma Agrupada de Afídeos;2
435;2017-01-26;4;3;Armadilha 0;Soma Agrupada de Afídeos;14
435;2017-02-02;5;3;Armadilha 0;Soma Agrupada de Afídeos;25
435;2017-02-09;6;3;Armadilha 0;Soma Agrupada de Afídeos;30
435;2017-02-16;7;3;Armadilha 0;Soma Agrupada de Afídeos;14
435;2017-02-23;8;3;Armadilha 0;Soma Agrupada de Afídeos;44
435;2017-03-02;9;3;Armadilha 0;Soma Agrupada de Afídeos;19
```

Figura 1: Dados de afídeos coletados em armadilhas amarelas expostas a campo.

```
[
  {
    "id": 127153849,
    "idEstacao": 9003000,
    "data": "2017-01-01",
    "temperaturaMinima": 18.08,
    "temperaturaMinimaDuvidosa": false,
    "temperaturaMinimaEstimada": true,
    "temperaturaMedia": 23.1,
    "temperaturaMediaDuvidosa": false,
    "temperaturaMediaEstimada": true,
    "temperaturaMaxima": 28.12,
    "temperaturaMaximaDuvidosa": false,
    "temperaturaMaximaEstimada": true,
    "precipitacao": 29.1,
    "precipitacaoDuvidosa": false,
    "precipitacaoEstimada": true,
    "editavel": false
  },

```

Figura 2: Dados climáticos do INMET.

forma semanal possuindo os seguintes campos: station, que armazena um identificador para uma determinada estação; date, referente a data em que o dado foi coletado; semana, semana do ano em que o dado foi coletado; id, refere-se a uma identificação de armadilha; description, campo para armazenar um nome para uma armadilha; desinsect, uma descrição da operação feita para chegar ao resultado da coluna total; por fim, total, responsável por armazenar a quantidade de afídeos encontrados na ocasião.

Outra base de dados utilizada no estudo é referente a dados climáticos que são armazenados e disponibilizados de forma pública pelo INMET (Instituto Nacional de Meteorologia).

A coleta dos dados climáticos é feita de forma diária e seus dados são disponibilizados em formato json, possuindo, conforme Fig. 2, os campos para armazenamento de: id, que armazena um identificador para os registros que ocorrem; idEstacao: campo para identificação da estação; data, armazena a data em que o dado é coletado; temperatura Mínima, temperaturaMedia e temperaturaMaxima, que representam, respectivamente, a temperatura mínima, média e máxima que aconteceu em um determinado dia; precipitacao, quantidade em milímetros de chuva que aconteceu em um dia.

As localidades e os períodos levados em consideração para realização deste trabalho podem ser vistos na Tabela 1.

Para verificar a correlação existente entre a ocorrência do fenômeno ENOS com a dinâmica populacional dos

**Tabela 1:** Localidades utilizadas e período cronológico dos Afídeos e os dados climáticos.

Localidade	Dados de Afídeos	Dados Climáticos
Coxilha	jan/2017 à dez/2020	jul/2018 à dez/2020
Passo Fundo	jan/2017 à dez/2020	jan/2017 à dez/2020

afídeos se fez necessária a descoberta dos períodos em que estes fenômenos aconteceram. A aquisição destes dados se deu através do site da NOAA (National Oceanic and Atmospheric Administration - (NOAA, 2022)), onde a temperatura do oceano é medida de forma mensal.

A fim de determinar que tipo de fenômeno aconteceu em um determinado ano foi levado em consideração a ocorrência dos fenômenos entre os períodos de agosto de um ano até julho do ano seguinte, assim como é feito e exposto no site do INPE (Instituto Nacional de Pesquisas Espaciais - (INPE, 2021)). A Tabela 2 expõe a ocorrência de fenômenos ENOS em seus respectivos períodos.

**Tabela 2:** Ocorrência de fenômenos ENOS e seus respectivos períodos.

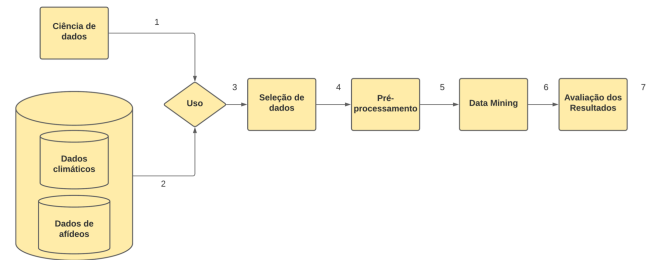
Período	Fenômeno ENOS
ago/2016 à jul/2017	La Niña
ago/2017 à jul/2018	La Niña
ago/2018 à jul/2019	El Niño
ago/2019 à jul/2020	Ano Neutro
ago/2020 à jul/2021	La Niña

### 3.2 Metodologia

Conforme Fig. 3, utilizando ciência de dados (passo 1) sobre as bases de dados (passo 2) a etapa de seleção de dados (passo 4) foi realizada, nela ocorreu a leitura dos dados climáticos e dados de afídeos. Na fase de pré-processamento (passo 5) ocorre a manipulação dos dados, nela são analisadas e tratadas possíveis falhas ou inconsistências. Foi desenvolvido um dataframe contendo dados de afídeos e dados climáticos para cada localidade e também um dataframe que concatena os dados correspondentes à todas as localidades estudadas. Também foram adicionadas duas novas colunas referentes a ocorrência do fenômeno ENOS e o número de dias em que choveu na semana para cada localidade.

A etapa de Data Mining (passo 6) é de exploração dos dados, onde ocorre a extração de conhecimento. Para isso foi feita a plotagem e análise de gráficos, utilização da matriz de correlação e o uso de técnicas de aprendizado de máquina com o objetivo de prever a quantidade de afídeos com base nas variáveis climáticas abordadas.

Ainda se tratando do passo 6, na análise exploratória as localidades estudadas foram analisadas de forma individual, isso foi feito para verificar se a correlação entre variáveis climáticas e quantidade de afídeos se faz verdadeira e não seja algo ocasional. Para as demais etapas, foi feita a concatenação dos dados, unindo dados referentes a todas as localidades em um único dataframe. Por fim, a última etapa é a de avaliação de resultados (passo 7), onde foi avaliada a existência de correlação entre dados climáti-

**Figura 3:** Modelo conceitual que representa a metodologia usada no trabalho.

cos e dados de afídeos e os resultados obtidos com modelos de predição desenvolvidos.

### 3.3 Seleção de Dados

A leitura dos arquivos contendo os dados foi feita utilizando funções disponibilizadas pela biblioteca Pandas, que permite executar a leitura de arquivos com extensões csv e json, transformando-os em um novo dataframe.

Sendo assim, obteve-se como resultado duas novas estruturas de dataframe para cada localidade, sendo que uma delas é referente aos dados de afídeos e outra tratando-se dos dados climáticos.

### 3.4 Pré-Processamento

Na etapa de pré-processamento foram removidas as colunas julgadas desnecessárias para o estudo dos dados e o treinamento dos modelos. Nos dados de afídeos foram eliminados os campos: id, estacao, descricao e desinsect. Para os dados climáticos, julgou-se necessária a permanência de todas as variáveis. Como os dados de afídeos são armazenados de forma semanal, optou-se por transformar o dataframe de dados climáticos em dados semanais. Para isso, foi feita uma média semanal dos campos temperaturaMinima, temperaturaMedia e temperaturaMaxima. No campo precipitacao, foi feita a soma do total de chuva que ocorreu na semana.

Para tornar possível as operações de média e soma, o primeiro passo foi gerar duas novas colunas no dataframe dos dados climáticos, uma contendo a semana do ano e outra contendo o ano em que o registro ocorreu. Isso foi feito utilizando a Pandas, que possibilita extrair essas informações a partir do campo que armazena a data de cada registro.

Por meio das colunas criadas e agrupamento de dados, foi possível executar operações de média e soma sobre os dados. A partir disto, foi criado um novo dataframe, que armazena os dados semanais com os campos de: ano, semana, temperaturaMinima, temperaturaMedia, temperaturaMaxima e precipitacao.

Pelo fato de algumas localidades conterem mais de uma armadilha, foi feita a média de ocorrência em cada semana para cada localidade, posterior a isso, estes dados foram concatenados ao novo dataframe utilizando laços de repetição, onde os mesmos são atribuídos a uma determinada linha caso a semana e o ano forem iguais nos dois datafra-

ano	semana	temperaturaMinima	temperaturaMedia	temperaturaMaxima	precipitacao	frqChuva	fenomeno	afideos	
0	2017	1	19.63	24.23	28.82	39.8	6.0	1.0	4.0
1	2017	2	18.13	23.35	28.56	23.2	4.0	1.0	2.0
2	2017	3	17.89	23.24	28.58	16.3	2.0	1.0	2.0
3	2017	4	15.50	21.10	26.70	19.0	2.0	1.0	14.0
4	2017	5	17.96	22.41	26.86	22.9	7.0	1.0	25.0
...	...	...	...	...	...	...	...	...	...
173	2020	47	12.55	19.50	26.45	6.3	5.0	1.0	34.0
174	2020	48	17.47	23.37	29.28	68.8	4.0	1.0	10.0
175	2020	49	16.25	21.21	26.17	51.1	5.0	1.0	5.0
176	2020	50	16.73	23.70	30.67	2.0	2.0	1.0	4.0
177	2020	51	18.29	23.14	27.98	13.3	3.0	1.0	11.0

**Figura 4:** Conjunto de dados (dataframe) final da localidade de Passo Fundo.

mes. Isso foi feito pois através de observações realizadas sobre a base de afídeos constatou-se a existência semanas em que nenhum dado é cadastrado. Após este passo, foram eliminados os registros cujo campo de afídeos era NaN, significando semanas em que a ocorrência de afídeos não foi cadastrada na base de dados de afídeos.

A fim de verificar a influência que os fenômenos climáticos ENOS têm sobre a dinâmica populacional dos afídeos e também com objetivo de melhorar a acurácia dos modelos treinados, foi adicionada uma coluna chamada fenomeno, que armazena o fenômeno climático que aconteceu em um determinado ano. Nesta nova coluna, o valor 0 (zero) corresponde à ocorrência do El niño, 1 para La niña e 2 para anos neutros, que significa que nenhum dos dois fenômenos aconteceram.

No intuito de extrair o máximo possível dos dados estudados e também, tendo como base o que disse (Cunha et al., 2001), onde argumenta que não só a quantidade mas também a ocorrência de chuva é maior em anos de El niño, foi adicionado em cada dataframe um campo chamado frqChuva, que armazena a quantidade de dias em que choveu na semana. Para a contabilização de um novo dia de chuva na semana foi levado em consideração qualquer tipo de valor maior que 0 (zero) para a coluna de precipitacao.

Portanto, como forma de exemplificação a Fig. 4 representa o dataframe final referente à localidade de Passo Fundo. Vale ressaltar que o mesmo processo foi realizado para as demais localidades estudadas. Também nesta etapa ocorreu a criação de um dataframe que concatena os dados de todas as localidades, que será utilizado posteriormente na matriz de correlação e treinamento de modelos de predição.

### 3.5 Data Mining - Processamento

Conforme citado na metodologia, nesta etapa é feita a extração de conhecimento e geração de inferência sobre os dados estudados.

Com o objetivo de extrair conhecimento, em um primeiro momento foi realizada a análise exploratória sobre os dados no intuito de verificar a correlação existente entre dados climáticos com a dinâmica populacional dos afídeos. Para isso, foram realizadas plotagens e análise de gráficos e também a utilização da matriz de correlação aplicando o coeficiente de correlação de Spearman.

Após isso, foram aplicadas técnicas de aprendizado de máquina sobre os dados no intuito de prever a população

de afídeos a partir de dados climáticos. As técnicas utilizadas neste processo foram Regressão Linear, Redes Neurais Artificiais, Árvore de decisão e random forest. A escolha das técnicas de predição foram feitas principalmente com base na leitura de artigos científicos.

Optou-se por implementar Regressão Linear pois Rodrigues (2013) investigaram o uso dessa técnica para realização de inferências e concluiu sua viabilidade para predição de dados. do Nascimento et al. (2018) citam ainda que o uso desta técnica pode ser feito quando predição exata dos valores não é necessária e sim, uma aproximação, o que se encaixa no objetivo do trabalho.

Foi detectada a não linearidade no comportamento populacional dos afídeos com base na análise exploratória que será exposta em seguida. Sendo assim, a técnica de Redes Neurais Artificiais foi testada com base no que foi exposto no referencial bibliográfico, onde é citado que o modelo de saída gerado a partir do treinamento de uma Rede Neural Artificial não é linear.

A Árvore de decisão foi implementada pois Dias (2002) classifica esta técnica como adequada para realizar predições de dados contendo variáveis numéricas. Neste sentido, através de pesquisas feitas na documentação da biblioteca Scikit-learn, julgou-se interessante a implementação da random forest pois a mesma faz uso de diversas Árvores de decisão, gerando como saída o dado com maior votação entre as árvores.

No intuito de melhorar a assertividade dos modelos desenvolvidos foi utilizado um recurso chamado GridSearchCV. Este recurso permite que diferentes configurações de treinamento sejam passadas como parâmetro, a partir disso são feitas todas combinações possíveis, retornando a que obteve melhor acurácia. O recurso GridSearchCV e as técnicas de predição supracitadas estão presentes na biblioteca Scikit-learn e suas aplicações foram desenvolvidas com base em estudos voltados sobre a documentação da biblioteca.

Para realizar o treinamento dos modelos de predição foram utilizados dados referente à todas localidades, sendo que 80% dos dados foram separados de forma aleatória para a etapa de treinamento e os 20% restantes foram utilizados para comparações e testar a acurácia dos modelos. As variáveis climáticas utilizadas para o treinamento de ambos modelos foram precipitacao, semana, temperaturaMinima, temperaturaMaxima, frqChuva e fenomeno.

#### 3.5.1 Análise Exploratória

O primeiro passo dado na etapa de análise exploratória foi avaliar o comportamento da população de afídeos durante o período estudado. Para isso foi feita uma análise mensal sobre os dados referente à localidade de Coxilha no período de jan/2019 à dez/2020 e Passo Fundo no período de jan/2017 à dez/2020. Estes períodos foram escolhidos por apresentarem maior frequência no armazenamento dos dados de afídeos.

Na Fig. 5 é exposto a dinâmica populacional de afídeos da localidade de Coxilha (A) e Passo Fundo (B). Nela, cada ano estudado é representado por uma linha no gráfico, sendo que, o eixo X representa os meses de cada ano e o eixo Y a média de afídeos encontrados em cada mês. A análise de ambos os gráficos permite perceber que a ocorrência de afídeos não ocorre de forma linear.



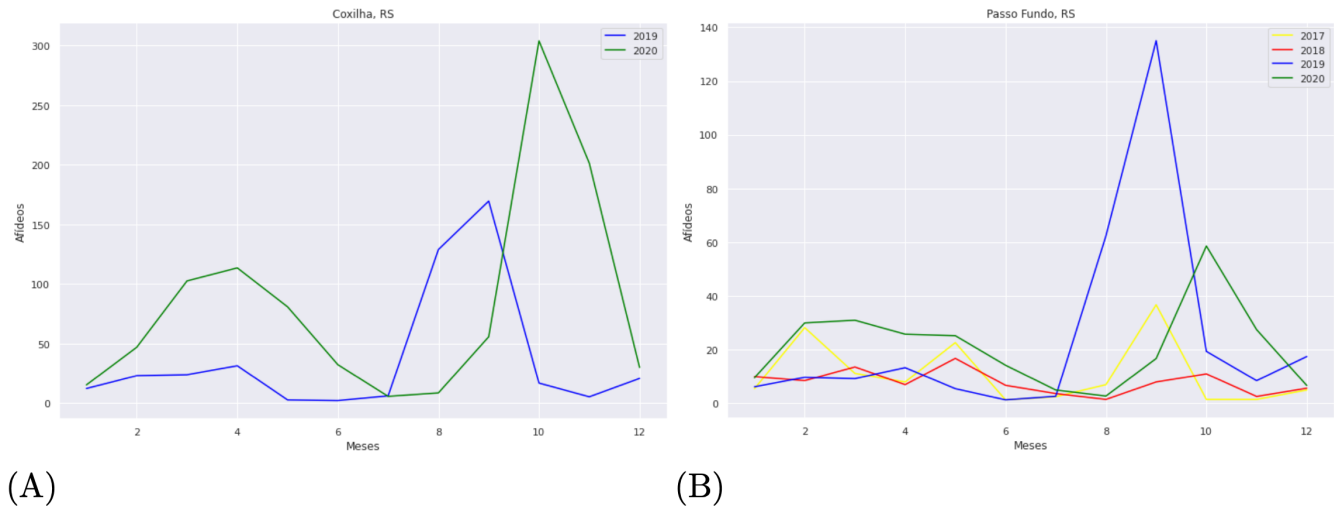


Figura 5: Comportamento mensal da população de afídeos.

Nos meses de junho, julho e agosto costuma ocorrer uma diminuição na quantidade populacional dos afídeos, este período corresponde a estação de inverno. Outro padrão que se encontrou foi que posterior a este período, nos meses de setembro e outubro, na maioria das ocasiões aconteceram picos na ocorrência destes pulgões, os meses fazem parte do final de inverno e se estende por grande parte da primavera.

Os padrões encontrados entre as épocas do ano e a população de afídeos estão de acordo com estudos feitos por Engel et al. (2021), que detectou que reduções significativas nos pulgões acontecem nos meses de junho, julho e agosto. Neste mesmo artigo é citado ainda os picos populacionais observados entre os meses de setembro e outubro. Tomé et al. (2013) corrobora expondo que os picos que acontecem nas estações de inverno e primavera representam de 89 a 95% das ocorrências de uma espécie de afídeos conhecida como *Rhopalosiphum padi* (Linnaeus, 1758) na localidade de Coxilha, RS.

Optou-se por verificar a correlação existente entre variáveis climáticas e a dinâmica populacional dos afídeos pois Oliveira (1971) conclui em seu estudo que fatores como temperatura e precipitações possuem correlação com a quantidade de afídeos. Oliveira (1971) argumenta ainda que anos de fenômenos *El Niño* e *La Niña* também possuem correlação com danos causados por pragas.

Sendo assim, foi realizada a análise de gráficos de forma semanal sobre cada uma das localidades em seus respectivos períodos expostos no estudo de caso. Esta análise pode ser observada na Fig. 6, onde cada linha de gráficos representa o estudo feito sobre determinada localidade. Nela, cada ponto presente nos gráficos representa o ponto de intersecção entre o eixo Y, que representa a ocorrência de afídeos na semana com o eixo X, que representa em diferentes gráficos as variáveis de temperaturaMedia, precipitacao, frqChuva e fenomeno.

Com base na análise dos gráficos da Fig. 6, notou-se que a população de afídeos diminui em semanas de temperaturas extremas e apresenta maior quantidade em semanas onde a temperatura média fica entre 15 e 22°C. O mesmo

padrão foi exposto em Pereira et al. (2016) e Tomé et al. (2013), onde a ocorrência de populações de afídeos acima da média foram detectadas em períodos em que a temperatura média se manteve entre 15 e 20°C. Esses mesmos estudos citam ainda picos na população de afídeos para semanas em que a precipitação se manteve abaixo de 20 mm para Pereira et al. (2016) e 30 mm para Tomé et al. (2013), conforme os gráficos expostos, os mesmos padrões podem ser observados.

Ainda se tratando da Fig. 6, através da análise do campo frqChuva este estudo supõe que a ocorrência de afídeos acima da média possui ambiente mais propício para acontecer quando chove entre 2 à 4 dias durante a semana. Quanto aos fenômenos ENOS, o gráfico expõe diminuição considerável na população dos pulgões em anos de *El Niño* (0) quando comparados a anos de *La Niña* (1) e anos neutros (2). Tal fato confirma a correlação exposta por Rosenzweig et al. (2001) entre a dinâmica populacional de afídeos com anos de fenômeno *El Niño* e *La Niña*.

Foi feito ainda o uso da matriz de correlação para verificar os níveis de correlação existentes entre as variáveis através de cálculos matemáticos. O cálculo de correlação de Spearman foi utilizado na matriz de correlação pois pode ser utilizado para avaliação de dados não lineares. Conforme a Fig. 6, as duas localidades apresentam o mesmo padrão na população de afídeos, sendo assim, optou-se por concatenar os dados das localidades estudadas para então aplicar a matriz de correlação.

Para fazer a análise da matriz presente na Fig. 7 basta associar a linha da variável desejada com a coluna de outra variável desejada, ou vice-versa. Sendo assim, o cálculo de Spearman mostra que a quantidade populacional de afídeos possui correlação com as variáveis temperaturaMínima (0,13), temperaturaMédia (0,15), temperaturaMáxima (0,17), precipitacao (0,1) e fenomeno (0,37).

O fato das correlações entre temperaturas e afídeos serem positivas mostra que grandes populações de afídeos tendem a ter dificuldades em ocorrer em temperaturas partindo de 0°C, conforme a temperatura aumenta esta probabilidade vai crescendo também. O mesmo ocorre com a

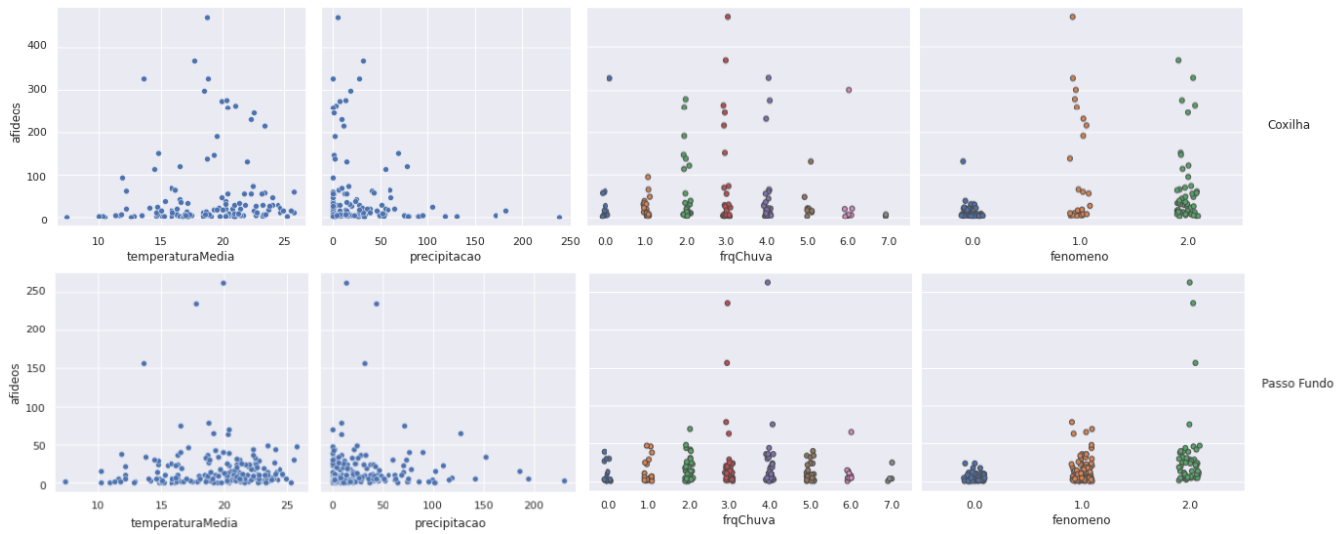


Figura 6: Gráficos relacionando dados climáticos e ocorrência de afídeos para cada localidade.

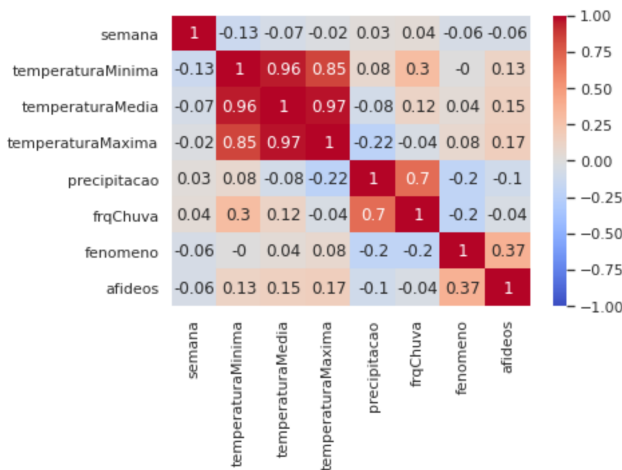


Figura 7: Matriz de correlação aplicada à dataframe contendo dados de todas localidades.

variável de precipitação, que por se tratar de uma correlação negativa, quer dizer que quanto mais milímetros de chuvas ocorrem por semana, menores são as populações de afídeos.

Segundo o cálculo de Spearman, a variável contendo fenômenos ENOS é a que apresenta maior correlação com a população destes insetos (0,37). Acredita-se que o motivo desta variável apresentar uma correlação consideravelmente maior em relação as outras é que a ocorrência destes fenômenos acaba alterando diversas variáveis climáticas não estudadas, como por exemplo umidade e os ventos, as quais alteram a quantidade populacional dos afídeos.

Vale ressaltar que o fato de que algumas variáveis não apresentam nenhum tipo de correlação na matriz não significa que as duas variáveis não tenham correlação entre si, pois, a não existência de correlação entre duas variáveis, não significa que não se verifique outro tipo de correlação

diferente da que foi aplicada, como por exemplo, a correlação linear ou exponencial.

### 3.5.2 Regressão Linear

O desenvolvimento da técnica de Regressão Linear se deu por meio da classe Linear- Regression presente na biblioteca do Scikit-learn. Conforme a documentação da biblioteca, a técnica de regressão linear não apresenta muitas configurações para o treinamento do modelo. Portanto o único passo executado para esta técnica foi o de treinamento sobre a base de dados que foi separada para etapa de treinamento dos modelos.

Para visualizar as previsões realizadas pelo modelo de regressão linear, a Fig. 8 mostra as ocorrências da base de dados de testes, cada ocorrência possui o valor real de afídeos (linha vermelha) e o valor predito pelo modelo (linha azul). O eixo X representa as ocorrências e o eixo Y a quantidade de afídeos. Utilizando o coeficiente de determinação  $R^2$  para verificar a acurácia do modelo de regressão linear sobre os dados separados para testes chegou-se a uma acurácia de 11,4%.

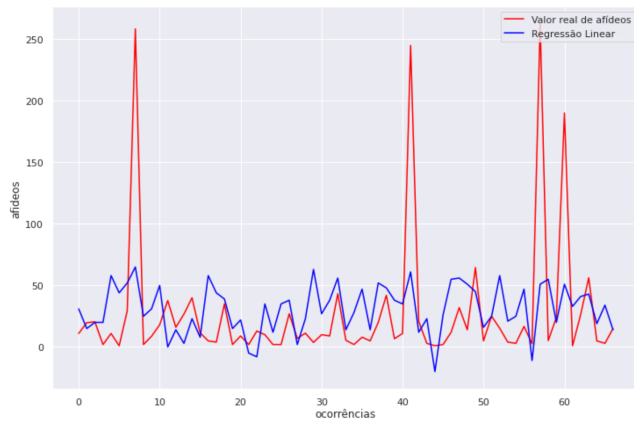
### 3.5.3 Redes Neurais Artificiais

Para implementar a técnica de Redes Neurais Artificiais foi utilizada a classe MLPRegressor da Scikit-learn. Para o treinamento do modelo foram feitas algumas configurações que tiveram como base os resultados do recurso GridSearchCV e testes realizados de forma manual.

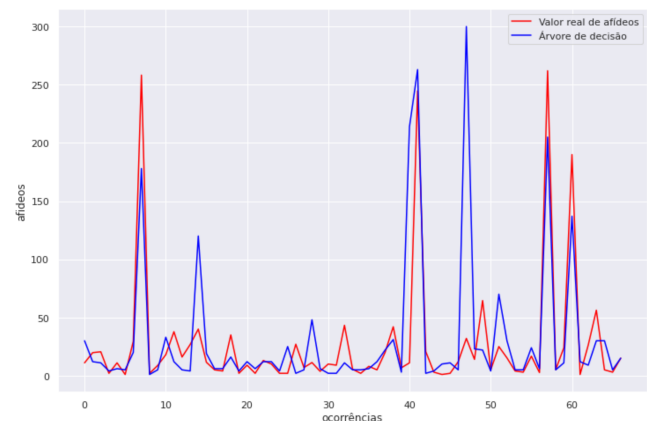
Ao fim do processo de testes para definir a melhor configuração, ficou definido que a rede deve possuir quatro camadas de 8 neurônios cada e o número máximo de iterações no processo de aprendizagem é de no máximo 5000 passos.

Como forma de visualizar as previsões desta técnica, a Fig. 9 mostra ocorrências dos dados de teste, cada ocorrência possui o valor real de afídeos (linha vermelha) e o valor predito pelo modelo (linha azul). O eixo X representa as ocorrências e o eixo Y a quantidade de afídeos. A acurácia obtida com coeficiente de determinação  $R^2$  foi de 26,4%.

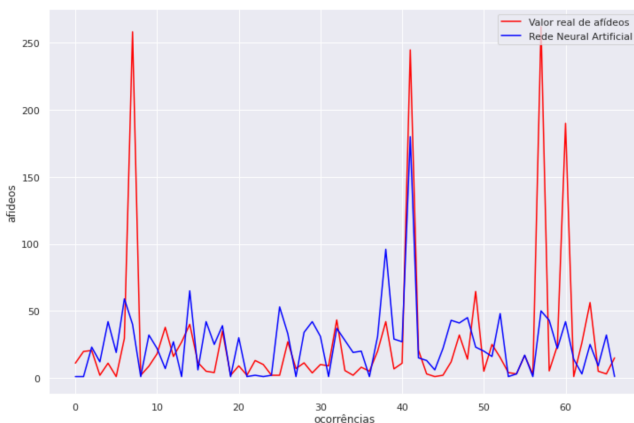




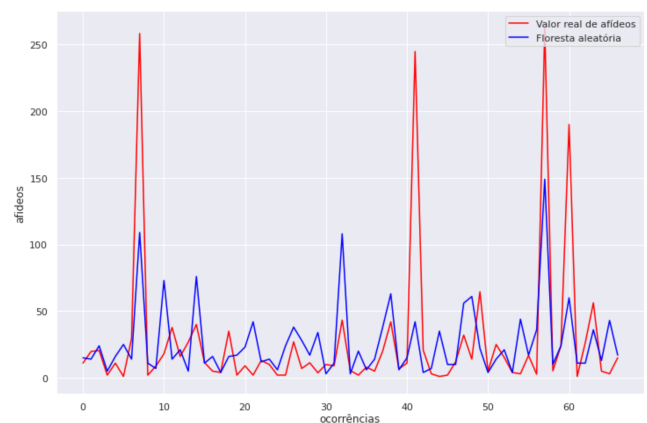
**Figura 8:** Ocorrências de afídeos em comparação à predição do modelo de Regressão Linear.



**Figura 10:** Ocorrências de afídeos em comparação à predição do modelo de Árvore de decisão.



**Figura 9:** Ocorrências de afídeos em comparação à predição do modelo de Redes Neurais Artificiais.



**Figura 11:** Ocorrências de afídeos em comparação à predição do modelo de random forest.

### 3.5.4 Árvore de decisão

A árvore de decisão foi implementada utilizando a classe `DecisionTreeRegressor`. Com base no que disse Lan et al. (2020), "o crescimento irrestrito da árvore de decisão pode facilmente levar ao sobreajuste e degradação de sua capacidade de generalização".

Sendo assim, para o processo de aprendizagem da árvore utilizou-se uma técnica conhecida como poda, que limita a profundidade máxima que a árvore terá, evitando assim a ocorrência do chamado *overfitting*. Após a realização de diversos testes manuais, definiu-se como oito a profundidade máxima para árvore de decisão.

A Fig. 10 mostra ocorrências da base de dados de testes, cada ocorrência possui o valor real de afídeos (linha vermelha) e o valor predito pelo modelo (linha azul). O eixo X representa as ocorrências e o eixo Y a quantidade de afídeos. Sobre a base de dados para teste este modelo chegou a acurácia de 29,3%.

### 3.5.5 Random forest

A classe responsável por implementar a random forest na biblioteca Scikit-learn é a `RandomForestRegressor`. Por

implementar o conceito de árvore, com auxílio do recurso `GridSearchCV` chegou-se a conclusão de que para o treinamento do modelo seria limitada a sete a profundidade máxima de cada uma das árvores.

Também definiu-se como dez o número máximo de árvores presentes na floresta, seis para o mínimo de amostras necessárias para dividir um nó interno e dois para o número mínimo de amostras para estar em um nó folha.

Na Fig. 11 é possível visualizar as ocorrências da base de dados de testes, cada ocorrência possui o valor real de afídeos (linha vermelha) e o valor predito pelo modelo (linha azul). O eixo X representa as ocorrências e o eixo Y a quantidade de afídeos. A acurácia da random forest sobre os dados separados para testes alcançou o valor de 41,4%.

## 3.6 Discussão dos Modelos Implementados

Para entender de que forma os dados foram utilizados pelos modelos para realizar as predições, a presente Eq. (1) refere-se ao modelo de regressão linear desenvolvido. Nela, o coeficiente linear da reta é representado pela constante  $-65,79$  enquanto os valores do coeficiente angular cor-

respondem respectivamente aos campos de: precipitação (-0,18), semana (0,26), temperaturaMaxima (5,89), temperaturaMinima (-6,88), frqChuva (8,15) e fenomeno (18,50).

$$Y = -65,79 - 0,18X_1 + 0,26X_2 + 5,89X_3 - 6,88X_4 + 8,15X_5 + 18,50X_6 \quad (1)$$

Através da análise da equação, acredita-se que o valor coeficiente linear da reta explica a predição de valores negativos para a população de afídeos, que podem ser observados na Fig. 8. Também é possível analisar a intensidade e sentido que o modelo previu para cada uma das variáveis utilizadas.

Para as técnicas de árvore de decisão e random forest existe um recurso que retorna a importância das variáveis utilizadas durante o treinamento, sendo que este é exibido em formato de porcentagem. Porém, para a técnica de rede neural artificial, nenhum recurso similar foi encontrado na documentação da biblioteca Scikit-learn.

Sendo assim, a árvore de decisão classificou como importância das variáveis: 22% para precipitação; 17% para semana; 06% para temperaturaMaxima; 34% para temperaturaMinima; 09% para frqChuva; 12% para fenomeno. Enquanto a técnica de random forest definiu como porcentagem de importância durante o treinamento: 10% para precipitação; 35% para semana; 14% para temperaturaMaxima; 20% para temperaturaMinima; 08% para frqChuva; 13% para fenomeno.

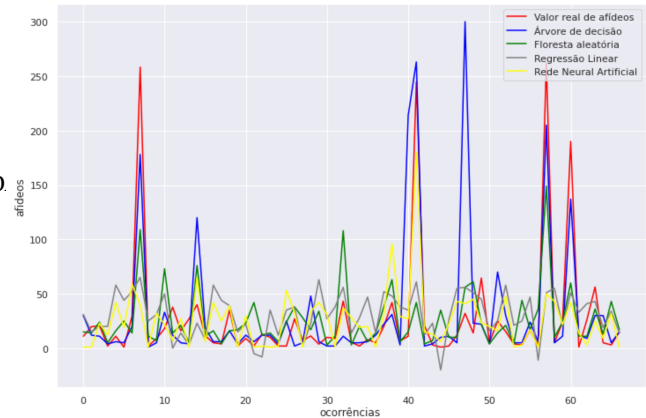
Como forma de comparar os resultados dos modelos desenvolvidos, a Tabela 3 mostra uma tabela contendo as dez primeiras ocorrências dos dados separados para teste (Valor real de afídeos) e os valores obtidos pelos modelos de predição desenvolvidos.

**Tabela 3:** Populações de afídeos em comparação a inferências realizadas pelos modelos.

Nro Afídeos	AD	RNA	RF	RL
11	30	1	15	31
19	12	1	14	15
20	11	23	24	20
2	4	12	5	20
11	6	42	16	58
1	5	19	25	44
29	20	59	14	52
258	178	40	109	65
2	1	1	11	25
8	5	32	7	31

AD: Árvore de Decisão;  
RNA: Rede Neural Artificial;  
RF: Random Forest;  
RL: Regressão Linear.

Se tem como destaque as técnicas de Árvore de decisão e random forest que, salvo excessões, foram capazes de acompanhar a dinâmica da população de afídeos, inclusive, na ocorrência da população de afídeos acima da média que aconteceu no registro de índice 7. Por outro lado, as técnicas de Redes neurais artificiais e Regressão linear obtiveram resultados inferiores quando comparadas



**Figura 12:** Ocorrências de afídeos em comparação à predição do modelo.

às outras duas técnicas estudadas.

A Fig. 12 apresenta a base de dados de testes completa em comparação aos valores preditos pelos modelos. Nela é possível analisar o mesmo padrão para random forest e Árvore de decisão, sendo que este último se mantém mais próximo da real ocorrência de afídeos quando comparada à random forest. Outro fator observado é a dificuldade apresentada pelas técnicas de Regressão linear e Redes Neurais Artificiais em prever a ocorrência da população de afídeos acima da média.

A técnica que apresentou maior acurácia no coeficiente de determinação  $R^2$  foi random forest (41,4%). Porém tanto random forest quanto Árvore de decisão apresentaram bons resultados de predições.

Para analisar individualmente o comportamento dos modelos durante o ano todo, na Fig. 13 são plotados gráficos para cada localidade e ano estudado. Nestes gráficos, a linha vermelha representa a população de afídeos que de fato aconteceu e as demais linhas os valores preditos pelos modelos. O eixo X representa as semanas do ano e o eixo Y a quantidade de afídeos. A localidade e ano analisado é exposto acima de cada um dos gráficos.

Além da aproximação da população real com a predita, foi levado em consideração a capalocalidade que as técnicas estudadas tem de se adaptar aos períodos de baixa na população de afídeos, junho e agosto (semanas 23 à 32 de cada ano) e os períodos de população de afídeos acima da média, entre os meses de setembro à outubro (semanas do ano 36 à 42).

Conforme calculado pelo coeficiente de determinação  $R^2$ , o modelo de regressão linear é o que apresenta o pior desempenho entre as técnicas abordadas. Em períodos normais para população de afídeo seu desempenho é razoável, porém o modelo não é capaz de prever diminuição ou aumento na população dos insetos. Também em diversos momentos o modelo prevê populações negativas de afídeos, o que inviabiliza sua aplicação.

Com exceções, a rede neural artificial foi capaz de prever alguns picos na população de afídeos, porém se mostrou inconstante em períodos onde a dinâmica da população destes insetos se mantém dentro da média. Apesar de acontecer com menos frequência quando comparado ao



Figura 13: Ocorrências semanais de afídeos em comparação à predição do modelo.

modelo de regressão linear, a rede neural artificial também prevê alguns casos onde a população de afídeo é negativa.

Os modelos de árvore de decisão e random forest são os que melhor representam os valores reais da população de afídeos. Ambos são capazes de identificar a redução populacional que estes insetos sofrem durante os meses de junho à agosto, bem como seu aumento entre os meses de setembro à outubro.

O modelo de Árvore de decisão aparenta ter maior aproximação com as ocorrências reais de afídeos quando comparado ao modelo de random forest, porém ambos são eficazes para prever os dados estudados.

Acredita-se ainda ser possível melhorar os resultados dos modelos com a adição de novas variáveis climáticas para o treinamento, pois, conforme citado por Pereira et al. (1999), variáveis como umidade relativa do ar, insolação, vento, balanço hídrico, radiação solar global também possuem correlação com a dinâmica populacional destes pulgões, porém não estão sendo abordadas neste estudo.

#### 4 Considerações Finais

O objetivo de organizar uma base de dados referente à regiões onde existem monitoramento de afídeos foi alcançado. O que possibilitou o estudo da correlação existente entre população de afídeos e dados climáticos, a fim de extrair conhecimentos sobre a dinâmica populacional dos afídeos com base em dados e fenômenos climáticos.

Através da matriz de correlação e também a visualização dos gráficos foi possível verificar a forma de como a alteração nas variáveis climáticas de temperatura, precipitação e fenômenos climáticos como *El Niño* e *La Niña* impactam na quantidade populacional destes pulgões. Conforme visto em alguns estudos citados nos resultados, também foi possível verificar períodos do ano em que as populações de afídeos aumentam (setembro e outubro) e diminuem (junho, julho e agosto).

Ainda através da análise de dados realizada no trabalho, este estudo propõe que a quantidade de dias de chuva que ocorrem em uma determinada semana também é um fator que possui correlação com a diminuição ou aumento na população de afídeos.

Quanto aos modelos desenvolvidos conclui-se que as técnicas que utilizam Árvore de decisão foram capazes de prever os dados estudados com maior eficiência quando comparados às técnicas de Regressão Linear e Redes Neurais Artificiais. Este fator pode ser verificado também através do coeficiente de determinação  $R^2$  aplicado sobre os dados separados para teste, onde Regressão Linear apresentou uma acurácia de 11,4%; Redes Neurais Artificiais a acurácia de 26,4%; Árvore de decisão acurácia de 29,3% e random forest acurácia de 41,4%.

Como trabalhos futuros, acredita-se que a acurácia dos modelos pode ser melhorada com a utilização de um conjunto maior de dados para o processo de treinamento, podendo estes dados serem referentes a mais localidades em que se tem o monitoramento de afídeos. Outra forma de melhorar os resultados de predição deve ser através da adição de novas variáveis que influenciam na população de afídeos como por exemplo: variáveis climáticas como umidade relativa do ar, insolação, vento, entre outras; inimigos

naturais dos afídeos como parasitóides e predadores.

Outro fator que deve ser levado em consideração para trabalhos futuros é a realização da discretização dos valores, que significa transformar os dados contínuos que foram estudados para dados categóricos, que abrangem determinados períodos de dados, podendo assim melhorar a assertividade dos modelos, pois os dados são classificados dentro de um conjunto finito de valores.

Por fim, com base nas considerações citadas, pode-se concluir que o objetivo de extrair conhecimentos sobre a dinâmica populacional dos afídeos em relação às variáveis climáticas foi alcançado, pois com o trabalho foi possível detectar situações em que a população de afídeos aumenta ou diminui conforme a variação dos dados climáticos estudados.

Para as técnicas de predição utilizadas, levando-se em consideração a análise realizada nos gráficos e os resultados obtidos no coeficiente de determinação  $R^2$ , técnicas que implementam Árvores de decisão como a própria Árvore de decisão e a random forest são mais apropriados para o tipo de dados estudado neste trabalho.

#### Referências

- Auad, A. M., Freita, S. D. and Barbosa, L. R. (2002). Ocorrência de afídeos em alface (*lactuca sativa* l.) em cultivo hidropônico, *Neotropical Entomology* 31(2): 335–339. <https://doi.org/10.1590/S1519-566X2002000200025>.
- CABI (2022). *Rhopalosiphum padi* (grain aphid), *CABI Compendium* **CABI Compendium**. <https://doi.org/10.1079/cabicompendium.47321>.
- Cunha, G. R., Dalmago, G. A., Estefanel, V., AldemirPasinato and Moreira, M. B. (2001). El niño – oscilação do sul e seus impactos sobre a cultura de cevada no Brasil, **09(1)**: 137–145. Disponível em <http://www.sbagro.org/files/biblioteca/1264.pdf>.
- Dhar, V. (2012). Data science and prediction, *Communications of the ACM* 56. <https://dx.doi.org/10.2139/ssrn.2086734>.
- Dias, M. (2002). Parâmetros na escolha de técnicas e ferramentas de mineração de dados, *Acta Scientiarum: Technology* 24. <https://doi.org/10.4025/actascitechnol.v24i0.2549>.
- do Nascimento, R. L. S., da Cruz Junior, G. G. and de Araújo Fagundes, R. A. (2018). Mineração de dados educacionais: Um estudo sobre indicadores da educação em bases de dados do inep, *Revista Novas Tecnologias na Educação* 16(1). <https://doi.org/10.22456/1679-1916.85989>.
- Engel, E., Lau, D., Godoy, W. A. C., Pasini, M. P. B., Malaquias, J. B., Santos, C. D. R., Pivato, J. and Pereira, P. R. V. d. S. (2021). Oscillation, synchrony, and multi-factor patterns between cereal aphids and parasitoid populations in southern Brazil, *Bulletin of Entomological Research*. <https://doi.org/10.1017/s0007485321000729>.
- Fernandes, G. A. G., Araújo, H. L. and Gomes, R. C. C. (2013). Sistema automatizado de aquisição de dados meteorológicos., *Jornada Acadêmica da UEG*.



- Finkler, C. L. L. (2013). Controle de insetos: Uma breve revisão, *Anais da Academia Pernambucana de Ciência Agrônômica* 8: 169–189. Disponível em <https://www.journals.ufrpe.br/index.php/apca/article/view/155>.
- Gassen, D. N. (1984). Insetos associados a cultura do trigo no Brasil, *Embrapa trigo-circular técnica (infoteca-e)*, EMBRAPA-CNPT. Disponível em <http://www.infoteca.cnptia.embrapa.br/infoteca/handle/doc/846603>.
- INPE (2021). Instituto nacional de pesquisas espaciais. Disponível em <https://www.gov.br/inpe/pt-br>.
- Kamilaris, A., Kartakoullis, A. and Prenafeta-Boldú, F. X. (2017). A review on the practice of big data analysis in agriculture, *Computers and Electronics in Agriculture* 143: 23–37. <https://doi.org/10.1016/j.compag.2017.09.037>.
- Lan, T., Hu, H., Jiang, C., Yang, G. and Zhao, Z. (2020). A comparative study of decision tree, random forest, and convolutional neural network for spread-f identification, *Advances in Space Research* 65(8): 2052–2061. <https://doi.org/10.1016/j.asr.2020.01.036>.
- Lazzaretti, A. T., Lau, D., Fernandes, J. M. C., Wiest, R., Bavaresco, J. L. B. and Schaefer, F. (2016). Trapsystem – uma aplicação para gerenciamento de dados coletados a partir de armadilhas de insetos, *In: Reunião da Comissão Brasileira de Pesquisa de Trigo e Triticale*. Disponível em <https://www.embrapa.br/busca-de-publicacoes/-/publicacao/1064573/trapsystem---uma-aplicacao-para-gerenciamento-de-dados-coletados-a-partir-de-armadilhas-de-insetos>.
- Lidiane, L. M. K., Kuang, H., Patrícia, B. C. and Milton, P. J. (2018). Análise fatorial por meio da matriz de correlação de Pearson e polícrica no campo das cisternas, *Engineering and Science* 7(1): 58–70. <https://doi.org/10.18607/ES201875266>.
- Navathe, E. . (ed.) (2018). *Sistemas de bancos de dados*, Pearson.
- NOAA (2022). National oceanic and atmospheric administration. Disponível em <https://www.noaa.gov>.
- Oliveira, A. M. (1971). Observações sobre a influência de fatores climáticos nas populações de afídeos em batata, *Pesquisa Agropecuária Brasileira* 6: 163–172. Disponível em <https://seer.sct.embrapa.br/index.php/pab/article/view/17666>.
- Pandas Development Group (2023). *Pandas – Python*, Pandas Community. Disponível em <https://pandas.pydata.org>.
- Pereira, A. B., Banzatto, D. A. and Furiatti, R. (1999). Influência de elementos climáticos na flutuação populacional de alados de myzus persicae (sulzer)(homoptera: Aphididae), *Revista Brasileira de Agrometeorologia*, 7(2): 219–225. Disponível em <http://www.sbagro.org/files/biblioteca/252.pdf>.
- Pereira, P. d. S., Lau, D. and Júnior, A. M. (2016). Dinâmica populacional de afídeos vetores de bydy: impactos ao rendimento de grãos em trigo, *In: Reunião da Comissão Brasileira de Pesquisa de Trigo e Triticale*. Disponível em <https://www.embrapa.br/busca-de-publicacoes/-/publicacao/1068229/dinamica-populacional-de-afideos-vetores-de-bydy-impactos-ao-rendimento-de-graos-em-trigo>.
- Python.org (2023). *What is Python? Executive Summary*, Python.org. Available at <https://www.python.org/doc/essays/blurb/>.
- Rodrigues, R. L.; Medeiros, F. P. D. G. A. S. (2013). Modelo de regressão linear aplicado à previsão de desempenho de estudantes em ambiente de aprendizagem, *In: Brazilian Symposium on Computers in Education (Simpósio Brasileiro de Informática na Educação-SBIE)*, Vol. 24. <https://doi.org/10.5753/cbie.sbie.2013.607>.
- Rosenzweig, C., Iglesias, A., Yang, X., Epstein, P. and Chivian, E. (2001). Climate change and extreme weather events; implications for food production, plant diseases, and pests, *Global Change and Human Health* 2: 90–104. <https://doi.org/10.1023/A:1015086831467>.
- Scikit Learning Development Group (2023). *Scikit-learn Machine Learning in Python*, Scikit-Learning Community. Disponível em <https://scikit-learn.org/stable/#>.
- Stern, D. L. (2008). Aphids, *Current Biology* 18(12): R504–R505. <https://doi.org/10.1016/j.cub.2008.03.034>.
- Teixeira Martins, M., Marangon, G., Costa, E., Silveira, B., Cubas, R. and Cavalli, J. (2019). Estimação da altura de plantios florestais de eucalipto por regressão e redes neurais artificiais, *BIOFIX Scientific Journal* 5: 141–152. <http://dx.doi.org/10.5380/biofix.v5i1.68839>.
- Toebe, J. (2014). *Um modelo baseado em agentes para o ciclo de vida dos insetos: aplicação na interação afídeo-planta-vírus*, Phd in Agronomy, Post-Graduate Program in Agronomy at University of Passo Fundo. Disponível em <https://secure.upf.br/pdf/2014JosueToebe.pdf>.
- Tomé, A. C., Lau, D., Pereira, P. R. V. d. S. and Marsaro Júnior, A. L. (2013). Dinâmica das populações de afídeos de cereais de inverno em coxilha/rs entre 2011 e setembro de 2013, *In: Mostra de Iniciação Científica da Embrapa Trigo*. Disponível em <https://www.embrapa.br/busca-de-publicacoes/-/publicacao/1005396/dinamica-das-populacoes-de-afideos-de-cereais-de-inverno-em-coxilha-entre-2011-e-setembro-de-2013>.
- Torres, J. D., Monteiro, I. O., Santos, J. R. and Ortiz, M. S. (2015). Aquisição de dados meteorológicos através da plataforma arduino: construção de baixo custo e análise de dados, *Scientia Plena* 11(2). Disponível em <https://www.scientiaplena.org.br/sp/article/view/1742>.
- Wiest, R., Salvadori, R., J., Fernandes, J. M., Lau, D., Pavan, W., Zanini, W. R. and Lazzaretti, A. T. (2021). Population growth of rhopalosiphum padi under different thermal regimes: an agent-based model approach, *Agricultural and Forest Entomology* 23(1): 59–69. <https://doi.org/10.1111/afe.12404>.