

UNIVERSIDADE FEDERAL DE UBERLÂNDIA

Juliana de Faria Paulo

**Análise de incidentes de segurança usando séries  
temporais e modelos ARIMA**

**Uberlândia, Brasil**

**2021**

UNIVERSIDADE FEDERAL DE UBERLÂNDIA

Juliana de Faria Paulo

**Análise de incidentes de segurança usando séries temporais e  
modelos ARIMA**

Monografia apresentada ao Curso de Bacharelado em Ciência da Computação da Universidade Federal de Uberlândia, como requisito parcial para obtenção do título de Bacharel em Ciência da Computação.

Orientador: Rodrigo Sanches Miani

Universidade Federal de Uberlândia – UFU

Faculdade de Computação

Curso de Bacharelado em Ciência da Computação

Uberlândia, Brasil

2021

Juliana de Faria Paulo

## **Análise de incidentes de segurança usando séries temporais e modelos ARIMA**

Monografia apresentada ao Curso de Bacharelado em Ciência da Computação da Universidade Federal de Uberlândia, como requisito parcial para obtenção do título de Bacharel em Ciência da Computação.

Trabalho aprovado. Uberlândia, Brasil, 18 de junho de 2021:

---

**Rodrigo Sanches Miani**  
Orientador

---

**Paulo Henrique Ribeiro Gabriel**  
Professor

---

**Julio Fernando Costa Santos**  
Professor

Uberlândia, Brasil  
2021

*Aos meus pais David Paulo e Maria do Carmo do Nascimento Faria Paulo  
e irmã Carolina de Faria Paulo*

# Agradecimentos

Agradeço primeiramente aos meus pais por sempre apoiarem e incentivarem meus estudos, principalmente depois que mudei de cidade para iniciar na faculdade. O apoio e conselho de vocês ajudaram a moldar tanto minha vida acadêmica quanto profissional.

À minha irmã pelo carinho, incentivo e apoio, inclusive nas madrugadas.

Aos meus amigos pelos momentos de descontração e compreensão durante esse momento.

A todos os professores que contribuíram para a minha formação acadêmica, desde a infância até a graduação.

Ao Professor Doutor Julio Fernando Costa Santos pelo aprendizado, disposição e paciência em tirar qualquer dúvida sobre a parte estatística do trabalho.

Ao meu orientador, Professor Doutor Rodrigo Miani, pelo apoio, confiança e cuidado nessa jornada. Só tenho a agradecer pelo seu jeito atencioso e descontraído que me guiou nesse trabalho.

# Resumo

Esse trabalho utiliza conceitos de análise de dados uma base pública sobre incidentes de segurança mantida de maneira independente, a *Hackmageddon*, a fim de se tentar prever comportamentos futuros dos mesmos. Ao longo da pesquisa, os dados foram transformados em cinco séries temporais, uma para todos os incidentes e quatro para cada tipo de incidente (Crime Cibernético, Espionagem Cibernética, Guerra Cibernética e Hacktivismo). A partir delas, foram utilizadas funções de autocorrelação, autocorrelação parcial, análise de gráficos e testes de validação para construir modelos ARIMA e, em seguida, elaborar previsões a partir deles. Os resultados encontrados indicam um desempenho semelhante entre os modelos que tiveram períodos de estimativa maiores se saíram melhores do que os que tiveram esse reduzido.

**Palavras-chave:** Segurança da informação, Modelos ARIMA, Incidentes de segurança, Previsão de incidentes.

# Abstract

This work uses data analysis concepts from an independently maintained public database on security incidents, the *Hackmageddon*, in order to try to predict future behavior of the same. During the research, the data were transformed into five time series, one for all incidents and four for each type of incident (Cyber Crime, Cyber Espionage, Cyber War and Hacktivism). From them, autocorrelation functions, partial autocorrelation functions, graph analysis and validation tests were used to build ARIMA models and elaborate predictions from them. The results indicated a similar performance among the models that had longer estimation periods performed better than those that used smaller periods.

**Keywords:** Information Security, ARIMA Models, Security incidents, Incident Forecast.

# Lista de ilustrações

Figura 1 – Exemplo de Série Temporal. Fonte: Do autor. . . . .	18
Figura 2 – Exemplo de tabela da base dados <i>Hackmageddon</i> . Fonte: Extraída de (PAS-SERI, 2018). . . . .	25
Figura 3 – Etapas de Definição de um Modelo ARIMA Fonte: Do autor. . . . .	26
Figura 4 – Exemplo do retorno da função <i>ggtstdisplay</i> do modelo. Fonte: Extraída de (HYNDMAN; ATHANASOPOULOS, 2018). . . . .	27
Figura 5 – Exemplo do retorno da função <i>autoplot</i> do modelo. Fonte: Extraída de (HYNDMAN; ATHANASOPOULOS, 2018). . . . .	28
Figura 6 – Exemplo do retorno da função <i>checkresiduals</i> do modelo. Fonte: Extraída de (HYNDMAN; ATHANASOPOULOS, 2018). . . . .	29
Figura 7 – Série temporal de Crimes Cibernéticos e seus ACF e PACF. Fonte: Do autor.	32
Figura 8 – Teste de Círculo Unitário para Série de Crimes Cibernéticos: ARIMA(1,0,1). Fonte: Do autor. . . . .	33
Figura 9 – Teste do Resíduo para o modelo estimado para Série: ARIMA(1,0,1). Fonte: Do autor. . . . .	34
Figura 10 – Crimes Cibernéticos - Valores Ajustados vs Dados Observados: ARIMA(1,0,1). Fonte: Do autor. . . . .	35
Figura 11 – Série temporal de Espionagem Cibernética e seus ACF e PACF. Fonte: Do autor. . . . .	35
Figura 12 – 1ª Diferença da Série temporal de Espionagem Cibernética e seus ACF e PACF. Fonte: Do autor. . . . .	36
Figura 13 – Teste de Círculo Unitário para o modelo da Série de Espionagem Cibernética: ARIMA(1,1,1). Fonte: Do autor. . . . .	37
Figura 14 – Teste do Resíduo para o modelo da Série de Espionagem Cibernética: ARIMA(1,1,1). Fonte: Do autor. . . . .	37
Figura 15 – Espionagem Cibernética - Valores Ajustados vs Dados Observados: ARIMA(1,1,1). Fonte: Do autor. . . . .	38
Figura 16 – Série temporal de Guerra Cibernética e seus ACF e PACF. Fonte: Do autor. .	38
Figura 17 – Teste de Círculo Unitário para o modelo da série de Guerra Cibernética: ARIMA(2,0,3). Fonte: Do autor. . . . .	39
Figura 18 – Teste do Resíduo para o modelo da Série de Guerra Cibernética: ARIMA(2,0,3). Fonte: Do autor. . . . .	40
Figura 19 – Guerra Cibernética - Valores Ajustados vs Dados Observados: ARIMA(2,0,3). Fonte: Do autor. . . . .	41
Figura 20 – Série temporal de Hacktivismo e seus ACF e PACF. Fonte: Do autor. . . . .	41

Figura 21 – 1ª Diferença da Série temporal de Hacktivismo e seus ACF e PACF. Fonte: Do autor. . . . .	42
Figura 22 – Teste de Círculo Unitário para o modelo da Série de Hacktivismo: ARIMA(1,1,1). Fonte: Do autor. . . . .	43
Figura 23 – Teste do Resíduo para o modelo da Série de Hacktivismo: ARIMA(1,1,1). Fonte: Do autor. . . . .	43
Figura 24 – Hacktivismo - Valores Ajustados vs Dados Observados: ARIMA(1,1,1). Fonte: Do autor. . . . .	44
Figura 25 – Série temporal do Total de incidentes e seus ACF e PACF. Fonte: Do autor. .	44
Figura 26 – Teste de Círculo Unitário para Série de Total de incidentes: ARIMA(1,0,1). Fonte: Do autor. . . . .	45
Figura 27 – Teste do Resíduo para o modelo da Série de Total de incidentes: ARIMA(1,0,1). Fonte: Do autor. . . . .	46
Figura 28 – Total de Incidentes - Valores Ajustados vs Dados Observados: ARIMA(1,1,1). Fonte: Do autor. . . . .	47
Figura 29 – Série temporal de Crimes Cibernéticos e seus ACF e PACF. Fonte: Do autor.	47
Figura 30 – 1ª Diferença da Série temporal de Crimes Cibernéticos e seus ACF e PACF. Fonte: Do autor. . . . .	48
Figura 31 – Teste de Círculo Unitário para Série de Crimes Cibernéticos: ARIMA(1,1,1). Fonte: Do autor. . . . .	49
Figura 32 – Teste do Resíduo para o modelo estimado para Série: ARIMA(1,1,1). Fonte: Do autor. . . . .	49
Figura 33 – Crimes Cibernéticos - Valores Ajustados vs Dados Observados: ARIMA(1,1,1). Fonte: Do autor. . . . .	50
Figura 34 – Série temporal de Espionagem Cibernética e seus ACF e PACF. Fonte: Do autor. . . . .	50
Figura 35 – 1ª Diferença da Série temporal de Espionagem Cibernética e seus ACF e PACF. Fonte: Do autor. . . . .	51
Figura 36 – Teste de Círculo Unitário para o modelo da Série de Espionagem Cibernética: ARIMA(1,1,1). Fonte: Do autor. . . . .	52
Figura 37 – Teste do Resíduo para o modelo da Série de Espionagem Cibernética: ARIMA(1,1,1). Fonte: Do autor. . . . .	52
Figura 38 – Espionagem Cibernética - Valores Ajustados vs Dados Observados: ARIMA(1,1,1). Fonte: Do autor. . . . .	53
Figura 39 – Série temporal de Guerra Cibernética e seus ACF e PACF. Fonte: Do autor. .	53
Figura 40 – Teste de Círculo Unitário para o modelo da série de Guerra Cibernética: ARIMA(2,0,3). Fonte: Do autor. . . . .	54
Figura 41 – Teste do Resíduo para o modelo da Série de Guerra Cibernética: ARIMA(2,0,3). Fonte: Do autor. . . . .	55

Figura 42 – Guerra Cibernética - Valores Ajustados vs Dados Observados: ARIMA(2,0,3). Fonte: Do autor. . . . .	56
Figura 43 – Série temporal de Hacktivismos e seus ACF e PACF. Fonte: Do autor. . . . .	56
Figura 44 – 1ª Diferença da Série temporal de Hacktivismos e seus ACF e PACF. Fonte: Do autor. . . . .	57
Figura 45 – Teste de Círculo Unitário para o modelo da Série de Hacktivismos: ARIMA(1,1,1). Fonte: Do autor. . . . .	58
Figura 46 – Teste do Resíduo para o modelo da Série de Hacktivismos: ARIMA(1,1,1). Fonte: Do autor. . . . .	58
Figura 47 – Hacktivismos - Valores Ajustados vs Dados Observados: ARIMA(1,1,1). Fonte: Do autor. . . . .	59
Figura 48 – Série temporal do Total de incidentes e seus ACF e PACF. Fonte: Do autor. . . . .	59
Figura 49 – 1ª Diferença da Série temporal de Total de incidentes e seus ACF e PACF. Fonte: Do autor. . . . .	60
Figura 50 – Teste de Círculo Unitário para Série de Total de incidentes: ARIMA(1,1,2). Fonte: Do autor. . . . .	61
Figura 51 – Teste do Resíduo para o modelo da Série de Total de incidentes: ARIMA(1,1,2). Fonte: Do autor. . . . .	61
Figura 52 – Total de Incidentes - Valores Ajustados vs Dados Observados: ARIMA(1,1,2). Fonte: Do autor. . . . .	62
Figura 53 – Crimes Cibernéticos - Previsão de 2018 a 2020 usando todo o período de estimativa de 2011 a 2017 . . . . .	63
Figura 54 – Espionagem Cibernética - Previsão de 2018 a 2020 usando todo o período de estimativa de 2011 a 2017. Fonte: Do autor. . . . .	63
Figura 55 – Guerra Cibernética - Previsão de 2018 a 2020 usando todo o período de estimativa de 2011 a 2017. Fonte: Do autor. . . . .	64
Figura 56 – Hacktivismos - Previsão de 2018 a 2020 usando todo o período de estimativa de 2011 a 2017. Fonte: Do autor. . . . .	64
Figura 57 – Total de incidentes - Previsão de 2018 a 2020 usando todo o período de estimativa de 2011 a 2017. Fonte: Do autor. . . . .	65
Figura 58 – Crimes Cibernéticos - Previsão de 2018 a 2020 usando todo o período de estimativa de 2014 a 2017. Fonte: Do autor. . . . .	66
Figura 59 – Espionagem Cibernética - Previsão de 2018 a 2020 usando todo o período de estimativa de 2014 a 2017. Fonte: Do autor. . . . .	66
Figura 60 – Guerra Cibernética - Previsão de 2018 a 2020 usando todo o período de estimativa de 2014 a 2017. Fonte: Do autor. . . . .	67
Figura 61 – Hacktivismos - Previsão de 2018 a 2020 usando todo o período de estimativa de 2014 a 2017. Fonte: Do autor. . . . .	67

Figura 62 – Total de incidentes - Previsão de 2018 a 2020 usando todo o período de  
estimativa de 2014 a 2017. Fonte: Do autor. . . . . 68

# Lista de tabelas

Tabela 1 – Resultados Teste ADF para Crimes Cibernéticos . . . . .	33
Tabela 2 – Resultados Teste Ljung-Box e Jarque-Bera para Crimes Cibernéticos . . . . .	34
Tabela 3 – Resultados Teste ADF para Espionagem Cibernética . . . . .	36
Tabela 4 – Resultados Teste Ljung-Box e Jarque-Bera para Espionagem Cibernética . . . . .	37
Tabela 5 – Resultados Teste ADF para Guerra Cibernética . . . . .	39
Tabela 6 – Resultados Teste Ljung-Box e Jarque-Bera para Guerra Cibernética . . . . .	40
Tabela 7 – Resultados Teste ADF para Hacktivismo . . . . .	42
Tabela 8 – Resultados Teste Ljung-Box e Jarque-Bera para Hacktivismo . . . . .	43
Tabela 9 – Resultados Teste ADF para Total de incidentes . . . . .	45
Tabela 10 – Resultados Teste Ljung-Box e Jarque-Bera para Total de incidentes . . . . .	46
Tabela 11 – Resultados Teste ADF para Crimes Cibernéticos . . . . .	48
Tabela 12 – Resultados Teste Ljung-Box e Jarque-Bera para Crimes Cibernéticos . . . . .	49
Tabela 13 – Resultados Teste ADF para Espionagem Cibernética . . . . .	51
Tabela 14 – Resultados Teste Ljung-Box e Jarque-Bera para Espionagem Cibernética . . . . .	52
Tabela 15 – Resultados Teste ADF para Guerra Cibernética . . . . .	54
Tabela 16 – Resultados Teste Ljung-Box e Jarque-Bera para Guerra Cibernética . . . . .	55
Tabela 17 – Resultados Teste ADF para Hacktivismo . . . . .	57
Tabela 18 – Resultados Teste Ljung-Box e Jarque-Bera para Hacktivismo . . . . .	58
Tabela 19 – Resultados Teste ADF para Total de incidentes . . . . .	60
Tabela 20 – Resultados Teste Ljung-Box e Jarque-Bera para Total de incidentes . . . . .	61
Tabela 21 – Comparativo entre medidas de erros dos diferentes períodos de estimativa . . . . .	69

# Lista de abreviaturas e siglas

ACF	Função de Autocorrelação ( <i>Autocorrelation Fuction</i> )
ACF1	Autocorrelação de erros no lag 1 ( <i>Autocorrelation of errors at lag 1</i> )
ADF	Dickey-Fuller Aumentado ( <i>Augmented Dickey–Fuller</i> )
ARIMA	Modelo Auto-regressivo Integrado de Médias Móveis ( <i>Autoregressive Integrated Moving Average</i> )
IoT	Internet das Coisas ( <i>Internet of Things</i> )
MAE	Erro Médio Absoluto ( <i>Mean Absolute Error</i> )
MAPE	Erro de porcentagem média absoluta ( <i>Mean Absolute Percentage Error</i> )
MASE	Erro Médio Absoluto Escalado ( <i>Mean Absolute Scaled Error</i> )
ME	Erro Médio ( <i>Mean Error</i> )
MPE	Erro de porcentagem média ( <i>Mean Percentage Error</i> )
OIT	<i>Office of Information Technology</i>
PACF	Função de Autocorrelação Parcial ( <i>Partial Autocorrelation Function</i> )
PRC	<i>Privacy Rights Clearinghouse</i>
RMSE	Raiz do erro quadrático médio ( <i>Root Mean Square Error</i> )
SARIMA	Modelo Auto-regressivo Integrado de Médias Móveis Sazonal ( <i>Seasonal Autoregressive Integrated Moving Average</i> )
VCDB	<i>VERIS Community Database</i>
WHID	<i>Web Hacking Incident Database</i>

# Sumário

<b>1</b>	<b>INTRODUÇÃO</b>	<b>15</b>
<b>1.1</b>	<b>Objetivos</b>	<b>16</b>
<b>1.2</b>	<b>Organização do trabalho</b>	<b>16</b>
<b>2</b>	<b>FUNDAMENTAÇÃO TEÓRICA</b>	<b>17</b>
<b>2.1</b>	<b>Segurança da Informação</b>	<b>17</b>
<b>2.2</b>	<b>Séries temporais</b>	<b>18</b>
<b>2.3</b>	<b>Modelo ARIMA</b>	<b>19</b>
<b>2.4</b>	<b>Modelos Autoregressivos</b>	<b>20</b>
<b>2.5</b>	<b>Modelos de Média Móvel</b>	<b>20</b>
<b>2.6</b>	<b>Modelo SARIMA</b>	<b>20</b>
<b>2.7</b>	<b>Trabalhos correlatos</b>	<b>21</b>
<b>3</b>	<b>DESENVOLVIMENTO</b>	<b>23</b>
<b>3.1</b>	<b>Descrição geral do trabalho</b>	<b>23</b>
<b>3.2</b>	<b>Investigação das bases de dados</b>	<b>23</b>
<b>3.3</b>	<b>Coleta dos dados</b>	<b>24</b>
<b>3.4</b>	<b>Modelagem dos dados</b>	<b>25</b>
3.4.1	Definição do Modelo ARIMA	26
3.4.2	Validação do Modelo	27
3.4.2.1	Teste do Círculo Unitário dos Coeficientes	27
3.4.2.2	Teste dos Resíduos	28
3.4.2.3	Teste de Jarque-Bera	29
<b>3.5</b>	<b>Realização de Previsões com o Modelo ARIMA</b>	<b>29</b>
<b>4</b>	<b>RESULTADOS</b>	<b>32</b>
<b>4.1</b>	<b>Definição do Modelo</b>	<b>32</b>
4.1.1	Período de estimativa de 2011 a 2017	32
4.1.1.1	Crimes Cibernéticos	32
4.1.1.2	Espionagem Cibernética	35
4.1.1.3	Guerra Cibernética	38
4.1.1.4	Hacktivismo	41
4.1.1.5	Total de incidentes	44
4.1.2	Período de estimativa de 2014 a 2017	47
4.1.2.1	Crimes Cibernéticos	47
4.1.2.2	Espionagem Cibernética	50

4.1.2.3	Guerra Cibernética . . . . .	53
4.1.2.4	Hacktivismo . . . . .	56
4.1.2.5	Total de incidentes . . . . .	59
<b>4.2</b>	<b>Previsões com o Modelo . . . . .</b>	<b>62</b>
4.2.1	Previsão de 2018 a 2020 com o período de estimativa de 2011 a 2017 . . . . .	62
4.2.2	Previsão de 2018 a 2020 com o período de estimativa de 2014 a 2017 . . . . .	65
4.2.3	Comparativo de desempenho entre os dois períodos de estimativa . . . . .	68
<b>5</b>	<b>CONCLUSÃO . . . . .</b>	<b>70</b>
	<b>REFERÊNCIAS . . . . .</b>	<b>71</b>

# 1 Introdução

A necessidade de estudar a área de Segurança da Informação tem ganhado relevância visto que as organizações e os indivíduos têm se tornado gradativamente dependentes da tecnologia, seja para armazenamento e processamento de suas informações, seja para o compartilhamento dessas a terceiros (FONTE, 2017). Porém, com a evolução das funcionalidades, também há evolução de atividades maliciosas que buscam danificar ou manipular os serviços existentes por meio da exploração de vulnerabilidades de rede e software.

A quantidade de organizações que sofrem ataques é bem significativa. Na primeira coleta da pesquisa em 2017, somente no Reino Unido, mais de quatro em dez empresas (43%) e duas em dez instituições de caridade (19%) sofreram uma violação da segurança cibernética ou ataque (VAIDYA, 2018) e, em 2020, as estatísticas anteriores foram para quase metade (46%) das empresas e para mais de um quarto das caridades (26%) (VAIDYA, 2020).

Os ataques vêm de formas diversas e chegam a causar sérias consequências. O *malware* Mirai infectou dispositivos IoT como câmeras e DVRs e os inseriu uma *botnet*<sup>1</sup>. Ele teve sua primeira aparição em 31 de Agosto de 2016 e em Novembro atingiu um pico de 600.000 dispositivos infectados (ANTONAKAKIS et al., 2017). Em 12 de Maio de 2017, o WannaCry, um tipo de código malicioso conhecido como *ransomware*, infectou mais de 230.000 computadores com sistema operacional Windows em 150 países. Tal código malicioso tem como objetivo cifrar dados dos usuários e, conseqüentemente, bloquear o seu acesso ao sistema e demandar pagamentos para reverter tal situação (EHRENFELD, 2017). Existe também o episódio de Agosto de 2013 da Yahoo!, no qual foi feito um roubo de informações que chegou a afetar todos as contas da plataforma na época, estimados em cerca de 3 bilhões (Jonathan Stempel; Jim Finkle, 2017; REUTERS, 2017).

Existem diversos trabalhos que analisam dados sobre incidentes de segurança. Algumas frentes de pesquisa coletam dados ao longo dos anos para fazer a modelagem e posterior análise dos incidentes. Condon, He e Cukier (2008) extrai seus dados de uma base do *Office of Information Technology* (OIT) da *University of Maryland* na qual cada incidente inserido é verificado, então não existem falsos positivos, enquanto Edwards, Hofmeyr e Forrest (2016) usa a base *Privacy Rights Clearinghouse* (PRC) na qual são registrados incidentes que se tornaram públicos. Kuypers, Maillart e Paté-Cornell (2017) organiza os dados coletados em níveis de severidade, medido pela força de trabalho usada para investigar e remediar uma vulnerabilidade.

Porém, mesmo com o alto interesse em pesquisar sobre como esses incidentes se comportam, os dados existentes para análise são somente os que requerem divulgação pública ou

---

<sup>1</sup> Botnet é uma rede de computadores interligados pela internet cada um rodando um ou mais bots com objetivo de executar uma determinada tarefa

recebem incentivos para tal e os de grande escala expostos pela mídia (KUYPERS; MAILLART; PATÉ-CORNELL, 2017). Isso leva ao surgimento de pesquisas que não dependem da base de dados do incidente em si, mas utilizam-se de outros atributos para previsão, como Liu et al. (2015b) e Liu et al. (2015a) que usam *Machine Learning* em conjunto com das próprias atividades maliciosas originadas de organização para análise.

## 1.1 Objetivos

O objetivo do trabalho é analisar e fazer previsões de incidentes de segurança usando uma base de dados pública. O foco da análise é a modelagem dos dados usando séries temporais, já que podem ser aplicadas em uma variedade de tipos de incidentes e usam movimentos passados para prever comportamentos futuros (CONDON; HE; CUKIER, 2008).

Nesse estudo foi utilizada a base de dados *Hackmageddon* (PASSERI, 2018) que mantém de forma independente vários incidentes de segurança organizados por mês e tipo de ataque. Utilizando como ferramenta para estudo o *R*, foram estimados modelos ARIMA (*Modelo Auto-Regressivo Integrado de Médias Móveis*) por lidarem bem com séries temporais e, posteriormente, foram feitas previsões usando os mesmos.

## 1.2 Organização do trabalho

Esse trabalho está estruturado da seguinte forma: no Capítulo 2 é apresentada uma fundamentação teórica necessária para o melhor entendimento do assunto. Os passos tomados para o desenvolvimento da pesquisa foram detalhados no Capítulo 3. No Capítulo 4, é feita a exibição e análise dos resultados encontrados. E, finalmente, no Capítulo 5, são feitas conclusões e sugestões de trabalhos futuros.

## 2 Fundamentação Teórica

Serão definidos nesse capítulo do trabalho alguns conceitos básicos sobre segurança da informação, incidentes de segurança, séries temporais, modelos ARIMA e trabalhos correlatos.

### 2.1 Segurança da Informação

Solms e Niekerk (2013) definem que o papel da segurança da informação é garantir a continuidade dos negócios e minimizar os danos a estes, limitando o impacto de incidentes de segurança. A área de segurança da informação consiste de medidas para desviar, prevenir, detectar e corrigir violações de segurança que envolvam a transmissão de informações (STALLINGS, 2014).

Existe uma série de conceitos que envolvem os objetivos fundamentais da segurança para serviços e dados computacionais, tais conceitos são conhecidos como a tríade CIA. A sigla CIA, do acrônimo em inglês, *confidenciability, integrity and availability* foi definida pelos padrões FIPS 199 da NIST como:

- Confidencialidade: preservar restrições sobre acesso e divulgação de informação, incluindo meios para proteger a privacidade de indivíduos e informações privadas.
- Integridade: prevenir-se contra a modificação ou destruição imprópria de informação, incluindo a irretratabilidade e autenticidade dela.
- Disponibilidade: assegurar acesso e uso rápido e confiável da informação.

Esses três itens são considerados objetivos a serem alcançados para se ter uma segurança bem estabelecida de dados e sistemas de informação (STALLINGS, 2014). Entretanto há frentes no campo da segurança que percebem que sejam considerados também conceitos adicionais para tal, como os seguintes:

- Autenticidade: a propriedade de ser genuíno e capaz de ser verificado e confiável.
- Responsabilização: a meta de segurança que gera o requisito para que ações de uma entidade sejam atribuídas exclusivamente a ela.

Quaisquer eventos que violem os requisitos citados acima são considerados ataques (STALLINGS, 2014) e podem ser catalogados como um incidente de segurança. Portanto, um incidente de segurança pode ser definido como qualquer evento adverso, confirmado ou sob suspeita, relacionado à segurança de sistemas de computação ou de redes de computadores

(Cert.br, 2020). Já Gualberto et al. (2013) definem incidente de segurança como um evento que tem alta probabilidade de impactar o negócio e segurança de uma organização E esses surgem a partir de explorações de vulnerabilidades de segurança. Exemplos de vulnerabilidade são falhas de projetos, na implementação ou configuração de programas, ou seja, qualquer condição que possa ser explorada em um ataque (Cert.br, 2020)

## 2.2 Séries temporais

De acordo com Hyndman e Athanasopoulos (2013), uma série temporal (como exemplificado na Figura 1) é um conjunto de observações, medidas de forma sistemática, segundo uma periodicidade regular, ou seja, o registro de observações ao longo ao tempo.

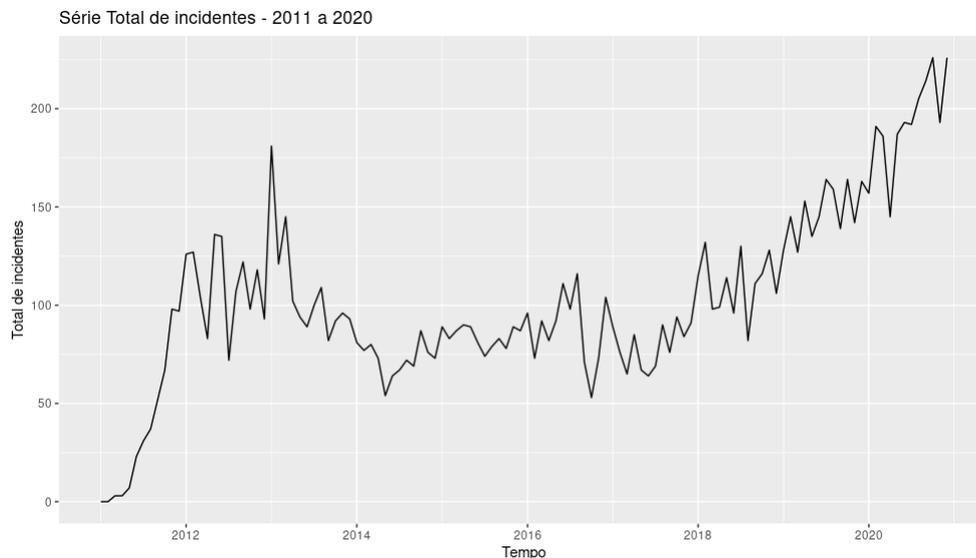


Figura 1 – Exemplo de Série Temporal. Fonte: Do autor.

Reis (2016) descreve um modelo clássico de série temporal com quatro componentes principais:

- Tendência (T): comportamento de longo prazo da série, que pode ser causada pelo crescimento demográfico, ou mudança gradual de hábitos de consumo, ou qualquer outro aspecto que afete a variável de interesse no longo prazo;
- Variações cíclicas ou ciclos (C): flutuações nos valores da variável com duração superior a um ano, e que se repetem com certa periodicidade, que podem ser resultado de variações da economia como períodos de crescimento ou recessão, ou fenômenos climáticos;
- Variações sazonais ou sazonalidade (S): flutuações nos valores da variável com duração inferior a um ano, e que se repetem todos os anos, geralmente em função das estações do ano (ou em função de feriados e festas populares, por exemplo);

- Variações irregulares (I), que são as flutuações inexplicáveis, resultado de fatos fortuitos e inesperados como catástrofes naturais, atentados terroristas ou decisões de governos.

Uma das aplicações da análise de séries temporais é identificar padrões não aleatórios de uma variável de interesse, e com a observação deste comportamento passado, fazer previsões sobre o futuro. Para se iniciar a estudar uma série temporal deve-se analisar duas características principais, detalhadas a seguir.

1. Estacionariedade: uma série é dita estacionária quando ela se desenvolve no tempo aleatoriamente ao redor de uma média constante, refletindo alguma forma de equilíbrio estável, ou seja, suas leis probabilísticas não se alteram ao longo do tempo (CRYER; CHAN, 2008).
2. Sazonalidade: existe quando uma série é influenciada por fatores sazonais, que ocorrem ano, mês ou dia da semana (HYNDMAN; ATHANASOPOULOS, 2013).

Uma das formas de se fazer isso é estimar um modelo ARIMA, que é descrito nas próximas seções, e usá-lo para prever os próximos eventos da série.

## 2.3 Modelo ARIMA

ARIMA é um acrônimo para o inglês *Auto Regressive Integrated Moving Average*, ou seja, ele é a combinação dos modelos de auto-regressão com os modelos de média móvel e diferenciação (HYNDMAN; ATHANASOPOULOS, 2018). O modelo pode ser descrito da seguinte forma:

$$y_t = c + \phi_1 y_{t-1} + \dots + \phi_p y_{t-p} + \theta_1 \varepsilon_{t-1} + \dots + \theta_q \varepsilon_{t-q} + \varepsilon_t \quad (2.1)$$

onde  $\varepsilon_t$  é o ruído branco,  $\phi_p$ , os coeficientes autoregressivos,  $\theta_q$ , os coeficientes de média móvel,  $y_t$ , os dados sobre os quais o ARIMA é aplicado e  $c$ , uma constante que determina a tendência da série. Isso caracteriza um **modelo ARIMA(p,d,q)**, onde:

$$\begin{array}{l} p = \text{ordem da parte auto-regressiva} \\ d = \text{grau de primeira diferenciação envolvida} \\ q = \text{ordem da parte de média móvel} \end{array}$$

Quando tem-se a modelagem do ARIMA, a diferença entre os valores observados da série e os valores ajustados do modelo é chamada de resíduo do modelo. Quando esse não apresenta autocorrelação e é normalmente distribuído, o mesmo é chamado de ruído branco (HYNDMAN; ATHANASOPOULOS, 2018).

A análise da estacionariedade é relevante para decidir se será usada a série em nível ou com uma de suas diferenças na estimativa do modelo, ou seja, para definir o valor do parâmetro  $d$ . Diferenciar a série pode ajudar a estabilizar a média da série temporal e, conseqüentemente, eliminar (ou reduzir) tendências e sazonalidades (HYNDMAN; ATHANASOPOULOS, 2018). Se existir um fator de sazonalidade sobre a série, modelo ARIMA sazonal deve ser usado, também conhecido como SARIMA.

## 2.4 Modelos Autoregressivos

Num modelo auto-regressivo, a previsão de uma variável de interesse é alcançada usando a combinação dos valores passados da mesma. Um modelo auto-regressivo de ordem  $p$  pode ser escrito da seguinte forma:

$$y_t = c + \phi_1 y_{t-1} + \phi_2 y_{t-2} + \dots + \phi_p y_{t-p} + \varepsilon_t \quad (2.2)$$

onde  $\varepsilon_t$  é o ruído branco,  $\phi_p$ , os coeficientes autoregressivos,  $y_t$ , os dados sobre os quais o ARIMA é aplicado e  $c$ , uma constante que determina a tendência da série. Isso caracteriza um **modelo AR( $p$ )**, um modelo auto-regressivo de ordem  $p$  (HYNDMAN; ATHANASOPOULOS, 2018).

## 2.5 Modelos de Média Móvel

Ao invés de usar valores passados da variável a ser prevista numa regressão, um modelo de média móvel usa erros de previsão passados em um modelo de regressão.

$$y_t = c + \varepsilon_t + \theta_1 \varepsilon_{t-1} + \dots + \theta_q \varepsilon_{t-q} + \varepsilon_t \quad (2.3)$$

onde  $\varepsilon_t$  é o ruído branco,  $\theta_q$ , os coeficientes de média móvel,  $y_t$ , os dados sobre os quais o ARIMA é aplicado e  $c$  uma constante que determina a tendência da série. Isso caracteriza um **modelo AM( $q$ )**, um modelo média móvel de ordem  $q$  (HYNDMAN; ATHANASOPOULOS, 2018).

## 2.6 Modelo SARIMA

O modelo sazonal ARIMA é formado incluindo termos sazonais ao ARIMA que já vimos. Ele é escrito da seguinte forma:

$$ARIMA \quad \underbrace{(p, d, q)}_{\text{Parte não sazonal}} \quad \underbrace{(P, D, Q)_m}_{\text{Parte sazonal}} \quad (2.4)$$

onde  $m$  = número de observações por ano. A parte sazonal do modelo consiste em termos que são similares aos componentes não sazonais do modelo, mas envolvem retrocesso do período sazonal (HYNDMAN; ATHANASOPOULOS, 2018).

## 2.7 Trabalhos correlatos

Essa seção da fundamentação teórica irá apresentar trabalhos que tiveram como objetivo a análise de eventos relacionados a segurança cibernética, com dados relacionados a vulnerabilidades de sistema ou sobre a exploração das mesmas e, posteriormente, mostrar as similaridades dos trabalhos apresentados com o trabalho proposto.

Condon, He e Cukier (2008) investigam incidentes de segurança utilizando séries temporais, já que as mesmas não dependem de fatores causais e podem ser aplicadas a qualquer tipo de incidente por só depender de eventos passados para prever comportamentos futuros. Os autores se mostraram mais focados em analisar se os dados tinham alguma tendência ou sazonalidade ao invés de somente tentar prever um número específico incidentes, visto que entender o que leva a exploração de uma vulnerabilidade pode ajudar a se realizar previsões sobre um total de incidentes. Também foi feito um comparativo entre as previsões de séries temporais e modelos de crescimento de confiabilidade de software, visto que, por mais que um incidente e uma falha de software não sejam a mesma coisa, existem similaridades. Os modelos de crescimento de confiabilidade de software geraram melhores previsões quando aplicados em dados de tipo individuais de incidentes, enquanto as séries temporais geraram melhores previsões na análise do total agregado resultante da combinação de todos os tipos de incidentes

Liu et al. (2015b) propõe o desenvolvimento de um método de previsão de incidentes de segurança com somente com as propriedades observadas externamente sobre a rede de uma organização. Para isso, foi utilizado um classificador de Floresta Randômica que alcançou resultados de relativamente alta acurácia, com 90% de resultados positivos verdadeiros e 10% de falsos positivos. Por mais que o aprendizado de máquina chegou num resultado muito positivo, ainda se ressalta que a comunidade acadêmica se beneficiaria bastante se existisse um meio mais sistemático e uniforme de reporte de incidentes de segurança.

Edwards, Hofmeyr e Forrest (2016) fazem a análise de tendências de violação de dados ao longo do tempo, levando em consideração violações maliciosas (atacantes visam o roubo das informações) ou negligentes (informações confidenciais são vazadas acidentalmente) com o uso de Modelos Lineares Generalizados Bayesianos. Eles chegam a conclusão que o número de ataques no período analisado se manteve constante, isso pode ser explicado pelo fato da base de dados utilizada depender de um conjunto de dados que foi publicamente reconhecido e, por isso, não conterem todos os vazamentos de dados que aconteceram ou pelas tratativas de segurança avançarem em conjunto com os níveis de ataques. Já Kuypers, Maillart e Paté-Cornell (2017) verificam a tendência de incidentes de seguranças baseados em nível de severidade e apesar

das grandes apostas associadas com a segurança cibernética, a falta de dados tem dificultado a análise, monitoração e previsão de ataques.

Xu et al. (2018) buscam responder se as violações de dados estão aumentando, diminuindo ou estabilizando ao longo do tempo, utilizando um modelo ARMA-GARCH para descrever a evolução de violações, sendo ARMA um modelo auto-regressivo de média móvel e GARCH um modelo de heterocedasticidade condicional auto-regressiva generalizado. Como resultado de sua pesquisa, tem-se uma piora em termos de frequência de violações de dados, porém não em termos de magnitude do dano causado.

Kesan e Zhang (2020) propõem uma maneira de categorizar incidentes de segurança baseado nas suas perdas. Os tipos de perda foram vistos como variáveis randômicas de Bernoulli e a partir dessas variáveis são usadas técnicas de agrupamento para separá-las em 4 categorias, que abrangem 90% da variação em tipos de perda e também analisada a dependência entre eles.

Apesar de Condon, He e Cukier (2008), Liu et al. (2015b), Edwards, Hofmeyr e Forrest (2016), Kuypers, Maillart e Paté-Cornell (2017) e Xu et al. (2018) se proporem a estudar tendências de incidentes de segurança ao longo do tempo, são observadas várias metodologias para se fazer isso. O trabalho proposto utiliza características semelhantes a alguns dos trabalhos apresentados. Foi feita uma análise baseada em categorias como em Kesan e Zhang (2020), porém, ao invés de ocorrer a separação baseada em severidade de perda, foi empregado o tipo de incidente. Também, como meio de se analisar e prever incidentes, foi usado uma espécie de modelo auto-regressivo de média móvel como em Xu et al. (2018) e Condon, He e Cukier (2008) ao contrário de Liu et al. (2015b), Edwards, Hofmeyr e Forrest (2016) e Kuypers, Maillart e Paté-Cornell (2017) que usam outros modelos estatísticos.

## 3 Desenvolvimento

Nessa capítulo será dado uma visão geral sobre as etapas tomadas para o desenvolvimento do trabalho e, posteriormente, detalhadas cada uma dessas.

### 3.1 Descrição geral do trabalho

O trabalho foi dividido e realizado nas seguintes etapas:

1. Investigação das bases de dados: busca de bases candidatas ao estudo, análise das mesmas e seleção da mais apropriada ao estudo.
2. Coleta dos dados da base escolhida: extração dos dados da fonte selecionada e processamento dos mesmos para ficarem mais fáceis de manipular em R;
3. Modelagem da base de dados usando séries temporais em conjunto com o modelo ARIMA: criação de gráficos e extração de métricas em R para escolher os parâmetros necessários para a definição do modelo ARIMA;
4. Realização de previsões de incidentes usando o modelo ARIMA: análise de gráficos e métricas geradas com o R a fim de determinar o desempenho dos modelos;
5. Discussão dos resultados obtidos nos passos anteriores.

Nas seções [3.2](#), [3.3](#), [3.4](#) e [3.5](#) serão explicadas com mais detalhes cada uma das etapas, exceto a última pois a discussão de resultados será feita ao longo dos capítulos [4](#) e [5](#).

### 3.2 Investigação das bases de dados

Dentro da primeira etapa da metodologia proposta, foram selecionadas e investigadas quatro bases de dados para avaliar qual é a mais apropriada ao trabalho.

Após uma investigação acerca de bases de dados públicas sobre incidentes de segurança, as seguintes bases foram encontradas; PRC, VCDB, WHID e *Hackmageddon*. As principais características são descritas a seguir:

*PRC - Privacy Rights Clearinghouse*

Contém informações sobre violações de dados que foram relatadas nos EUA, sendo que os dados possuem: data de divulgação da violação, nome da entidade responsável pelos dados, classificação do tipo de violação, número total de dados violados, localização da entidade, informação da fonte dos dados e descrição da violação ([PRC, 2019](#)).

#### *VCDB - VERIS Community Database*

Tem como objetivo é coletar e disseminar informações de violação de dados para todas as violações de dados divulgadas publicamente, seus dados são codificados no formato VERIS (um conjunto de métricas que descrevem incidentes de segurança em um ambiente estruturado e repetitivo) e também é fornecida uma visualização interativa disponível para uso público dos dados ([VERIZON, 2019](#)).

#### *WHID - Web Hacking Incident Database*

Rastreia apenas incidentes de segurança relatados pela mídia que podem ser associados a um aplicativo de página *Web*, ou seja, vulnerabilidades de sistema operacionais ou de rede não são inclusos na base. Entre esses incidentes tem-se dois tipos: violação e divulgação, sendo o primeiro tipo um incidente que o site foi comprometido e o segundo um incidentes em que um pesquisador publicou uma vulnerabilidade de algum site ([WASC, 2019](#)).

#### *Hackmageddon - Information Security Timelines and Statistics*

Contabiliza ataques de fato realizados, não vulnerabilidades. Os dados são armazenados em tabelas, sendo possível cadastrar informações sobre um ataque a partir de um formulário no site de controle do banco ([PASSERI, 2018](#)).

A base de dados selecionada para a análise foi o *Hackmageddon* por seus dados serem independentes e não dependerem de divulgação pública, por se tratar de uma base de incidentes de segurança e não só vulnerabilidades e, também, classificar seus incidentes por tipo.

### 3.3 Coleta dos dados

A partir da base escolhida, nessa segunda etapa foi feita a análise, pré-processamento e posterior coleta dos dados.

O *Hackmageddon* ([PASSERI, 2018](#)), disponibiliza os dados em forma de várias tabelas em arquivos *.xlsx*. As tabelas contém todos os ataques de segurança de um determinado período quinzenal e para cada ataque são guardadas as seguintes informações: Data, Autor, Alvo, Descrição, Forma do Ataque, Categoria do Alvo, Categoria do Ataque e País (como mostrado na Figura 2. Nesse trabalho, somente será usado para a análise a quantidade de incidentes e as categorias de ataque.

ID	Date	Author	Target	Description	Attack	Target Class	Attack Class	Country	Link	Tags
1	12/2/2019	?	RiverKids Pediatric Home Health	RiverKids Pediatric Home Health is affected by a hacking incident.	Unknown	Q Human health and social work activities	CC	US	<a href="https://www.bleepingcomputer.com/news/security/criminals-pull-hard-before-xmas-attack-us-health-industry/">https://www.bleepingcomputer.com/news/security/criminals-pull-hard-before-xmas-attack-us-health-industry/</a>	RiverKids Pediatric Home Health
2	12/11/2019	?	Arrigo Automotive Group	Arrigo Automotive Group is hit by a ransomware attack costing up to \$250,000.	Malware	H Transportation and storage	CC	US	<a href="https://heliontechnologies.com/2020/01/02/automotive-news-ransomware-is-an-escalating-problem-for-dealers/">https://heliontechnologies.com/2020/01/02/automotive-news-ransomware-is-an-escalating-problem-for-dealers/</a>	Arrigo Automotive Group, ransomware

Figura 2 – Exemplo de tabela da base dados *Hackmageddon*. Fonte: Extraída de (PASSERI, 2018).

Foi selecionado para o estudo o período de 2011 a 2020 no qual os dados já estão classificados em quatro categorias de ataque, sendo elas as seguintes:

- *Cyber Crime* (Crime Cibernético): toda a atividade criminosa que se utiliza de um computador ou rede de computadores como instrumento base do ataque;
- *Cyber Warfare* (Guerra Cibernética): envolve as ações de um estado-nação ou organização internacional para atacar e tentar danificar computadores ou redes de informação de outra nação;
- *Cyber Espionage* (Espionagem Cibernética): obter segredos e informações sem a permissão e conhecimento do detentor das informações;
- *Hacktivism* (Hacktivismo): uso da tecnologia para promover uma agenda política ou uma mudança social.

Foi feita a coleta de todas as tabelas referente ao período do estudo e, a fim de facilitar a modelagem, os dados foram pré processados, tomando os passos a seguir:

1. Junção das tabelas quinzenais de cada ano para a melhor contabilização do número de ataques;
2. Produção de tabelas com a quantidade de incidentes para cada mês do período analisado, sendo cinco tabelas no total, uma para cada categoria de ataque descrita anteriormente e uma tabela para todos os incidentes;
3. Conversão das tabelas em arquivos .csv, para melhor administração em R.

### 3.4 Modelagem dos dados

Já na terceira etapa, os dados resultantes da coleta do passo anterior foram usados para a construção de modelos auto-regressivos de médias móveis genéricos (ARIMA).

### 3.4.1 Definição do Modelo ARIMA

Os modelos foram criados com o auxílio do R, um ambiente de software livre para computação estatística e criação de gráficos. Para definir os parâmetros do ARIMA, foram usados os dados da série de 2011 a 2017 e para validação dos modelos foi utilizado o período de 2018 a 2020.

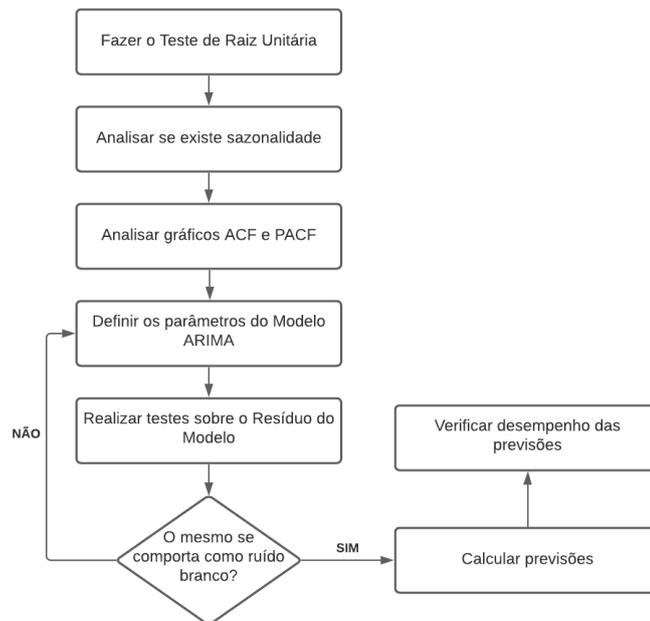


Figura 3 – Etapas de Definição de um Modelo ARIMA Fonte: Do autor.

A Figura 3 mostra, de maneira geral, quais são as etapas da construção e validação de um modelo ARIMA. Um dos primeiros passos a se tomar é fazer os testes de raiz unitária para verificar a estacionariedade da série.

Para fazer esse testes foi utilizado o teste de *Dickey-Fuller* Aumentado, também conhecido como teste ADF, que tem como objetivo checar se existe raiz unitária numa sequência (DICKKEY; FULLER, 1979). O mesmo usa como hipótese nula a estacionariedade da série, ou seja, para um valor baixo ( $p < 0.05$ ) a série é estacionária e não tem raiz unitária e valores maiores que definido ( $p > 0.05$ ) indicam a presença de uma raiz unitária na sequência, logo a série é não estacionária (MUSHTAQ, 2011). A função do R correspondente a esse teste é a *adfTest*. Se a série não for estacionária, deve-se repetir o teste para a sua diferença para ver se ela é apta a ser usada para definir o modelo.

É importante também fazer a análise da existência ou não de tendências sazonais que podem influenciar as séries usadas. Para isso foi usada uma função do R chamada *seas*, *Seasonal Adjustment by X-13ARIMA-SEATS*, que tenta fazer um ajuste sazonal na série. Em casos de sucesso, isso indica que a série tem uma tendência sazonal e o SARIMA é modelo mais indicado para descrever a série. Se nenhum ajuste for possível, não há sazonalidade e um ARIMA é

suficiente (SAX; EDELBUETTEL, 2018).

Para a definição dos parâmetros  $p$  e  $q$  do modelo ARIMA, foram analisadas as funções de autocorrelação parcial (PACF) e total (ACF) das séries e, caso necessário, das suas diferenças. É possível a visualização dos gráficos dessas funções usando a função *ggtsdisplay* no R (como evidenciado na Figura 4). Com a análise desses, foram estimados os parâmetros não sazonais e, caso necessário, sazonais do ARIMA.

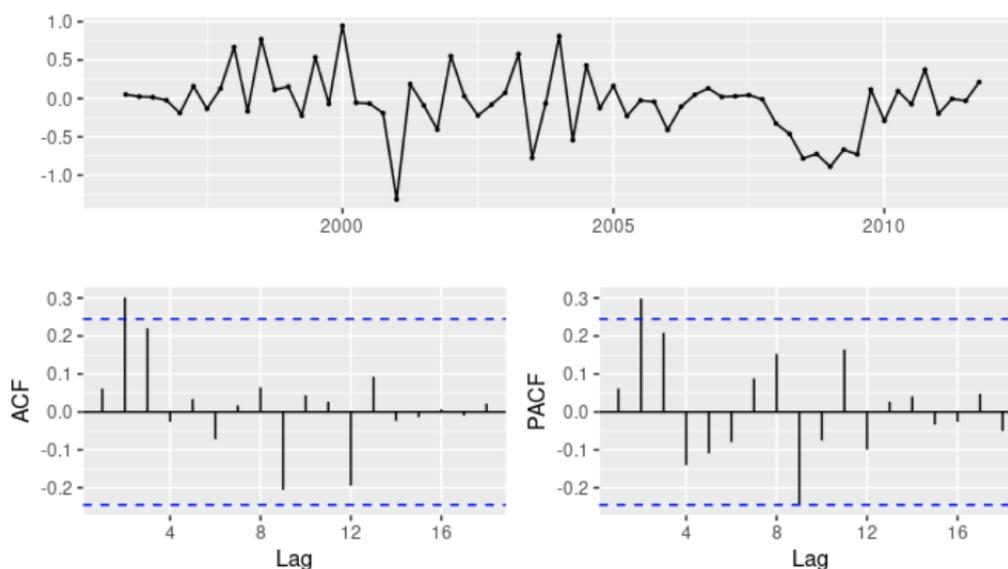


Figura 4 – Exemplo do retorno da função *ggtsdisplay* do modelo. Fonte: Extraída de (HYNDMAN; ATHANASOPOULOS, 2018).

### 3.4.2 Validação do Modelo

Com os parâmetros do modelo ARIMA definidos, o modelo foi gerado com a função *Arima* do R. Posteriormente, seus coeficientes foram testados com o Teste do Círculo Unitário dos Coeficientes para confirmar e os mesmos são estáveis (não explosivos), caso contrário eles poderiam impactar a estacionariedade do modelo. Também foram feitos testes sobre seus resíduos para verificar se os mesmos se comportam como ruído branco, entre eles a análise do gráfico ACF, o teste Ljung-Box e o teste Jarque-Bera. Com esses testes é possível definir se a série teve a maioria de suas particularidades capturadas pelo modelo que foi definido, caso contrário, deve-se rever os parâmetros do modelo.

#### 3.4.2.1 Teste do Círculo Unitário dos Coeficientes

Esse teste é facilmente verificado no R, com o *autoplot* do modelo é gerado um gráfico como o da Figura 5.

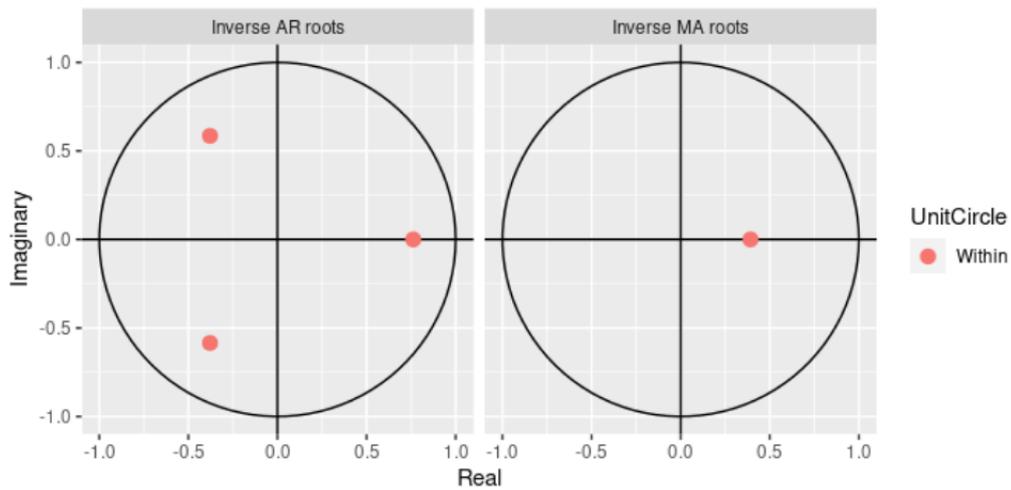


Figura 5 – Exemplo do retorno da função *autoplot* do modelo. Fonte: Extraída de (HYNDMAN; ATHANASOPOULOS, 2018).

Os pontos vermelhos no gráfico a esquerda correspondem às raízes dos polinômios  $\phi$  (parte auto-regressiva do modelo) enquanto o ponto no gráfico da direita corresponde à raiz de  $\theta$  (parte de média móvel do modelo). Como todos os pontos estão dentro do círculo unitário, o modelo estimado é estacionário.

#### 3.4.2.2 Teste dos Resíduos

A função *checkresiduals* do R constrói o gráfico ACF e o histograma dos resíduos: quando a maioria das autocorrelações, ou todas elas, estão dos limites mostrados no gráfico do ACF indica que o resíduo está se comportando como ruído branco.

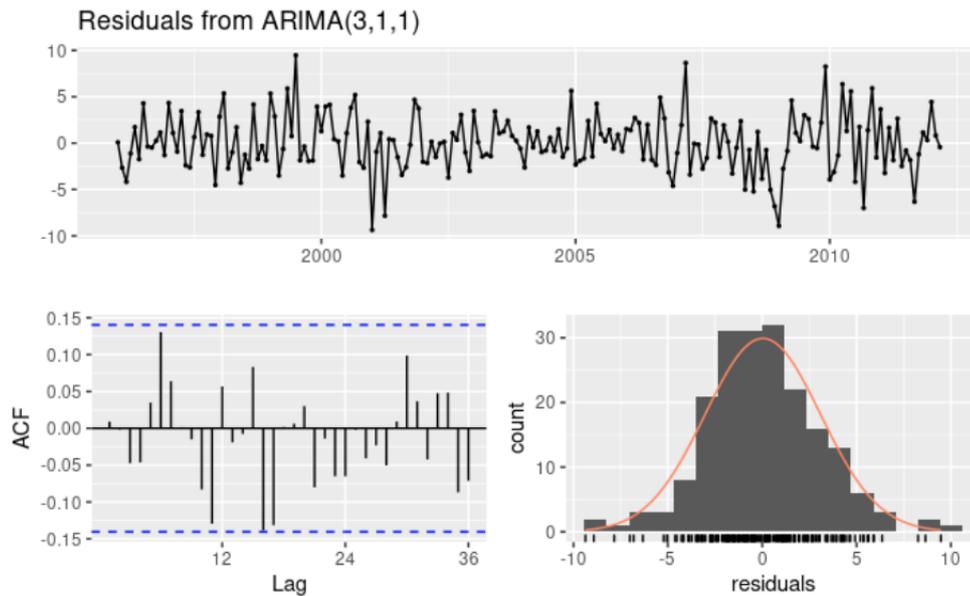


Figura 6 – Exemplo do retorno da função *checkresiduals* do modelo. Fonte: Extraída de (HYNDMAN; ATHANASOPOULOS, 2018).

O teste Ljung-Box, também executado na função *checkresiduals*, verifica se os dados são independentemente distribuídos. O mesmo usa como hipótese nula que os dados não são independentes, ou seja, para valores pequenos ( $p < 0.05$ ) os dados estão correlacionados e valores de maiores do que o limiar estimado para o teste ( $p > 0.05$ ) há a independência dos mesmos.

#### 3.4.2.3 Teste de Jarque-Bera

O teste Jarque-Bera se propõe a analisar a normalidade dos dados analisados, ele foi executado para verificar se os resíduos seguem uma distribuição normal. O mesmo usa como hipótese nula que os dados são normalmente distribuídos, ou seja, para resultados muito pequenos ( $p < 0.05$ ) é rejeitada a normalidade (JARQUE; BERA, 1980). A função do R correspondente a esse teste é *jarque.bera.test(residuals(modeloARIMA))*.

### 3.5 Realização de Previsões com o Modelo ARIMA

Como o modelo foi estimado usando os dados de 2011 a 2017, foram selecionados os anos de 2018 a 2020 para realizar previsões a fim de validar o mesmo a partir dos parâmetros obtidos anteriormente. As previsões foram calculadas seguindo os passos a seguir:

1. Estimação do modelo;
2. Armazenamento dos Parâmetros obtidos da Estimação;
3. Execução do Arima novamente, porém passando os coeficientes obtidos anteriormente e a janela temporal de teste;

#### 4. Criamos uma série prevista a partir do modelo gerado no passo anterior.

As funções que foram utilizadas em R referentes aos passos anteriores foram evidenciadas no algoritmo 1 a seguir.

---

##### Pseudocódigo - Geração da Previsão

---

```

modeloARIMA  $\leftarrow$  Arima(ts-treinamento, ordem)
coeficientes  $\leftarrow$  coefficients(modeloARIMA)
modeloARIMApredicao  $\leftarrow$  Arima(ts-teste, ordem, coeficientes)
predicao  $\leftarrow$  fitted(modeloARIMApredicao)

```

---

O parâmetro *ts-treinamento* consiste na série temporal referente ao período definido para a estimação do modelo e a *ordem* se refere aos parâmetros (p,d,q) do modelo Arima, esses parâmetros são passados na função *Arima()* para definir o modelo que será armazenado numa variável *modeloARIMA*. A função *coefficients()* recebe a variável *modeloARIMA* definida no passo anterior e os coeficientes do modelo calculados são guardados numa outra variável, a *coeficientes*. A variável *modeloARIMApredicao* guarda o resultado da terceira função na qual são usados o parâmetro *ts-teste* que se refere a série do período definido para a validação do modelo, a mesma *ordem* da primeira função e os *coeficientes* calculados anteriormente. O último passo é gerar uma série com os valores de previsão do modelo com a função *fitted()* e para isso ela recebe a variável *modeloARIMApredicao* e, por fim, numa variável *predicao* fica armazenada essa série.

Após a geração a série de previsão, a mesma foi colocada em um gráfico em conjunto com os dados observados e, posteriormente, foi medida sua acurácia, feita pela função *accuracy* do R, na qual é feita uma comparação da série gerada com a que foi separada para a validação do modelo. Essa função tem como saída as métricas: Erro Médio (ME), Raiz do erro quadrático médio (RMSE), Erro Médio Absoluto (MAE), Erro de porcentagem média (MPE), Erro de porcentagem média absoluta (MAPE), Erro Médio Absoluto Escalado (MASE) e Autocorrelação de erros no lag 1 (ACF1). Para o trabalho, foram levadas em consideração para análise somente as métricas ME, MAE e MAPE por serem se mais fácil entendimento.

A métrica de Erro Médio (ME) consiste simplesmente na média da série de erros de previsão enquanto a Erro Médio Absoluto (MAE) é dada pela média dos erros em valores absolutos. Finalmente, o Erro de porcentagem média absoluta (MAPE), pertence ao grupo de medidas que se baseiam no erro percentual de previsão e é dado pela média do erro percentual em valor absoluto. Consequentemente estas medidas têm a vantagem de não serem dependentes de escala e podem assim serem utilizadas para comparar o poder de previsão utilizando dados com diferentes escalas (LIMA, 2015).

A ideia é que se faça o comparativo entre as medidas dos modelos, quando se há a minimização dessas medidas significa que o modelo da medida de menor valor de comportou melhor. Em especial o MAPE expressa a porcentagem média dos erros, sua desvantagem é que

se existem séries que passam pelo valor zero, seu valor tenderá ao infinito (MIRANDA, 2014) e para essas não é possível verificar o percentual de erro do modelo, somente comparar se ele foi melhor que outro.

## 4 Resultados

Nesse capítulo serão mostrados os resultados da execução das etapas de Definição de Modelo e Previsões descritas no Capítulo 3.

### 4.1 Definição do Modelo

Foram definidos modelos ARIMA usando dois períodos de teste para se analisar qual teria melhor performance, ou seja, com qual período se alcançaria modelos com previsões de maior desempenho. E se, mesmo com um período menor, se o modelo ainda faria previsões relevantes. Em ambos períodos de estimação, a função *seas()* recomendou não fazer o ajuste sazonal nas séries pois não existem picos sazonais visualmente significativos, então todas foram modeladas usando o modelo ARIMA sem sazonalidade. Para séries que já são estacionárias, devido a análise dos gráficos ACF e PACF e pelo teste ADF com suas variáveis em nível, não foi realizado o teste ADF para sua 1ª Diferença nem calculado os gráficos de ACF e PACF para a mesma.

#### 4.1.1 Período de estimativa de 2011 a 2017

A seguir seguem os resultados dos passos de definição do modelo para as 4 categorias de incidentes e para o total de incidentes usando o período de treinamento de 2011 a 2017.

##### 4.1.1.1 Crimes Cibernéticos

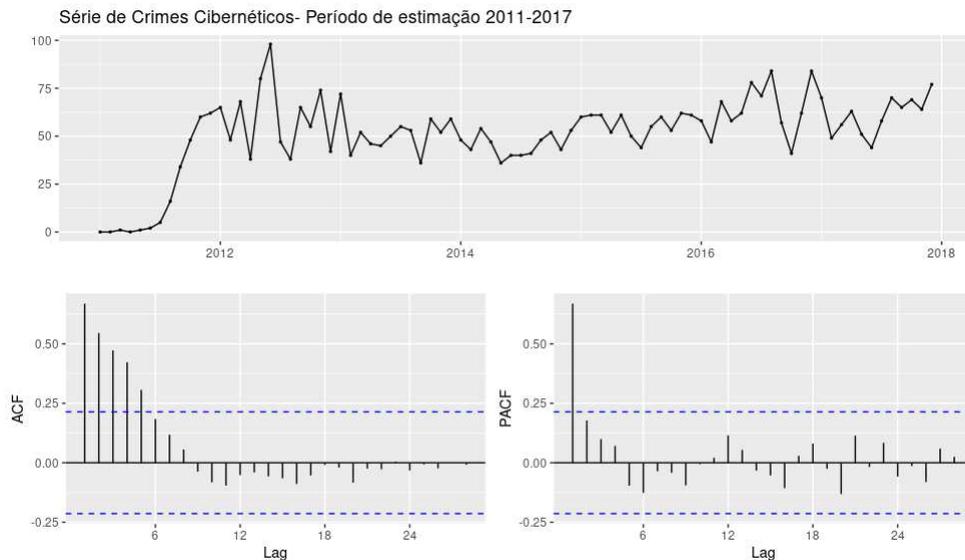


Figura 7 – Série temporal de Crimes Cibernéticos e seus ACF e PACF. Fonte: Do autor.

Tabela 1 – Resultados Teste ADF para Crimes Cibernéticos

	Série em Nível	1ª Diferença
$p$	0.020	-
Dickey-Fuller	-3.296	-

Os gráficos das Figuras 7 mostram a possibilidade da série não ser estacionária, mas com o resultado do teste ADF ( $p < 0.05$ ) evidenciado na Tabela 1 foi verificado que a série em nível já é estacionária. Com isso, parâmetro de  $d$  do modelo é igual a 0 e, como o intervalo de confiança do ACF e PACF da série em nível ser rompido logo no primeiro lag (Figura 7), foi definido o valor 1 tanto para o parâmetro  $q$  quanto para o parâmetro  $p$ , resultando num modelo ARIMA(1,0,1). Em sequência, verificou-se a estabilidade dos coeficientes do modelo com o Teste do Círculo Unitário dos Coeficientes (Figura 8) e a normalidade de seus resíduos vendo se os valores de lag no gráfico ACF dos mesmos estão em maior parte dentro do intervalo de confiança (Figura 9) e analisando os resultados do Teste Jarque-Bera e o Teste Ljung-Box (Tabela 2).

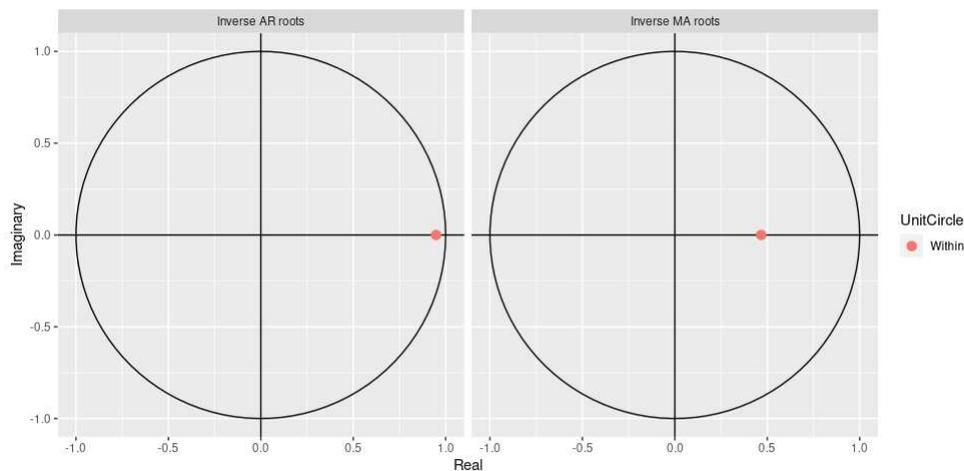


Figura 8 – Teste de Círculo Unitário para Série de Crimes Cibernéticos: ARIMA(1,0,1). Fonte: Do autor.

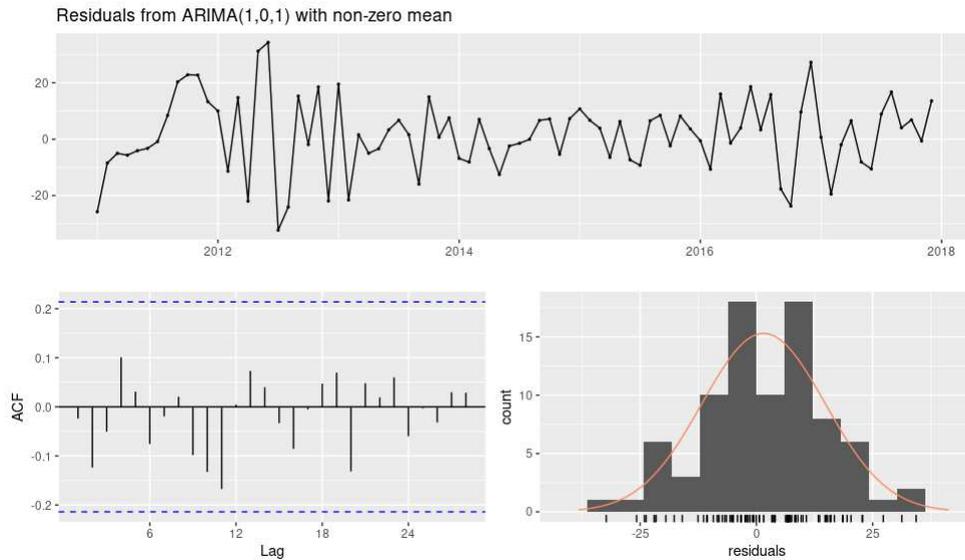


Figura 9 – Teste do Resíduo para o modelo estimado para Série: ARIMA(1,0,1). Fonte: Do autor.

Tabela 2 – Resultados Teste Ljung-Box e Jarque-Bera para Crimes Cibernéticos

	ARIMA(1,0,1)
Ljung-Box	$p = 0.743$
Jarque-Bera	$p = 0.902$

No teste de Círculo Unitário dos Coeficientes da Figura 8, os valores dos coeficientes estarem dentro do círculo unitário indica que os mesmos são estáveis, confirmando que o modelo estimado é estacionário. No gráfico ACF contidos na Figura9, observa-se que todos os valores de lag estão dentro do intervalo de confiança e na Tabela 2 como ambos os valores  $p$  rejeitam a hipótese nula dos resíduos não serem normalmente distribuídos ( $p > 0.05$ ), pode-se assumir que o resíduo do modelo se comporta como ruído branco.

Na Figura 10, tem-se a curva dos valores ajustados do modelo consegue acompanhar as tendências os dados observados de forma satisfatória.

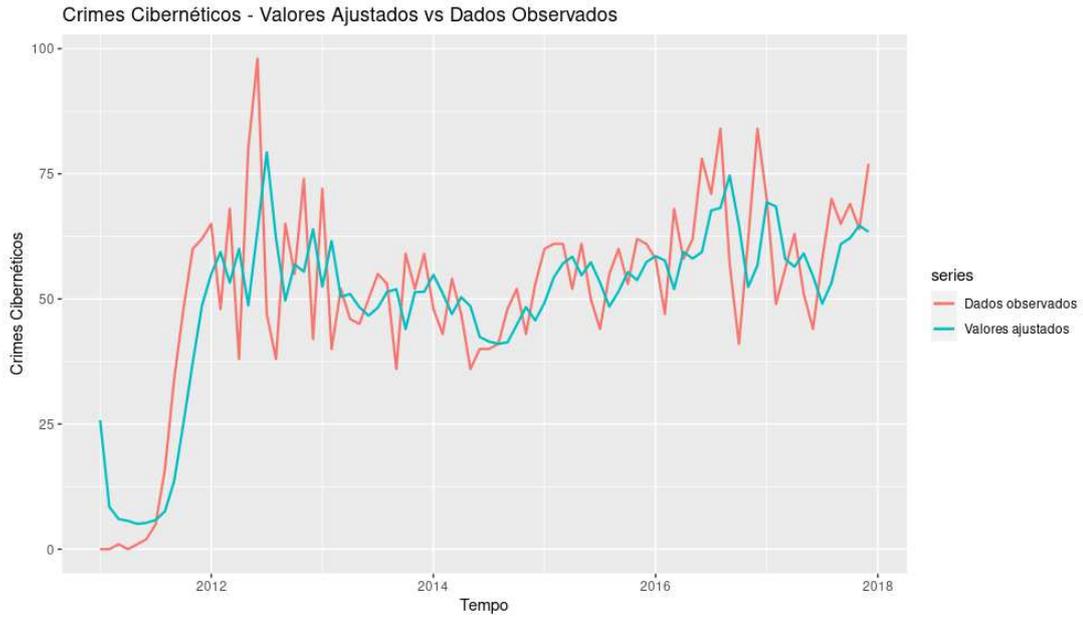


Figura 10 – Crimes Cibernéticos - Valores Ajustados vs Dados Observados: ARIMA(1,0,1).  
 Fonte: Do autor.

4.1.1.2 Espionagem Cibernética

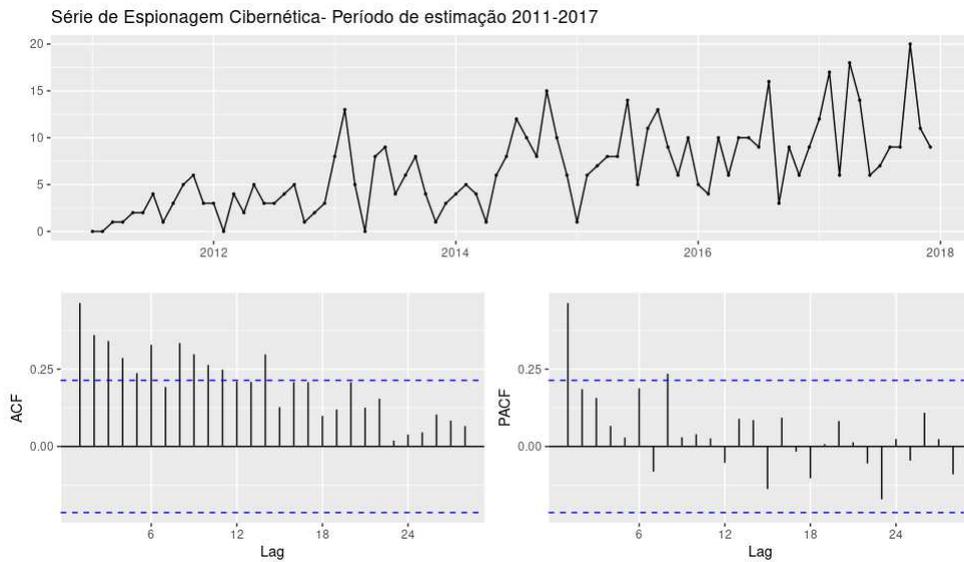


Figura 11 – Série temporal de Espionagem Cibernética e seus ACF e PACF. Fonte: Do autor.

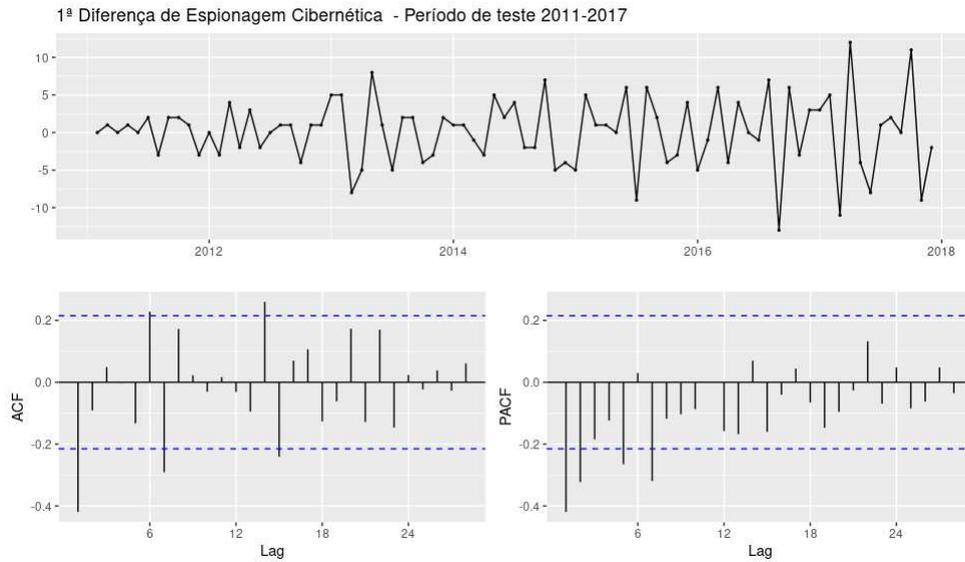


Figura 12 – 1ª Diferença da Série temporal de Espionagem Cibernética e seus ACF e PACF.  
 Fonte: Do autor.

Tabela 3 – Resultados Teste ADF para Espionagem Cibernética

	Série em Nível	1ª Diferença
$p$	0.131	0.01
Dickey-Fuller	-2.503	-6.564

O gráfico ACF da Figura 11 mostra um declínio lento dos valores de lag, o que indica a necessidade de se trabalhar com a sua diferença. Analisando os resultados dos testes ADF na Tabela 3, isso é confirmado pela série em nível não ser estacionária ( $p > 0.05$ ) e pela sua 1ª Diferença ser estacionária ( $p < 0.05$ ). Com isso, parâmetro de  $d$  do modelo é igual a 1 e, como o intervalo de confiança do ACF e PACF da 1ª Diferença da série (Figura 12) ser rompido logo no primeiro lag, foi definido o valor 1 tanto para o parâmetro  $q$  quanto para o parâmetro  $p$ , resultando num modelo ARIMA(1,1,1). Em sequência, verificou-se a estabilidade dos coeficientes do modelo com o Teste do Círculo Unitário dos Coeficientes (Figura 13) e a normalidade de seus resíduos vendo se os valores de lag no gráfico ACF dos mesmos estão em maior parte dentro do intervalo de confiança (Figura 14) e analisando os resultados do Teste Jarque-Bera e o Teste Ljung-Box (Tabela 4).

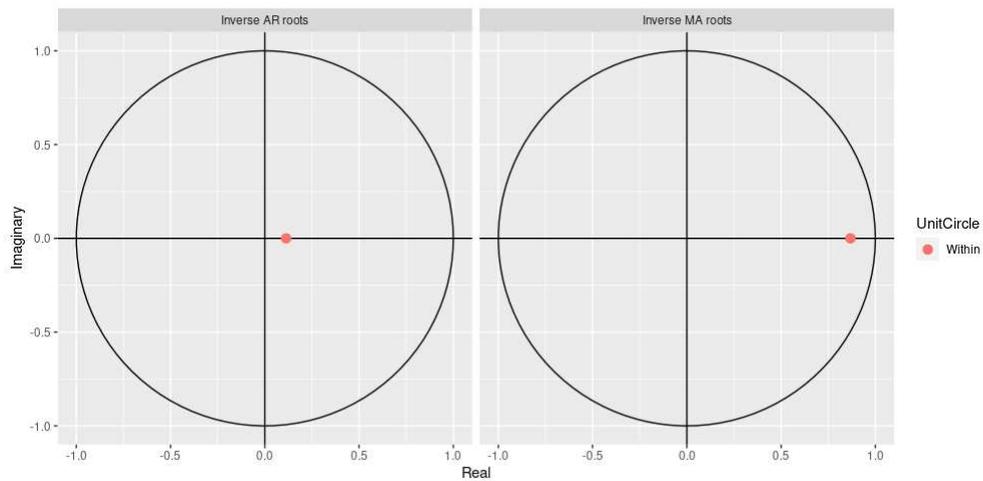


Figura 13 – Teste de Círculo Unitário para o modelo da Série de Espionagem Cibernética: ARIMA(1,1,1). Fonte: Do autor.

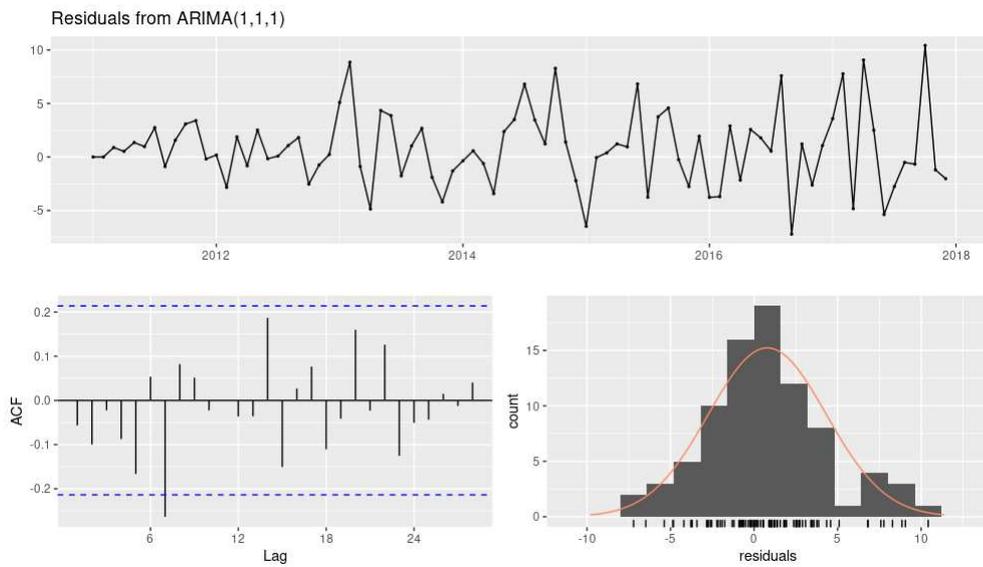


Figura 14 – Teste do Resíduo para o modelo da Série de Espionagem Cibernética: ARIMA(1,1,1). Fonte: Do autor.

Tabela 4 – Resultados Teste Ljung-Box e Jarque-Bera para Espionagem Cibernética

	ARIMA(1,1,1)
Ljung-Box	$p = 0.234$
Jarque-Bera	$p = 0.209$

No teste de Círculo Unitário dos Coeficientes da Figura 13, os valores dos coeficientes estarem dentro do círculo unitário indica que os mesmos são estáveis, confirmando que o modelo estimado é estacionário. No gráfico ACF contidos na Figura 14, observa-se que todos os valores

de lag estão dentro do intervalo de confiança e na Tabela 4 como ambos os valores  $p$  rejeitam a hipótese nula dos resíduos não serem normalmente distribuídos ( $p > 0.05$ ), pode-se assumir que o resíduo do modelo se comporta como ruído branco.

Na Figura 15, tem-se o gráfico dos valores ajustados do modelo versus os dados observados. Nota-se que por mais que os valores ajustados acompanhem a tendência da série, eles ficam como uma curva média entre os picos dos dados observados.

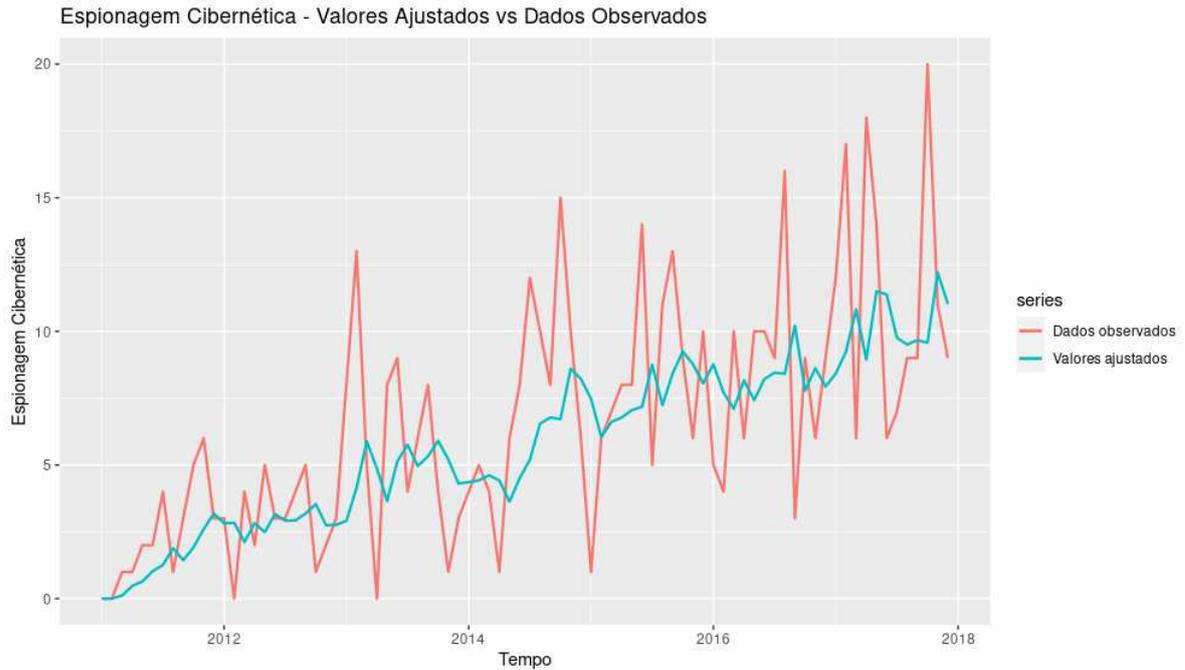


Figura 15 – Espionagem Cibernética - Valores Ajustados vs Dados Observados: ARIMA(1,1,1).  
 Fonte: Do autor.

#### 4.1.1.3 Guerra Cibernética

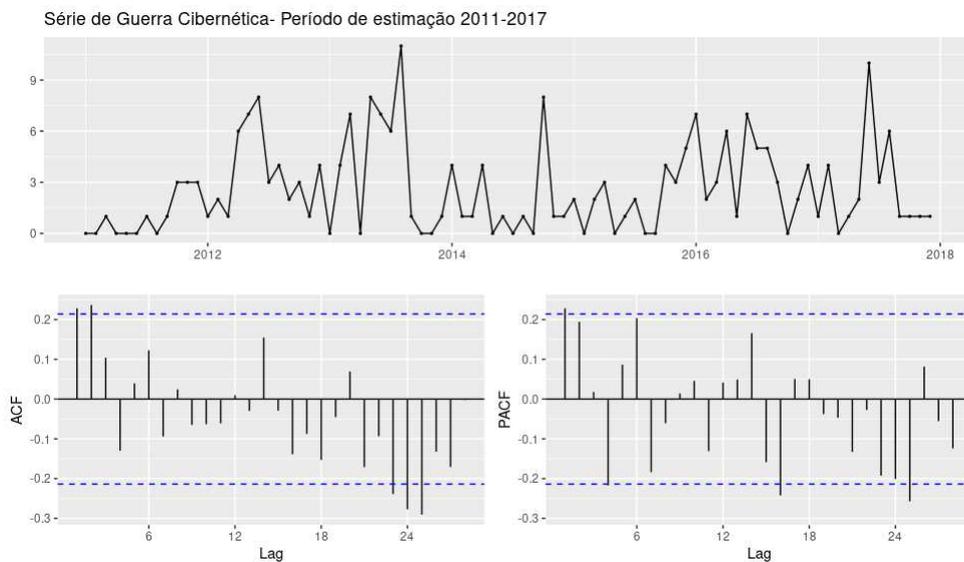


Figura 16 – Série temporal de Guerra Cibernética e seus ACF e PACF. Fonte: Do autor.

Tabela 5 – Resultados Teste ADF para Guerra Cibernética

	Série em Nível	1ª Diferença
$p$	0.01	-
Dickey-Fuller	-4.500	-

No gráfico ACF da Figura 16 já pode-se observar que os valores de lag caem rapidamente, indicando a estacionariedade da série e com o resultado do teste ADF ( $p < 0.05$ ) evidenciado na Tabela 5 foi confirmado que a série em nível já é estacionária, com isso, parâmetro de  $d$  do modelo é igual a 0. Por mais que se há a quebra dos lags do ACF e PACF logo no primeiro valor, indicando a possibilidade de adotar os valores  $p$  e  $q$  igual a 1, foi necessário ajustar esses valores até que se alcançasse a normalidade dos resíduos. Para isso, os valores de lag no gráfico ACF dos mesmos tem que estar em maior parte dentro do intervalo de confiança (Figura 18) e os resultados do Teste Jarque-Bera e o Teste Ljung-Box tem que ser maior que o limiar definido ( $p > 0.05$ ) (Tabela 6). Quando foram cumpridos esses requisitos, foi definido um ARIMA(2,0,3).

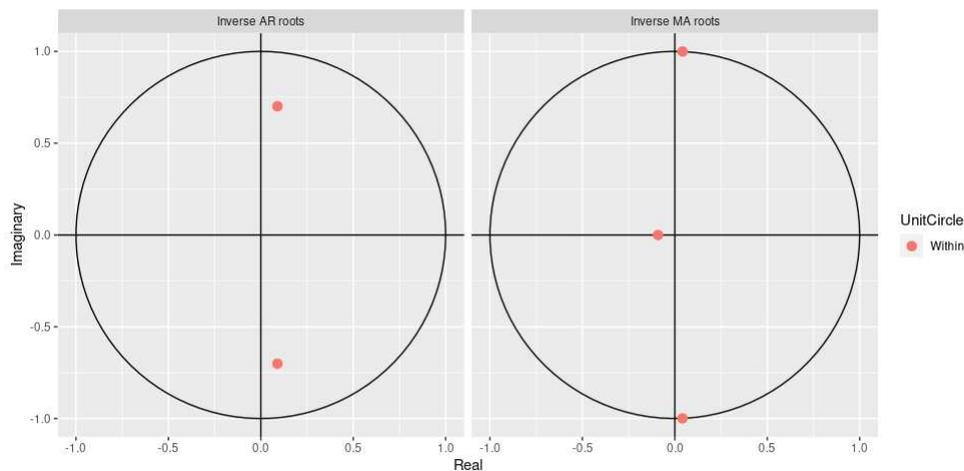


Figura 17 – Teste de Círculo Unitário para o modelo da série de Guerra Cibernética: ARIMA(2,0,3). Fonte: Do autor.

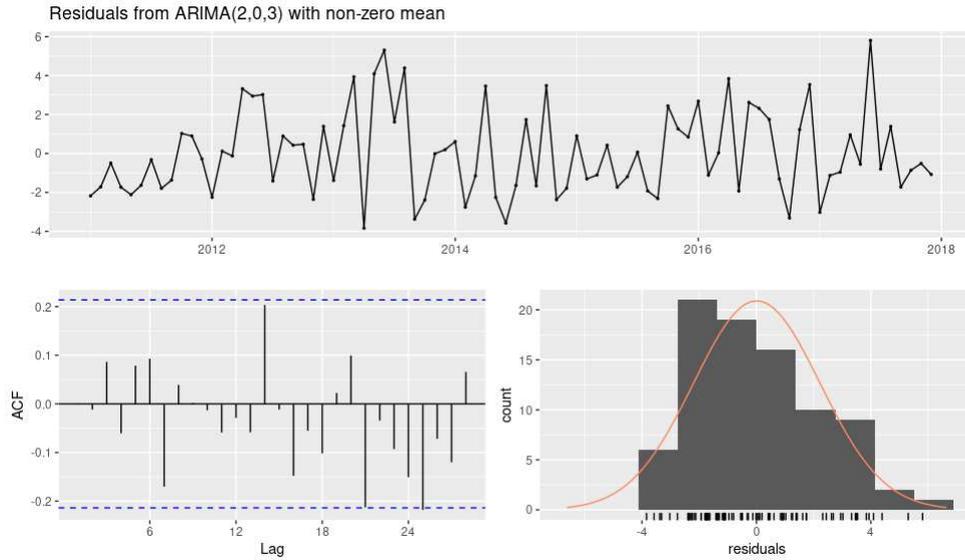


Figura 18 – Teste do Resíduo para o modelo da Série de Guerra Cibernética: ARIMA(2,0,3).  
 Fonte: Do autor.

Tabela 6 – Resultados Teste Ljung-Box e Jarque-Bera para Guerra Cibernética

	ARIMA(2,0,3)
Ljung-Box	$p = 0.296$
Jarque-Bera	$p = 0.075$

No teste de Círculo Unitário dos Coeficientes da Figura 17, os valores dos coeficientes estarem dentro do círculo unitário indica que os mesmos são estáveis, confirmando que o modelo estimado é estacionário. No gráfico ACF contidos na Figura 18, observa-se que todos os valores de lag estão dentro do intervalo de confiança e na Tabela 6 como ambos os valores  $p$  rejeitam a hipótese nula dos resíduos não serem normalmente distribuídos ( $p > 0.05$ ), pode-se assumir que o resíduo do modelo se comporta como ruído branco.

Na Figura 19, tem-se o gráfico dos valores ajustados do modelo versus os dados observados. Nota-se que os valores ajustados conseguem acompanhar algumas tendências da série, mas em vários pontos parecem errar bastante.

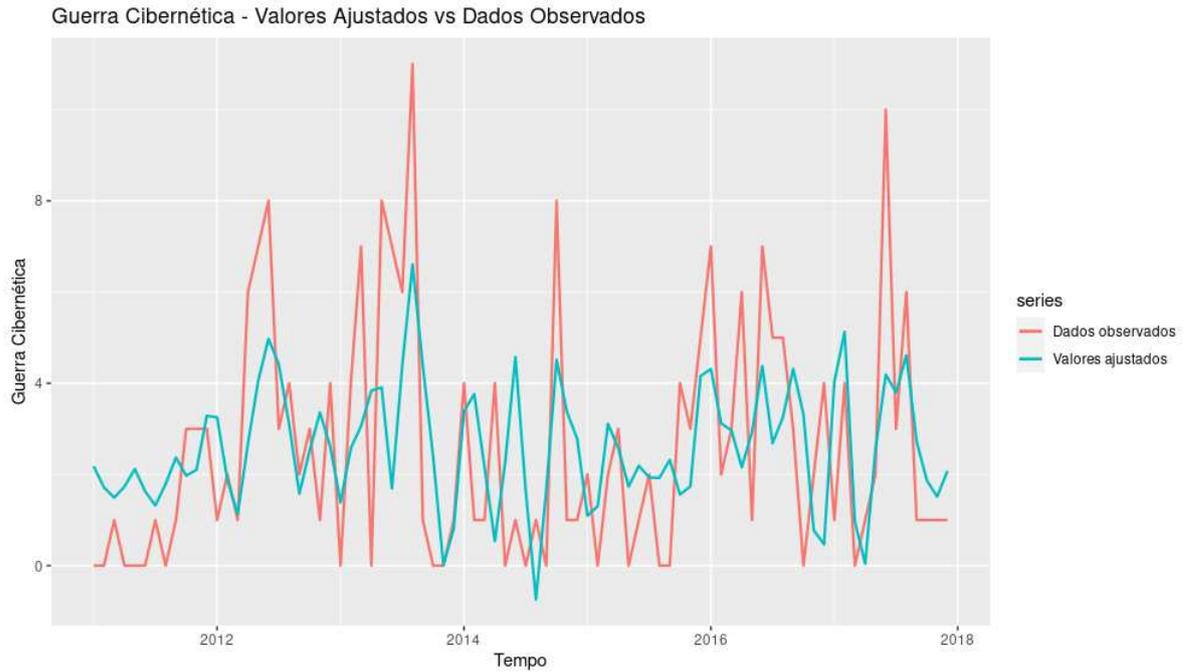


Figura 19 – Guerra Cibernética - Valores Ajustados vs Dados Observados: ARIMA(2,0,3). Fonte: Do autor.

#### 4.1.1.4 Hacktivismo

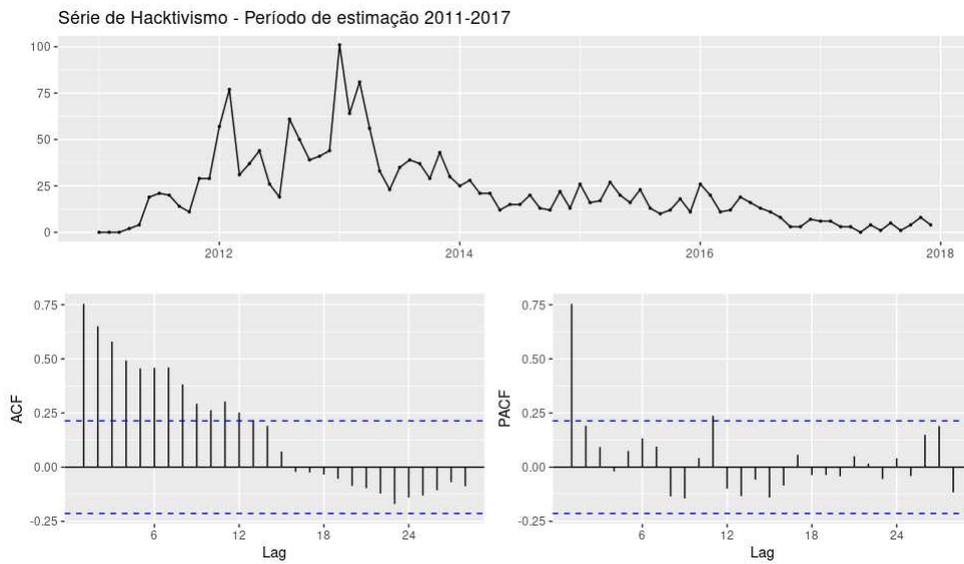


Figura 20 – Série temporal de Hacktivismo e seus ACF e PACF. Fonte: Do autor.

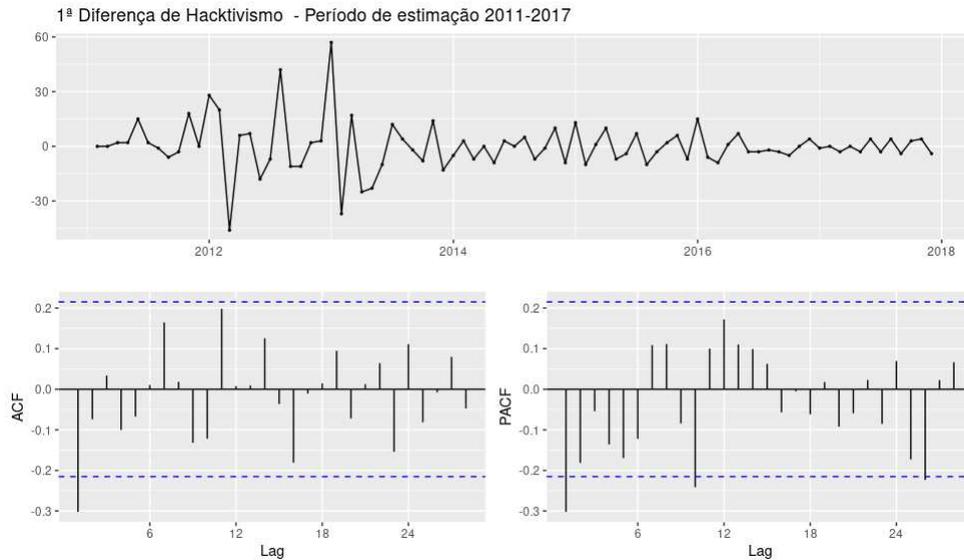


Figura 21 – 1ª Diferença da Série temporal de Hacktivismo e seus ACF e PACF. Fonte: Do autor.

Tabela 7 – Resultados Teste ADF para Hacktivismo

	Série em Nível	1ª Diferença
$p$	0.271	0.01
Dickey-Fuller	-2.122	-5.842

O gráfico ACF da Figura 20 mostra um declínio lento dos valores de lag, o que indica a necessidade de se trabalhar com a sua diferença. Analisando os resultados dos testes ADF na Tabela 7, isso é confirmado pela série em nível não ser estacionária ( $p > 0.05$ ) e pela sua 1ª Diferença ser ( $p < 0.05$ ). Com isso, parâmetro de  $d$  do modelo é igual a 1 e, como o intervalo de confiança do ACF e PACF da 1ª Diferença da série (Figura 21) ser rompido logo no primeiro lag, foi definido o valor 1 tanto para o parâmetro  $q$  quanto para o parâmetro  $p$ , resultando num modelo ARIMA(1,1,1). Em sequência, verificou-se a estabilidade dos coeficientes do modelo com o Teste do Círculo Unitário dos Coeficientes (Figura 22) e a normalidade de seus resíduos vendo se os valores de lag no gráfico ACF dos mesmos estão em maior parte dentro do intervalo de confiança (Figura 23) e analisando os resultados do Teste Jarque-Bera e o Teste Ljung-Box (Tabela 8).

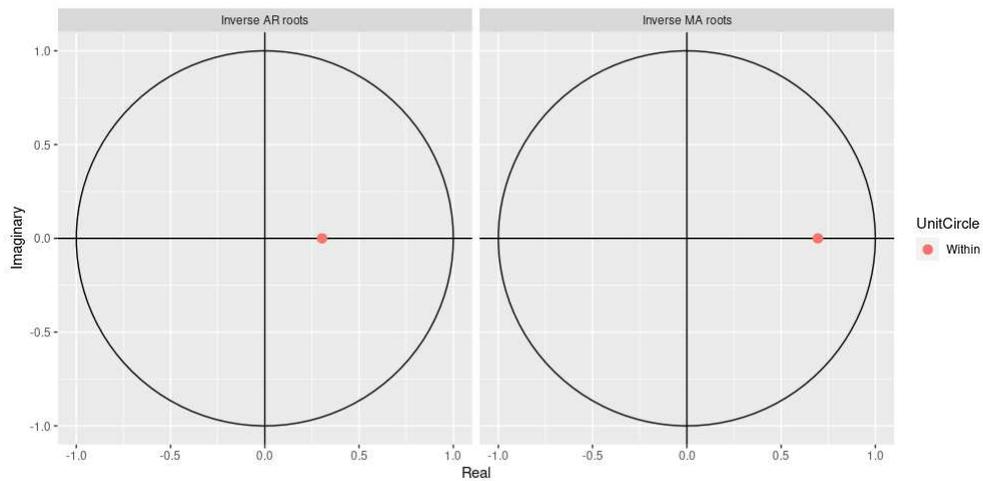


Figura 22 – Teste de Círculo Unitário para o modelo da Série de Hacktivism: ARIMA(1,1,1).  
 Fonte: Do autor.

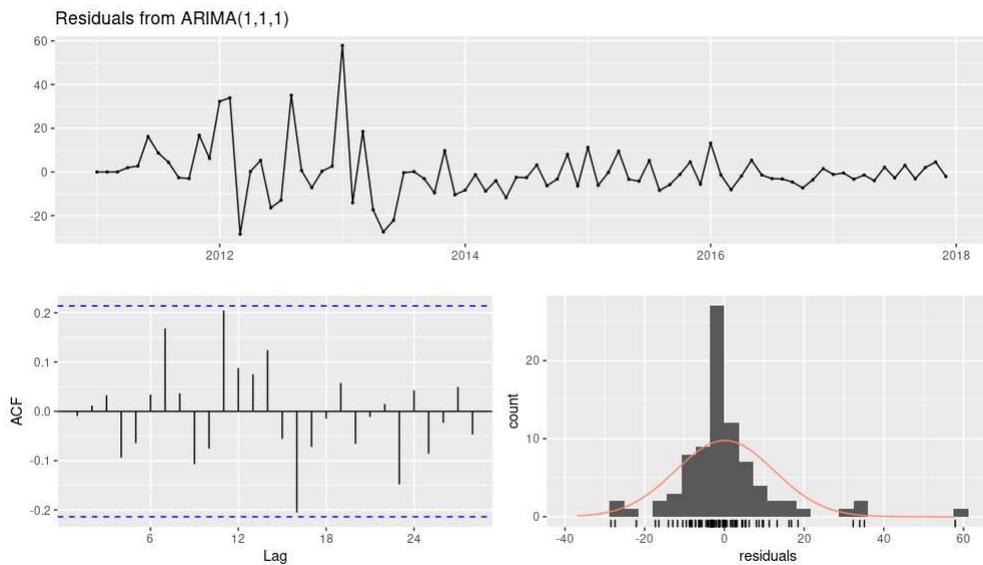


Figura 23 – Teste do Resíduo para o modelo da Série de Hacktivism: ARIMA(1,1,1). Fonte: Do autor.

Tabela 8 – Resultados Teste Ljung-Box e Jarque-Bera para Hacktivism

	ARIMA(1,1,1)
Ljung-Box	$p = 0.248$
Jarque-Bera	$p = 2.2e-16$

No teste de Círculo Unitário dos Coeficientes da Figura 22, eles se mostraram estáveis, confirmando que o modelo estimado é estacionário. Por mais que o teste Jarque-Bera que consta na Tabela 8 não ter sido validado ( $p < 0.05$ ), os valores de lag do ACF da Figura 23 estarem dentro

do intervalo de confiança e o valor de teste Ljung-Box (Tabela 8) rejeitar a hipótese nula de que os resíduos não são normalmente distribuídos ( $p > 0.05$ ), pode-se assumir que o resíduo do modelo se comporta como ruído branco.

Na Figura 24, tem-se o gráfico dos valores ajustados do modelo versus os dados observados. Nota-se que a curva de valores ajustados acompanha a curva dos dados reais de forma notável.

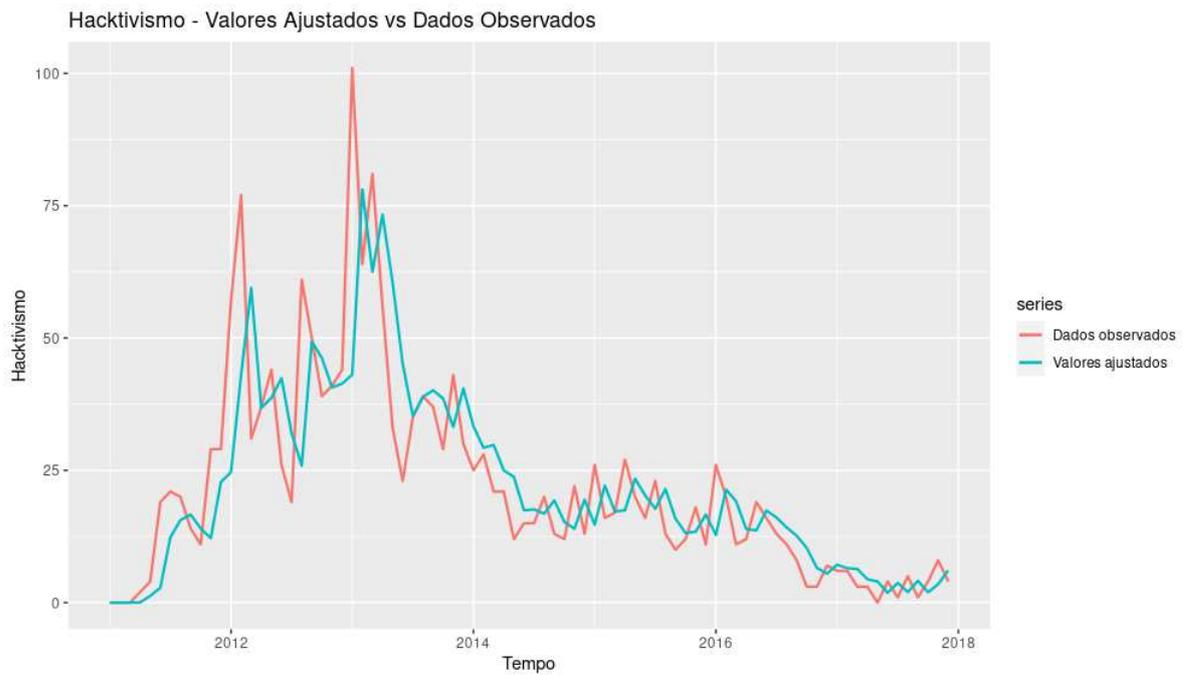


Figura 24 – Hackingismo - Valores Ajustados vs Dados Observados: ARIMA(1,1,1). Fonte: Do autor.

#### 4.1.1.5 Total de incidentes

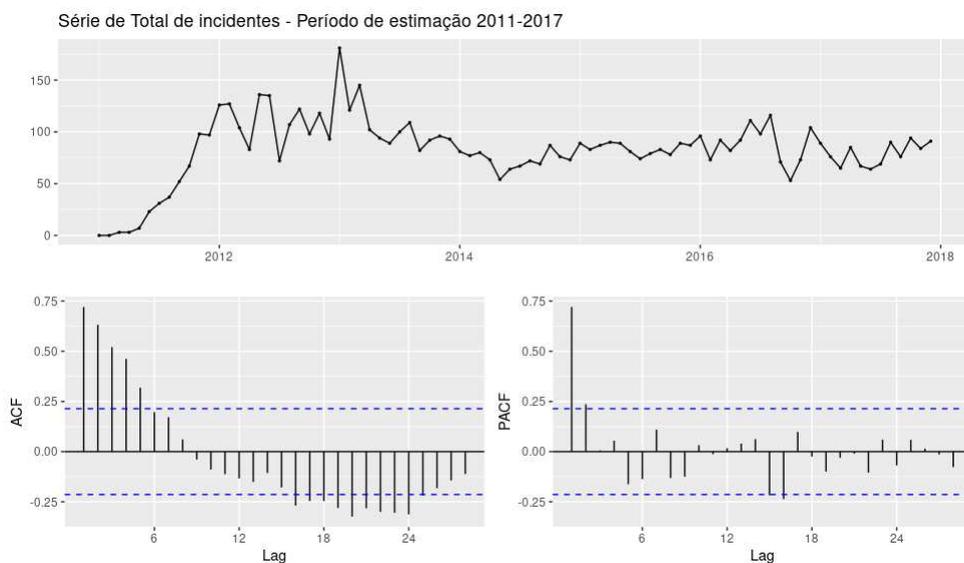


Figura 25 – Série temporal do Total de incidentes e seus ACF e PACF. Fonte: Do autor.

Tabela 9 – Resultados Teste ADF para Total de incidentes

	Série em Nível	1ª Diferença
$p$	0.015	-
Dickey-Fuller	-3.418	-

Os gráficos da Figura 25 mostram a possibilidade da série não ser estacionária, mas com o resultado do teste ADF ( $p < 0.05$ ) evidenciado na Tabela 9 foi verificado que a série em nível já é estacionária. Com isso, parâmetro de  $d$  do modelo é igual a 0 e, como o intervalo de confiança do ACF e PACF da série em nível ser rompido logo no primeiro lag, foi definido o valor 1 tanto para o parâmetro  $q$  quanto para o parâmetro  $p$ , resultando num modelo ARIMA(1,0,1). Em sequência, verificou-se a estabilidade dos coeficientes do modelo com o Teste do Círculo Unitário dos Coeficientes (Figura 26) e a normalidade de seus resíduos vendo se os valores de lag no gráfico ACF dos mesmos estão em maior parte dentro do intervalo de confiança (Figura 27) e analisando os resultados do Teste Jarque-Bera e o Teste Ljung-Box (Tabela 10).

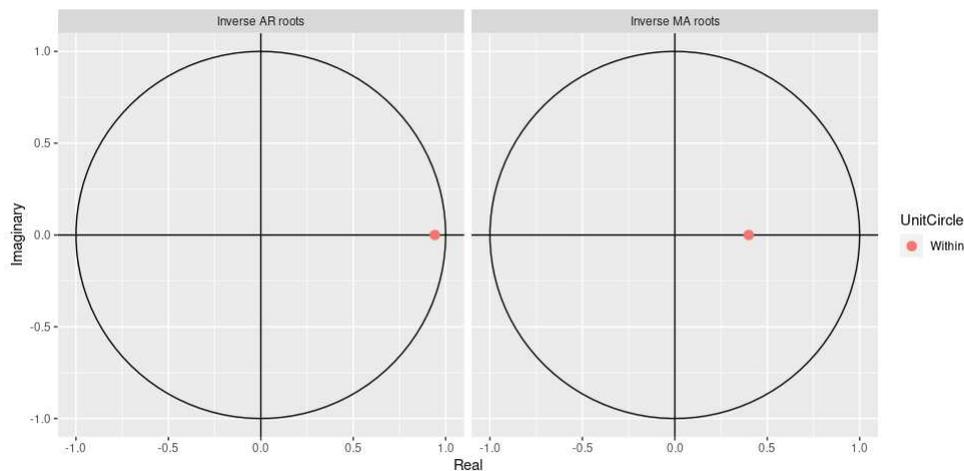


Figura 26 – Teste de Círculo Unitário para Série de Total de incidentes: ARIMA(1,0,1). Fonte: Do autor.

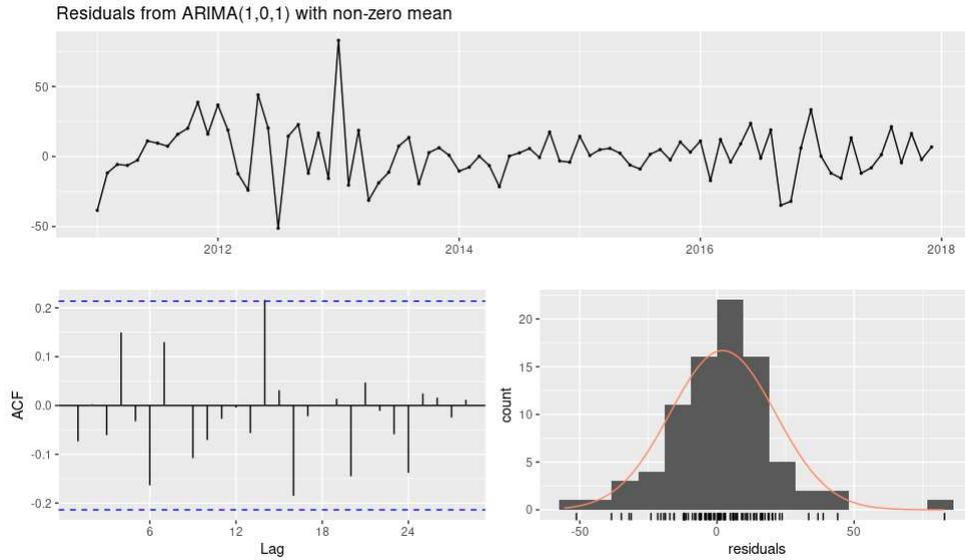


Figura 27 – Teste do Resíduo para o modelo da Série de Total de incidentes: ARIMA(1,0,1).  
 Fonte: Do autor.

Tabela 10 – Resultados Teste Ljung-Box e Jarque-Bera para Total de incidentes

	ARIMA(1,0,1)
Ljung-Box	$p = 0.225$
Jarque-Bera	$p = 3.561e-09$

No teste de Círculo Unitário dos Coeficientes da Figura 26, eles se mostraram estáveis, confirmando que o modelo estimado é estacionário. Por mais que o teste Jarque-Bera que consta na Tabela 10 não ter sido validado ( $p < 0.05$ ), os valores de lag do ACF da Figura 27 estarem dentro do intervalo de confiança e o valor de teste Ljung-Box (Tabela 10) rejeitar a hipótese nula de que os resíduos não são normalmente distribuídos ( $p > 0.05$ ), pode-se assumir que o resíduo do modelo se comporta como ruído branco.

Na Figura 28, tem-se o gráfico dos valores ajustados do modelo versus os dados observados. Nota-se que a curva do modelo acompanha as tendências da série de uma forma notável.

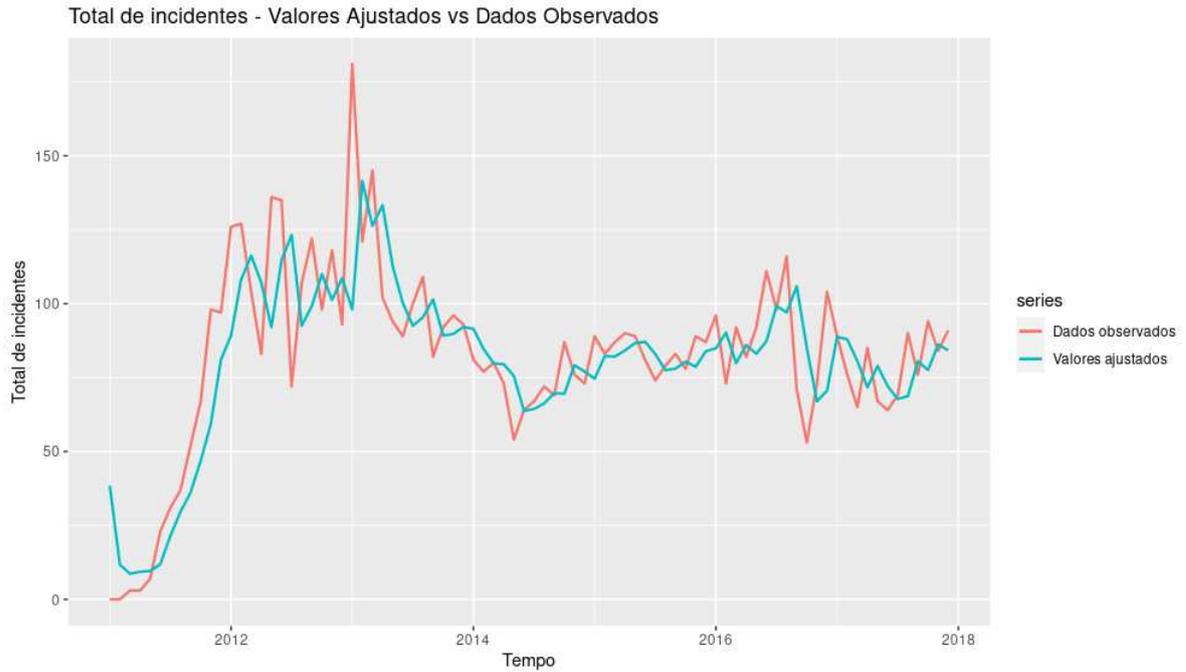


Figura 28 – Total de Incidentes - Valores Ajustados vs Dados Observados: ARIMA(1,1,1). Fonte: Do autor.

#### 4.1.2 Período de estimativa de 2014 a 2017

A seguir seguem os resultados dos passos de definição do modelo para as 4 categorias de ataques cibernéticos e para o total de incidentes usando o período de treinamento de 2014 a 2017.

##### 4.1.2.1 Crimes Cibernéticos

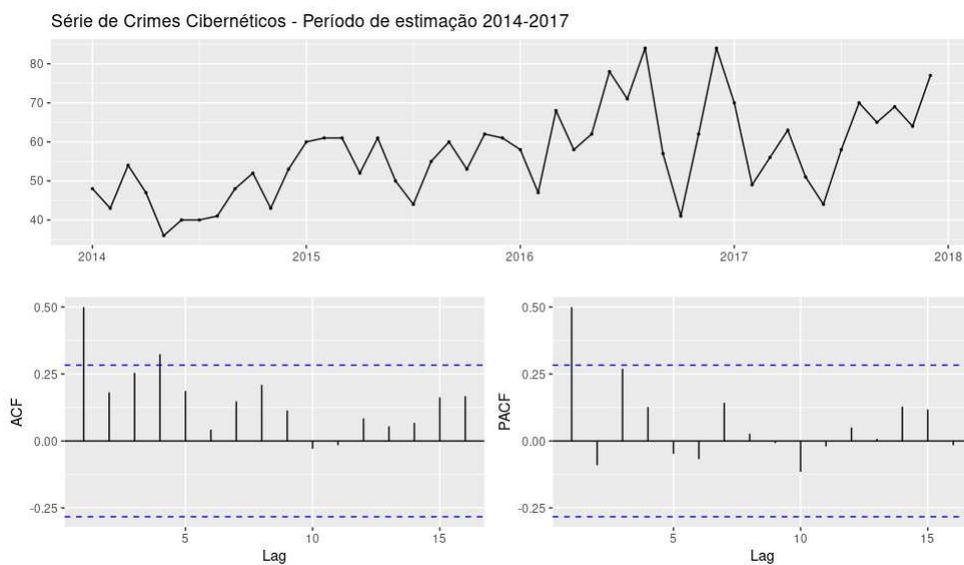


Figura 29 – Série temporal de Crimes Cibernéticos e seus ACF e PACF. Fonte: Do autor.

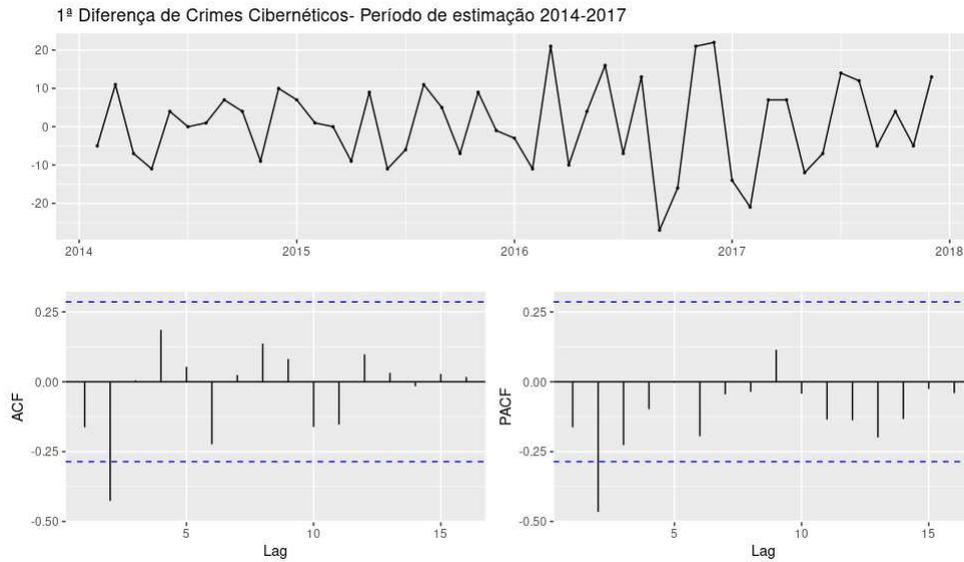


Figura 30 – 1ª Diferença da Série temporal de Crimes Cibernéticos e seus ACF e PACF. Fonte: Do autor.

Tabela 11 – Resultados Teste ADF para Crimes Cibernéticos

	Série em Nível	1ª Diferença
$p$	0.464	0.010
Dickey-Fuller	-1.600	-4.890

O gráfico ACF da Figura 29 mostra um declínio lento dos valores de lag, o que indica a necessidade de se trabalhar com a sua diferença. Analisando os resultados dos testes ADF na Tabela 11, isso é confirmado pela série em nível não ser estacionária ( $p > 0.05$ ) e pela sua 1ª Diferença ser ( $p < 0.05$ ). Com isso, parâmetro de  $d$  do modelo é igual a 1 e, como o intervalo de confiança do ACF e PACF da 1ª Diferença da série (Figura 30) ser rompido logo no primeiro lag, foi definido o valor 1 tanto para o parâmetro  $q$  quanto para o parâmetro  $p$ , resultando num modelo ARIMA(1,1,1). Em sequência, verificou-se a estabilidade dos coeficientes do modelo com o Teste do Círculo Unitário dos Coeficientes (Figura 31) e a normalidade de seus resíduos vendo se os valores de lag no gráfico ACF dos mesmos estão em maior parte dentro do intervalo de confiança (Figura 32) e analisando os resultados do Teste Jarque-Bera e o Teste Ljung-Box (Tabela 12).

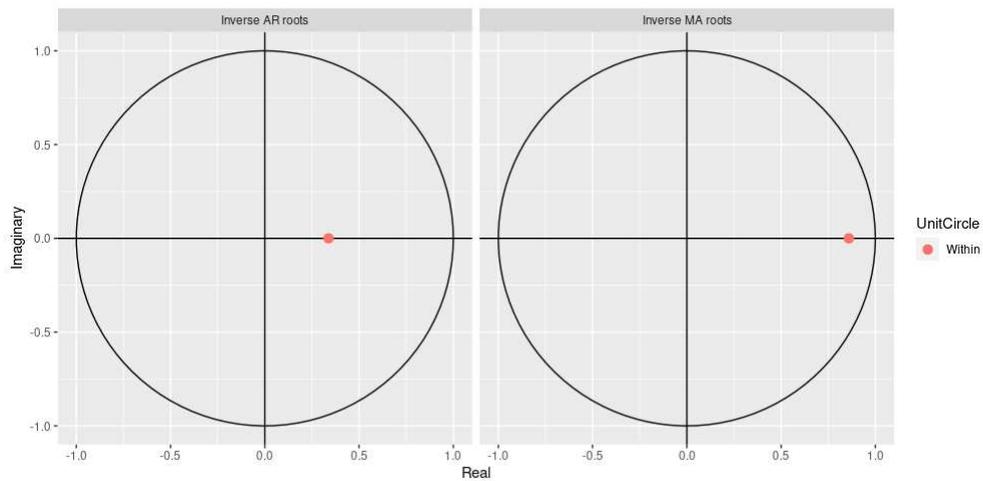


Figura 31 – Teste de Círculo Unitário para Série de Crimes Cibernéticos: ARIMA(1,1,1). Fonte: Do autor.

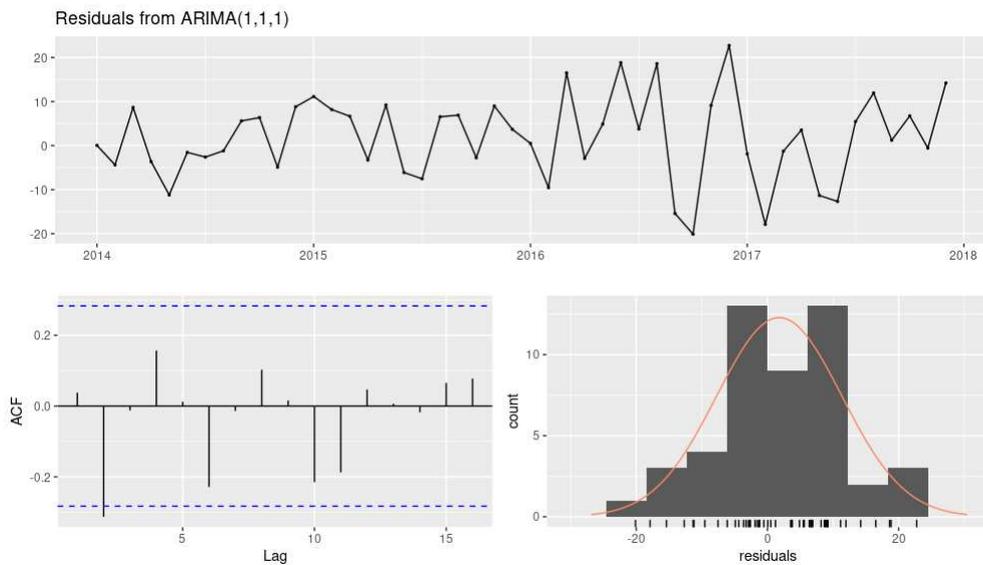


Figura 32 – Teste do Resíduo para o modelo estimado para Série: ARIMA(1,1,1). Fonte: Do autor.

Tabela 12 – Resultados Teste Ljung-Box e Jarque-Bera para Crimes Cibernéticos

	ARIMA(1,1,1)
Ljung-Box	$p = 0.109$
Jarque-Bera	$p = 0.877$

No teste de Círculo Unitário dos Coeficientes da Figura 31, os valores dos coeficientes estarem dentro do círculo unitário indica que os mesmos são estáveis, confirmando que o modelo estimado é estacionário. No gráfico ACF contidos na Figura 32, observa-se que todos os valores

de lag estão dentro do intervalo de confiança e na Tabela 12 como ambos os valores  $p$  rejeitam a hipótese nula dos testes de os resíduos não serem normalmente distribuídos ( $p > 0.05$ ), pode-se assumir que o resíduo do modelo se comporta como ruído branco.

Na Figura 33, tem-se o gráfico dos valores ajustados do modelo versus os dados observados, por mais que os valores ajustados parecem acompanhar uma tendência, eles ainda ficam um pouco distantes dos dados observados.

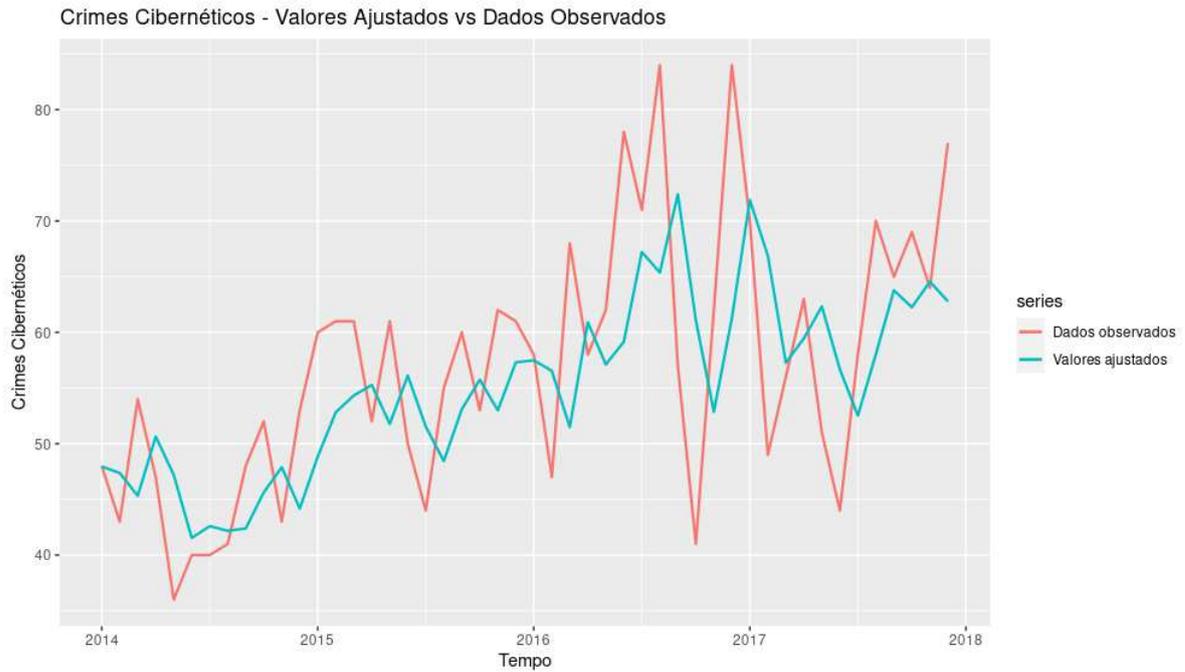


Figura 33 – Crimes Cibernéticos - Valores Ajustados vs Dados Observados: ARIMA(1,1,1).  
Fonte: Do autor.

#### 4.1.2.2 Espionagem Cibernética

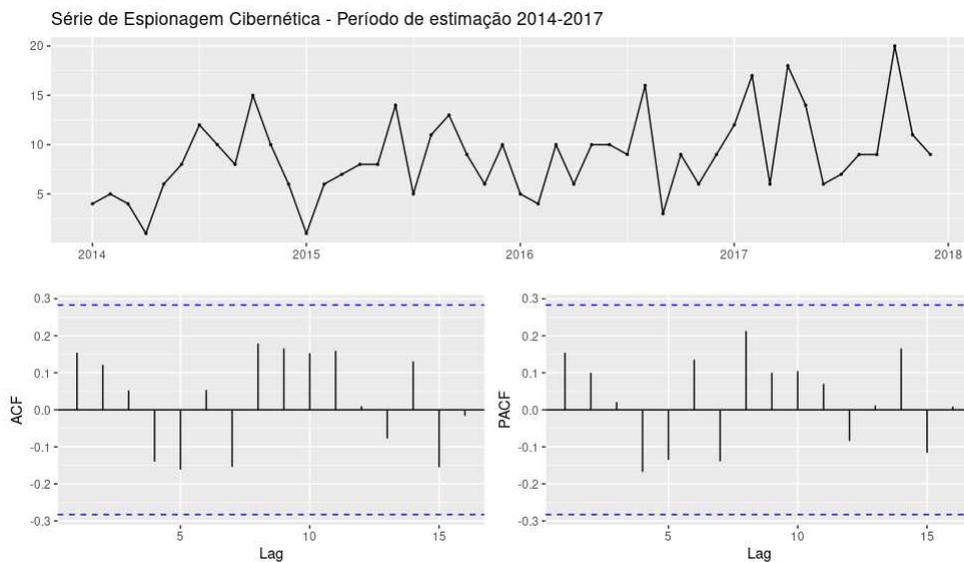


Figura 34 – Série temporal de Espionagem Cibernética e seus ACF e PACF. Fonte: Do autor.

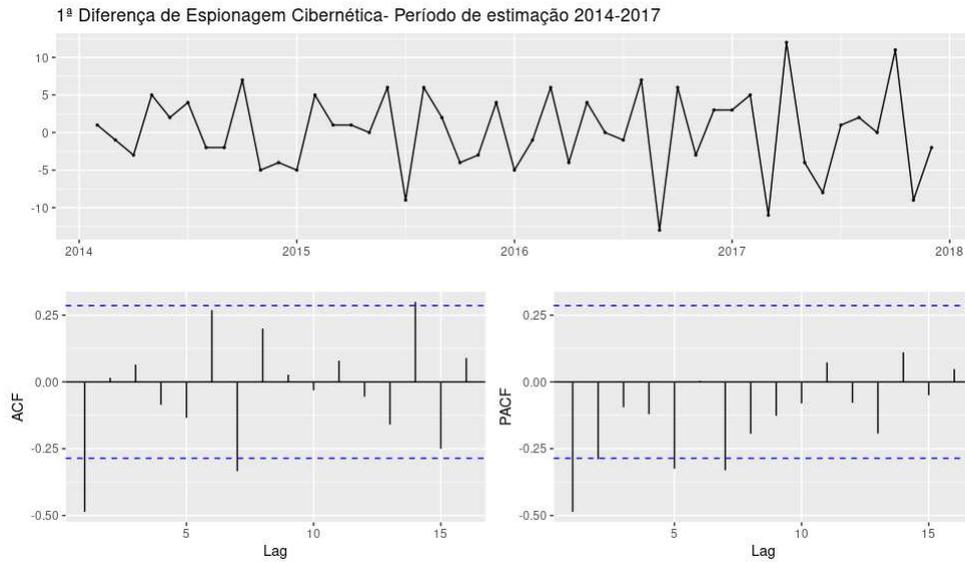


Figura 35 – 1ª Diferença da Série temporal de Espionagem Cibernética e seus ACF e PACF. Fonte: Do autor.

Tabela 13 – Resultados Teste ADF para Espionagem Cibernética

	Série em Nível	1ª Diferença
$p$	0.01	0.01
Dickey-Fuller	-3.835	-4.378

Por mais que não há uma quebra do intervalo de confiança no gráfico ACF e PACF das Figura 34 e Figura 35 e, analisando os resultados dos testes ADF na Tabela 13, ambas serem estacionárias ( $p < 0.05$ ), foi escolhido se trabalhar com a diferenciação da série, pois os modelos candidatos sem diferenciação não alcançaram a normalidade dos resíduos. Com isso, parâmetro de  $d$  do modelo é igual a 1 e, como o intervalo de confiança do ACF e PACF da 1ª Diferença da série (Figura 35) ser rompido logo no primeiro lag, foi definido o valor 1 tanto para o parâmetro  $q$  quanto para o parâmetro  $p$ , resultando num modelo ARIMA(1,1,1). Em sequência, verificou-se a estabilidade dos coeficientes do modelo com o Teste do Círculo Unitário dos Coeficientes (Figura 36) e a normalidade de seus resíduos vendo se os valores de lag no gráfico ACF dos mesmos estão em maior parte dentro do intervalo de confiança (Figura 37) e analisando os resultados do Teste Jarque-Bera e o Teste Ljung-Box (Tabela 14).

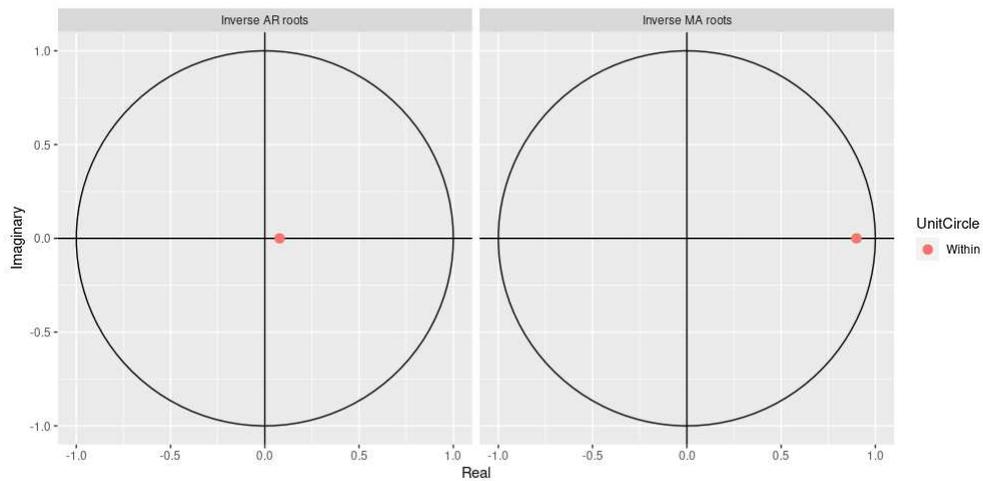


Figura 36 – Teste de Círculo Unitário para o modelo da Série de Espionagem Cibernética: ARIMA(1,1,1). Fonte: Do autor.

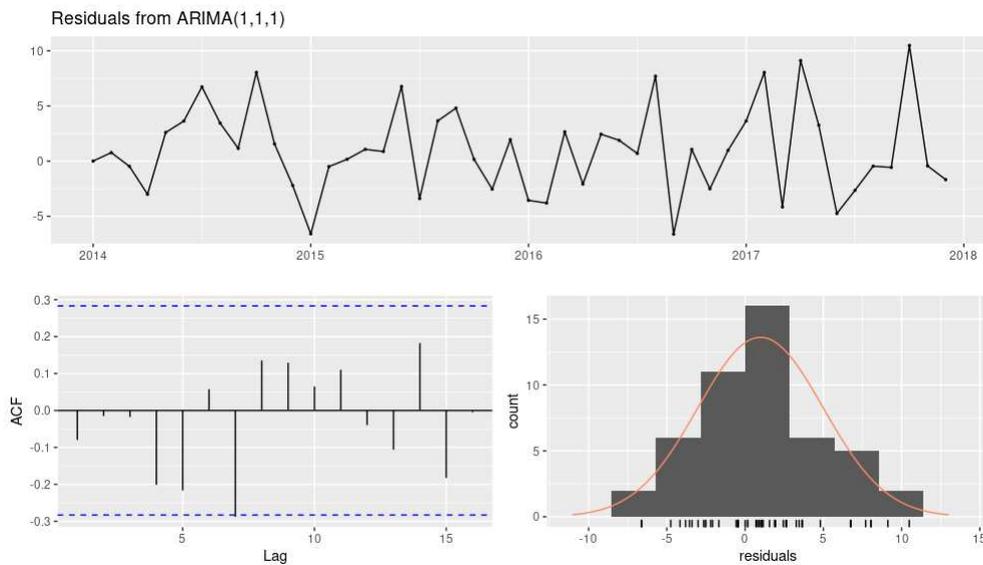


Figura 37 – Teste do Resíduo para o modelo da Série de Espionagem Cibernética: ARIMA(1,1,1). Fonte: Do autor.

Tabela 14 – Resultados Teste Ljung-Box e Jarque-Bera para Espionagem Cibernética

	ARIMA(1,1,1)
Ljung-Box	$p = 0.126$
Jarque-Bera	$p = 0.527$

No teste de Círculo Unitário dos Coeficientes da Figura 36, os valores dos coeficientes estarem dentro do círculo unitário indica que os mesmos são estáveis, confirmando que o modelo estimado é estacionário. No gráfico ACF contidos na Figura 37, observa-se que todos os valores

de lag estão dentro do intervalo de confiança e na Tabela 14 como ambos os valores p rejeitam a hipótese nula dos testes de os resíduos não serem normalmente distribuídos ( $p > 0.05$ ), pode-se assumir que o resíduo do modelo se comporta como ruído branco

Na Figura 38, tem-se o gráfico dos valores ajustados do modelo versus os dados observados. Nota-se que por mais que os valores ajustados acompanhem a tendência da série, eles ficam como uma curva média entre os picos dos dados observados.

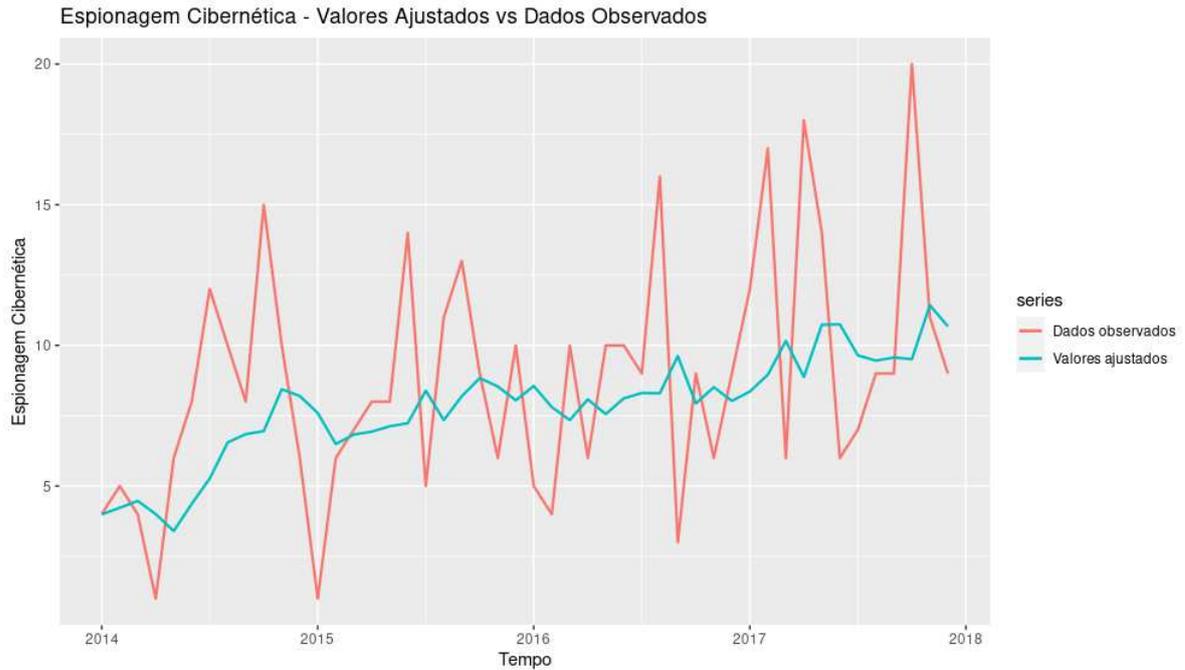


Figura 38 – Espionagem Cibernética - Valores Ajustados vs Dados Observados: ARIMA(1,1,1).  
Fonte: Do autor.

#### 4.1.2.3 Guerra Cibernética

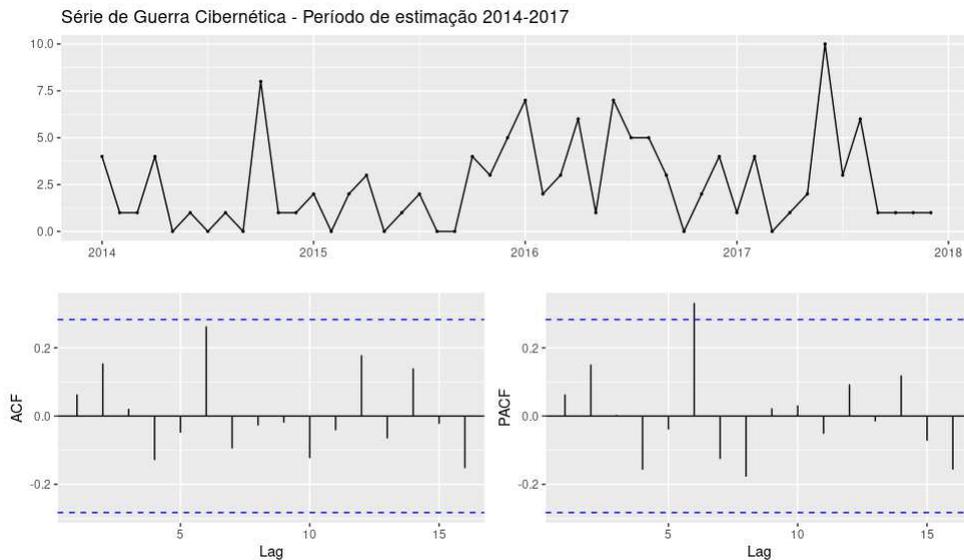


Figura 39 – Série temporal de Guerra Cibernética e seus ACF e PACF. Fonte: Do autor.

Tabela 15 – Resultados Teste ADF para Guerra Cibernética

	Série em Nível	1ª Diferença
$p$	0.026	-
Dickey-Fuller	-3.223	-

No gráfico ACF da Figura 39 já pode-se observar que os valores de lag já estão dentro do intervalo de confiança e que não há uma tendência visível, indicando a estacionariedade da série e com o resultado do teste ADF ( $p < 0.05$ ) evidenciado na Tabela 15 isso foi confirmado, com isso, parâmetro de  $d$  do modelo é igual a 0. Como os lags ACF e PACF da Figura 39 não indicam valores claros para  $p$  e  $q$ , foi necessário ajustar os mesmos até que se alcançasse a normalidade dos resíduos. Para isso, os valores de lag no gráfico ACF dos mesmos tem que estar em maior parte dentro do intervalo de confiança (Figura 41) e os resultados do Teste Jarque-Bera e o Teste Ljung-Box tem que ser maior que o limiar definido ( $p > 0.05$ ) (Tabela 16). Quando foram cumpridos esses requisitos, foi definido um ARIMA(2,0,3).

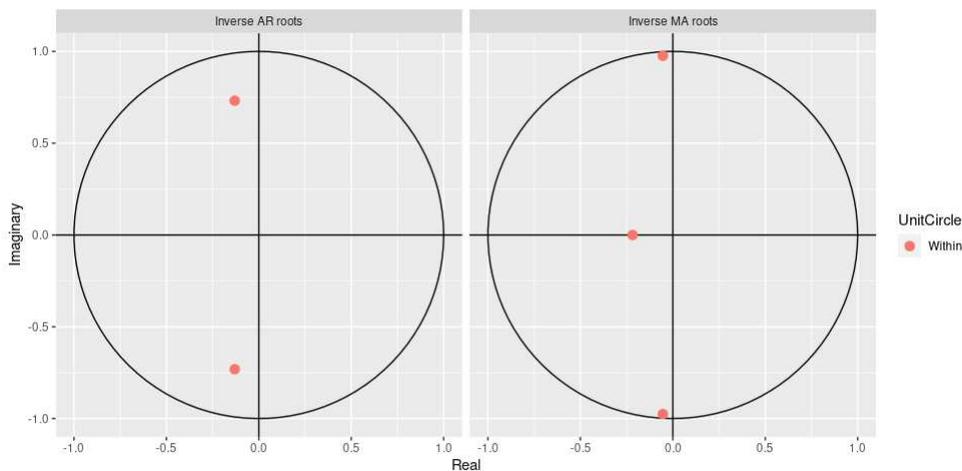


Figura 40 – Teste de Círculo Unitário para o modelo da série de Guerra Cibernética: ARIMA(2,0,3). Fonte: Do autor.

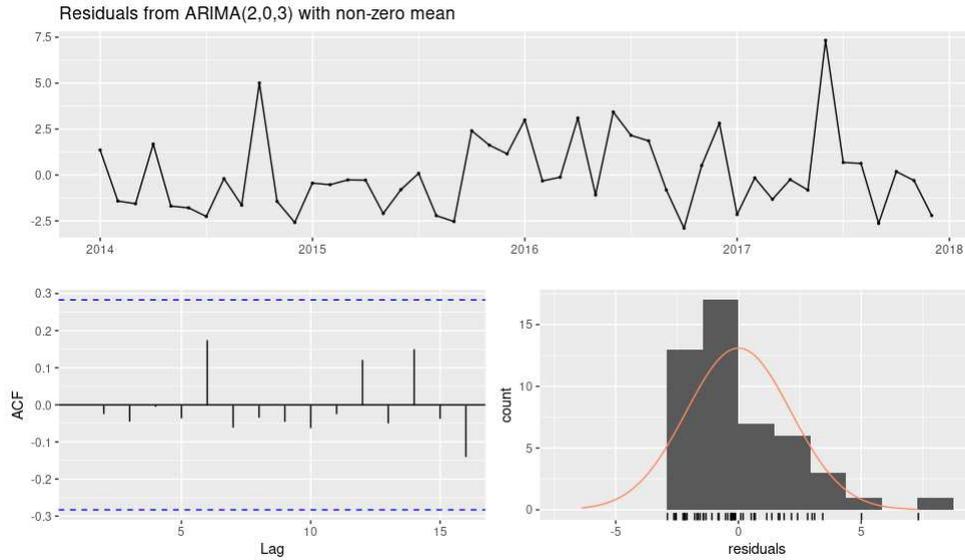


Figura 41 – Teste do Resíduo para o modelo da Série de Guerra Cibernética: ARIMA(2,0,3).  
 Fonte: Do autor.

Tabela 16 – Resultados Teste Ljung-Box e Jarque-Bera para Guerra Cibernética

	ARIMA(2,0,3)
Ljung-Box	$p = 0.623$
Jarque-Bera	$p = 4.806e-4$

No teste de Círculo Unitário dos Coeficientes da Figura 40, eles se mostraram estáveis, confirmando que o modelo estimado é estacionário. Por mais que o teste Jarque-Bera que consta na Tabela 16 não ter sido validado ( $p < 0.05$ ), os valores de lag do ACF da Figura 41 estarem dentro do intervalo de confiança e o valor de teste Ljung-Box (Tabela 16) rejeitar a hipótese nula de que os resíduos não são normalmente distribuídos ( $p > 0.05$ ), pode-se assumir que o resíduo do modelo se comporta como ruído branco.

Na Figura 42, tem-se o gráfico dos valores ajustados do modelo versus os dados observados. Percebe-se que os valores ajustados parecem acompanhar uma tendência, mas ainda ficam distantes dos dados observados.

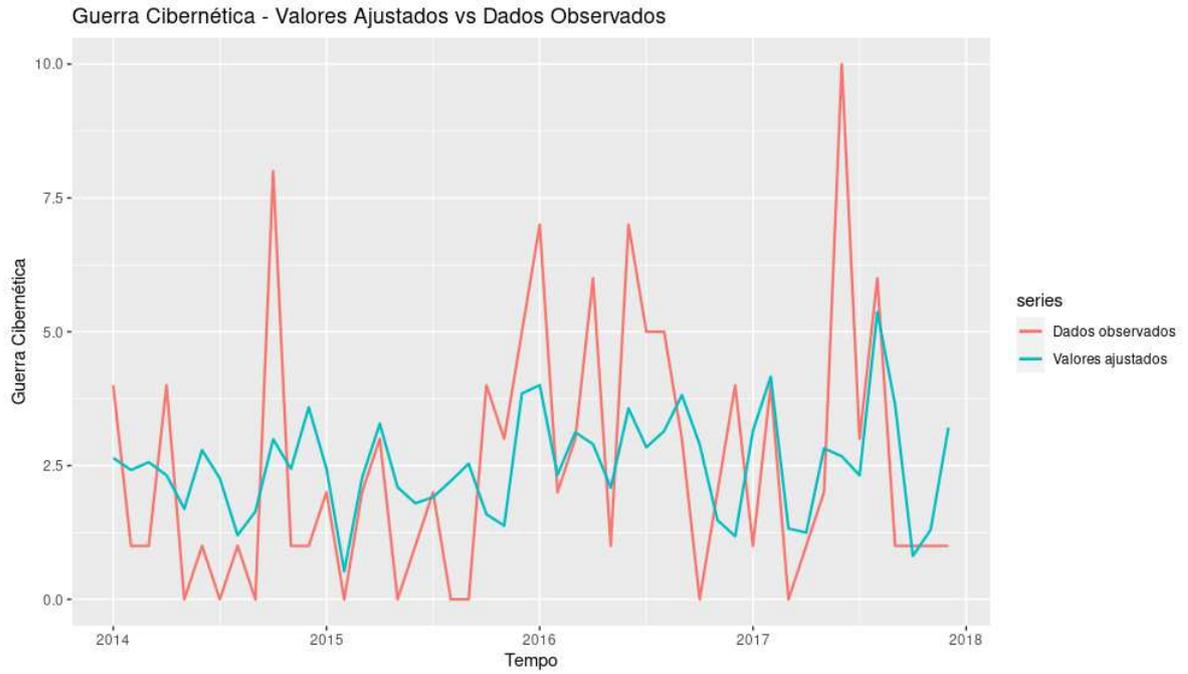


Figura 42 – Guerra Cibernética - Valores Ajustados vs Dados Observados: ARIMA(2,0,3). Fonte: Do autor.

#### 4.1.2.4 Hacking

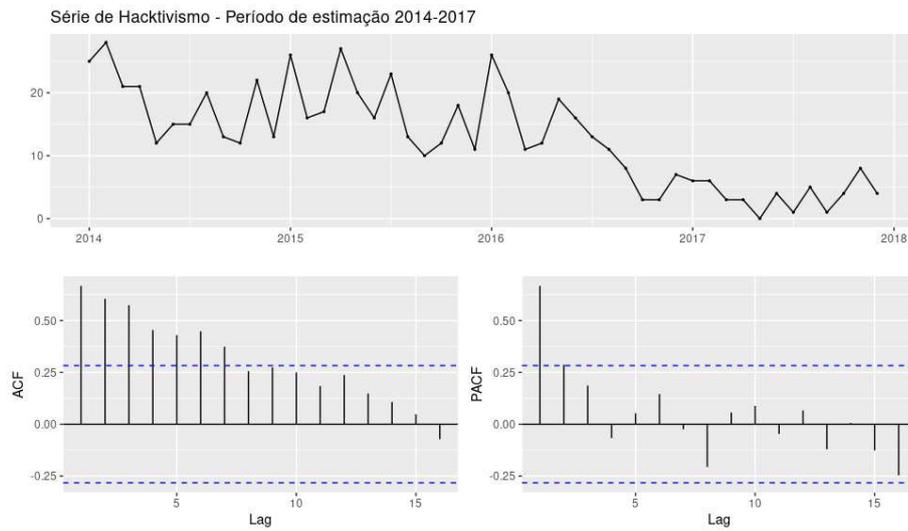


Figura 43 – Série temporal de Hacking e seus ACF e PACF. Fonte: Do autor.

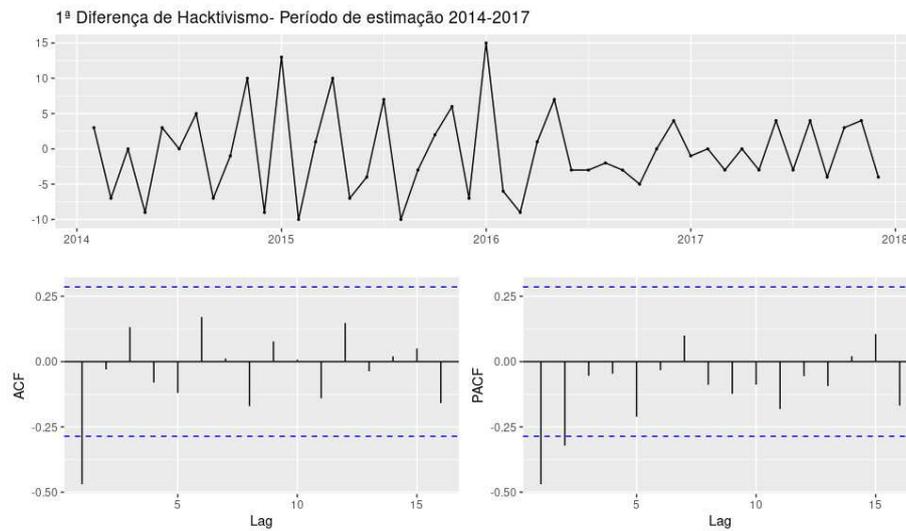


Figura 44 – 1ª Diferença da Série temporal de Hactivismo e seus ACF e PACF. Fonte: Do autor.

Tabela 17 – Resultados Teste ADF para Hactivismo

	Série em Nível	1ª Diferença
$p$	0.514	0.01
Dickey-Fuller	-1.461	-4.391

O gráfico ACF da Figura 43 mostra um declínio lento dos valores de lag, o que indica a necessidade de se trabalhar com a sua diferença. Analisando os resultados dos testes ADF na Tabela 17, isso é confirmado pela série em nível não ser estacionária ( $p > 0.05$ ) e pela sua 1ª Diferença ser ( $p < 0.05$ ). Com isso, parâmetro de  $d$  do modelo é igual a 1 e, como o intervalo de confiança do ACF e PACF da 1ª Diferença da série (Figura 44) ser rompido logo no primeiro lag, foi definido o valor 1 tanto para o parâmetro  $q$  quanto para o parâmetro  $p$ , resultando num modelo ARIMA(1,1,1). Em sequência, verificou-se a estabilidade dos coeficientes do modelo com o Teste do Círculo Unitário dos Coeficientes (Figura 45) e a normalidade de seus resíduos vendo se os valores de lag no gráfico ACF dos mesmos estão em maior parte dentro do intervalo de confiança (Figura 46) e analisando os resultados do Teste Jarque-Bera e o Teste Ljung-Box (Tabela 18).

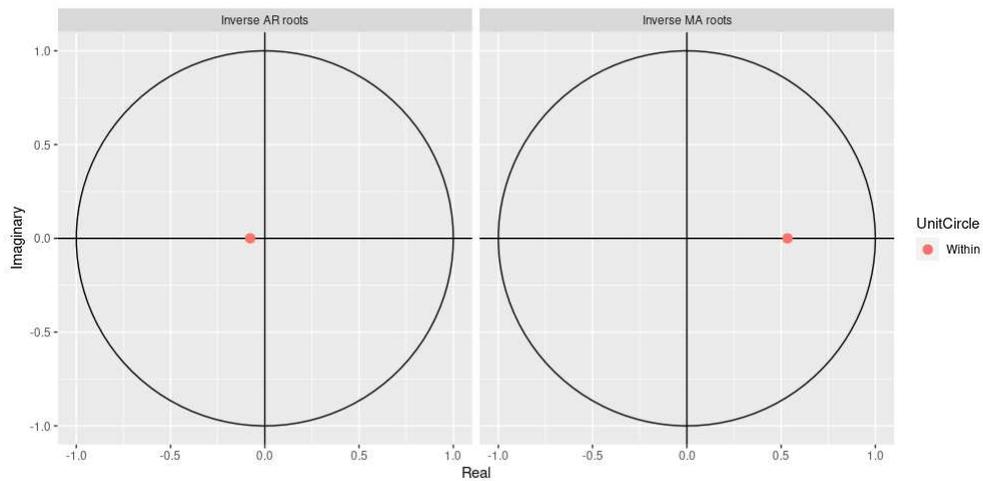


Figura 45 – Teste de Círculo Unitário para o modelo da Série de Hacktivism: ARIMA(1,1,1).  
 Fonte: Do autor.

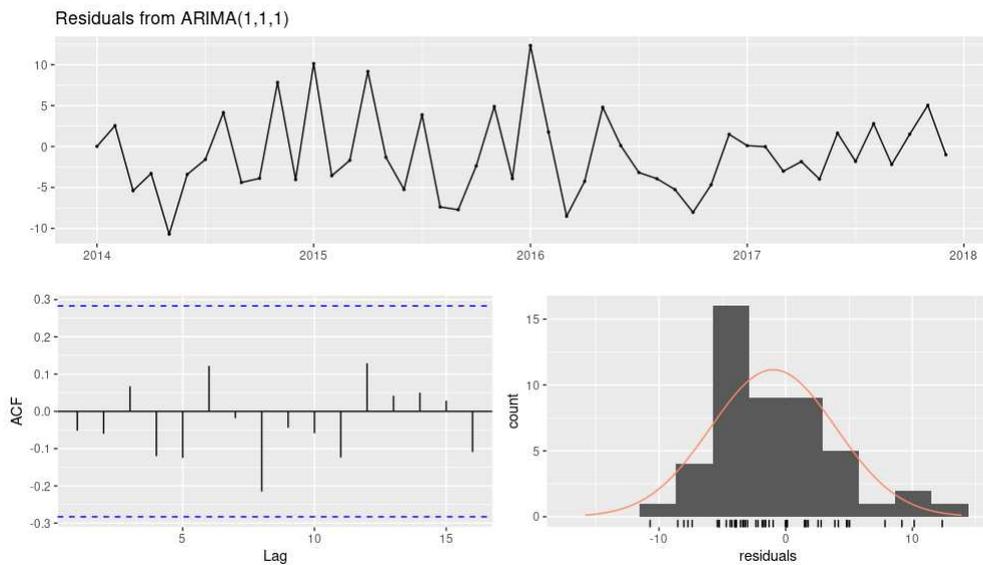


Figura 46 – Teste do Resíduo para o modelo da Série de Hacktivism: ARIMA(1,1,1). Fonte: Do autor.

Tabela 18 – Resultados Teste Ljung-Box e Jarque-Bera para Hacktivism

	ARIMA(1,1,1)
Ljung-Box	$p = 0.624$
Jarque-Bera	$p = 0.168$

No teste de Círculo Unitário dos Coeficientes da Figura 45, os valores dos coeficientes estarem dentro do círculo unitário indica que os mesmos são estáveis, confirmando que o modelo estimado é estacionário. No gráfico ACF contidos na Figura 46, observa-se que todos os valores

de lag estão dentro do intervalo de confiança e na Tabela 18 como ambos os valores  $p$  rejeitam a hipótese nula dos testes de os resíduos não serem normalmente distribuídos ( $p > 0.05$ ), pode-se assumir que o resíduo do modelo se comporta como ruído branco.

Na Figura 47, tem-se o gráfico dos valores ajustados do modelo versus os dados observados. Observa-se que os valores ajustados parecem acompanhar uma tendência, mas ainda ficam um pouco distantes dos dados picos observados.

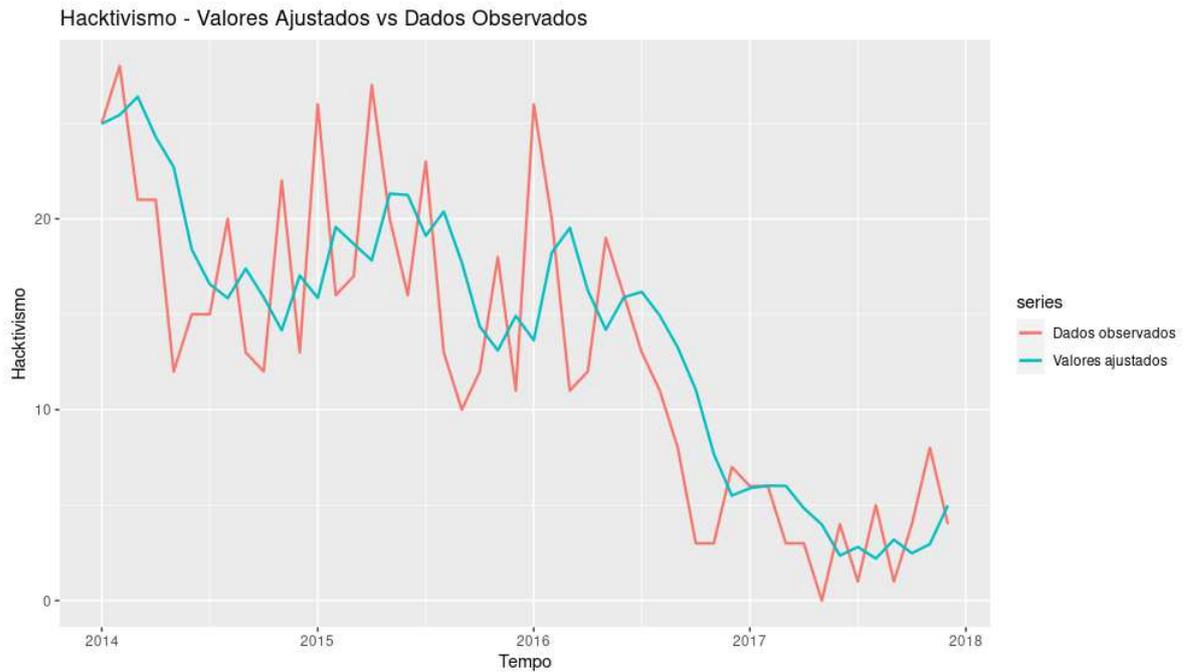


Figura 47 – Hacktivismo - Valores Ajustados vs Dados Observados: ARIMA(1,1,1). Fonte: Do autor.

#### 4.1.2.5 Total de incidentes

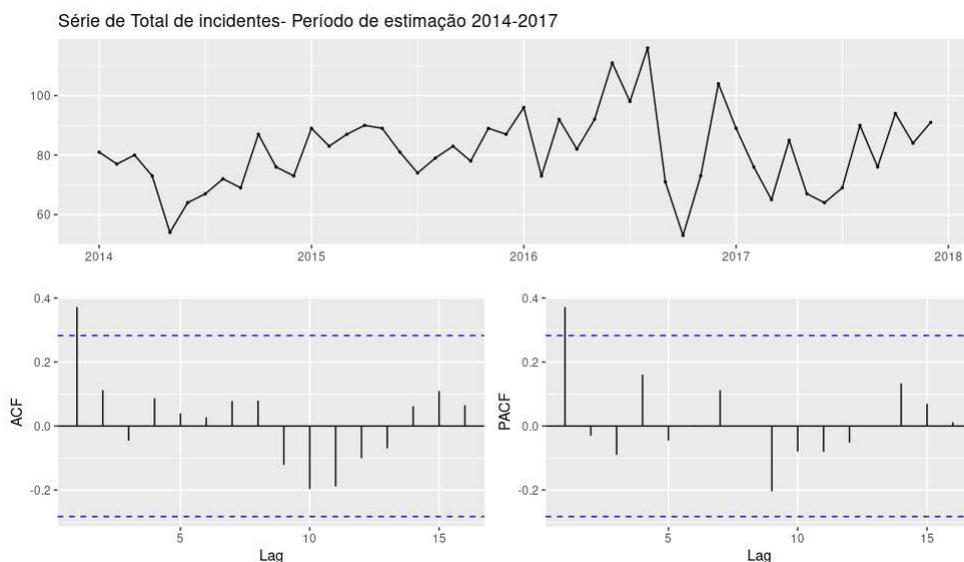


Figura 48 – Série temporal do Total de incidentes e seus ACF e PACF. Fonte: Do autor.

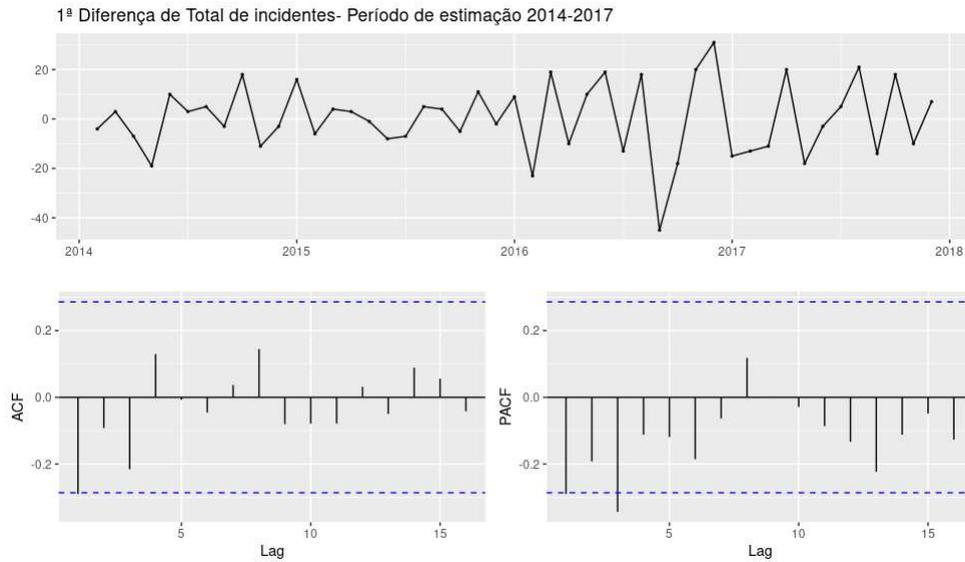


Figura 49 – 1ª Diferença da Série temporal de Total de incidentes e seus ACF e PACF. Fonte: Do autor.

Tabela 19 – Resultados Teste ADF para Total de incidentes

	Série em Nível	1ª Diferença
$p$	0.134	0.01
Dickey-Fuller	-2.511	-5.172

O gráfico ACF da Figura 48 indica a necessidade de se trabalhar com a sua diferença. Analisando os resultados dos testes ADF na Tabela 19, isso é confirmado pela série em nível não ser estacionária ( $p > 0.05$ ) e pela sua 1ª Diferença ser ( $p < 0.05$ ), com isso, parâmetro de  $d$  do modelo é igual a 1. Por mais que se há a quebra dos lags 1 do ACF e 3 do PACF da Figura 49 logo no primeiro valor, a normalidade dos resíduos foi se alcançada com num modelo ARIMA(1,1,2). Em sequência, verificou-se a estabilidade dos coeficientes do modelo com o Teste do Círculo Unitário dos Coeficientes (Figura 50) e a normalidade de seus resíduos vendo se os valores de lag no gráfico ACF dos mesmos estão em maior parte dentro do intervalo de confiança (Figura 51) e analisando os resultados do Teste Jarque-Bera e o Teste Ljung-Box (Tabela 20).

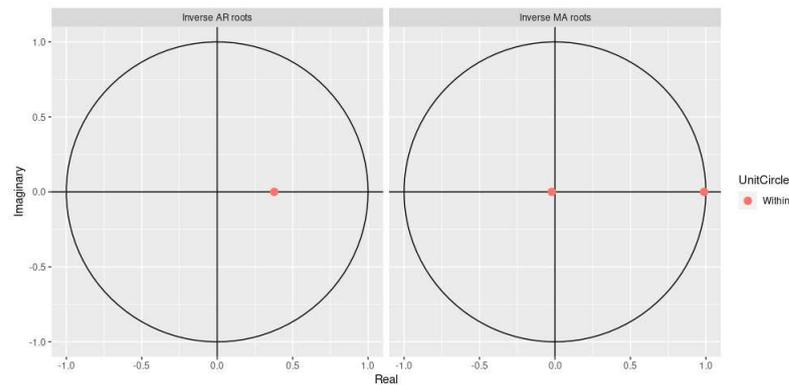


Figura 50 – Teste de Círculo Unitário para Série de Total de incidentes: ARIMA(1,1,2). Fonte: Do autor.

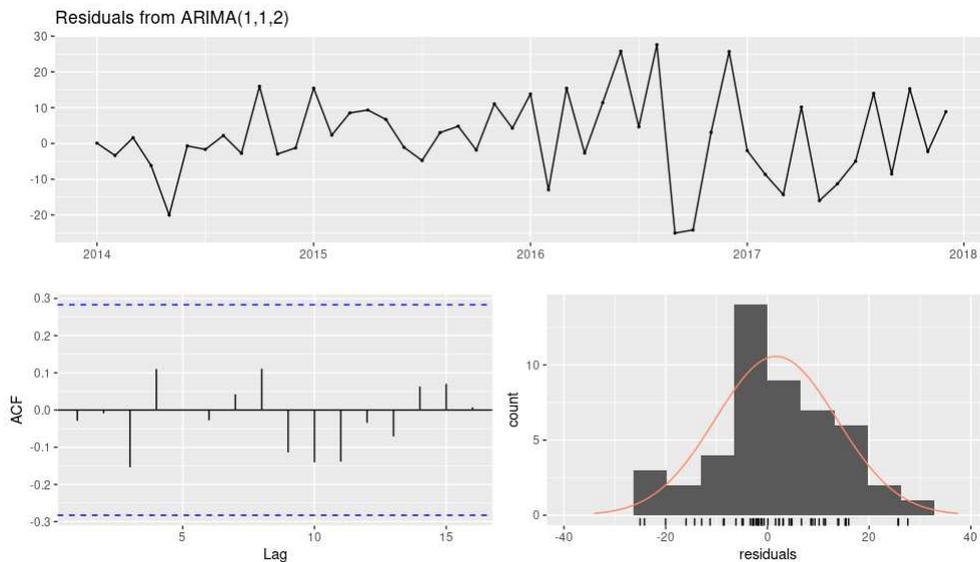


Figura 51 – Teste do Resíduo para o modelo da Série de Total de incidentes: ARIMA(1,1,2). Fonte: Do autor.

Tabela 20 – Resultados Teste Ljung-Box e Jarque-Bera para Total de incidentes

	ARIMA(1,1,2)
Ljung-Box	$p = 0.672$
Jarque-Bera	$p = 0.996$

No teste de Círculo Unitário dos Coeficientes da Figura 50, os valores dos coeficientes estarem dentro do círculo unitário indica que os mesmos são estáveis, confirmando que o modelo estimado é estacionário. No gráfico ACF contidos na Figura 51, observa-se que todos os valores de lag estão dentro do intervalo de confiança e na Tabela 20 como ambos os valores  $p$  rejeitam a hipótese nula dos testes de os resíduos não serem normalmente distribuídos ( $p > 0.05$ ), pode-se assumir que o resíduo do modelo se comporta como ruído branco.

Na Figura 52, tem-se o gráfico dos valores ajustados do modelo versus os dados observados. Observa-se que os valores ajustados parecem acompanhar a tendência da série, mas ficando distantes dos picos dos dados observados.

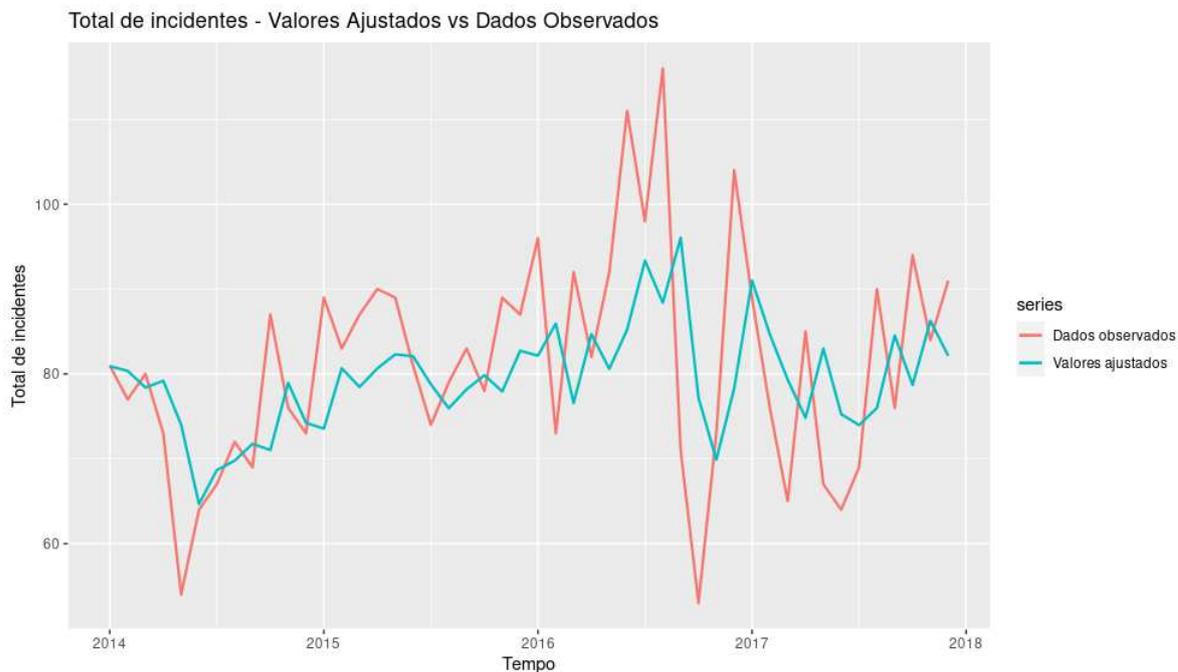


Figura 52 – Total de Incidentes - Valores Ajustados vs Dados Observados: ARIMA(1,1,2). Fonte: Do autor.

Com os modelos definidos e validados, foram feitas as previsões com os mesmos.

## 4.2 Previsões com o Modelo

Nessa seção será apresentada as previsões estimadas pelos modelos e, posteriormente, a comparação de acurácia entre cada período de estimativa.

A previsão foi feita para cada tipo de incidente e para o total de incidentes usando o período de estimativa de 2011 a 2017 e usando o período de estimativa 2014 a 2017 para prever os anos entre 2018 e 2020.

### 4.2.1 Previsão de 2018 a 2020 com o período de estimativa de 2011 a 2017

Nas Figuras de Crimes Cibernéticos 53 e Total de incidentes 57, percebe-se que os valores ajustados da previsão acompanham os dados observados e ficam bem próximos dos dados reais. Na Figura de Espionagem Cibernética 54, a curva de dados ajustados se encontra no valor médio dos picos da série. As previsões para o Hacktivismo na Figura 56 parecem seguir bem as tendências dos dados observados, por mais que se mostre um pouco distante. Pela Figura

55 da Guerra Cibernética, parece que há o pior ajuste entre os valores ajustados da previsão e a curva real dos dados, pois nem parece acompanhar a tendência da mesma.

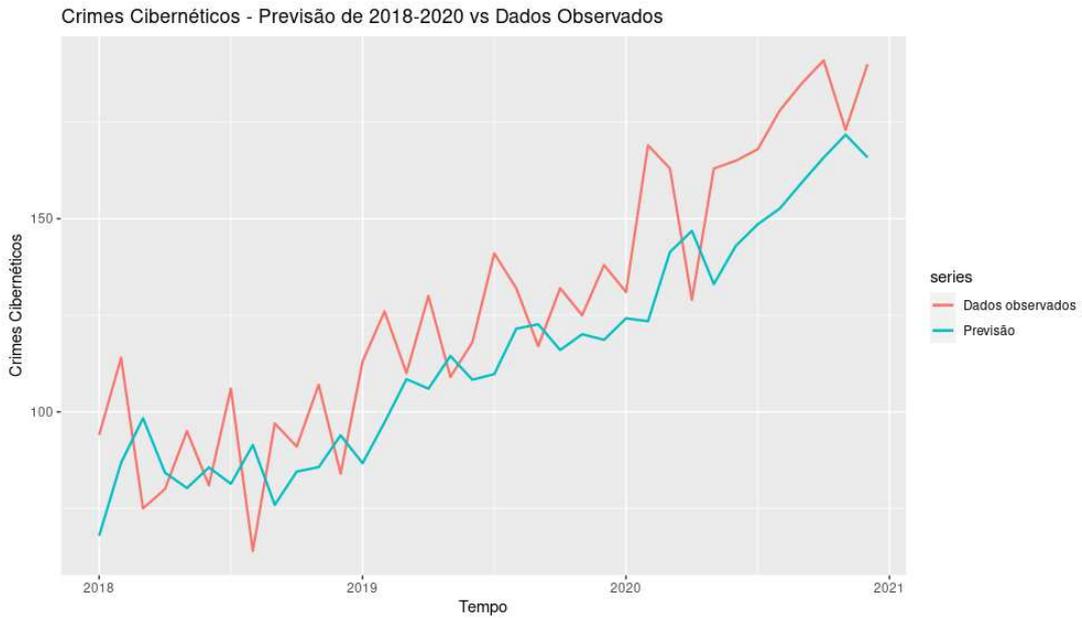


Figura 53 – Crimes Cibernéticos - Previsão de 2018 a 2020 usando todo o período de estimativa de 2011 a 2017

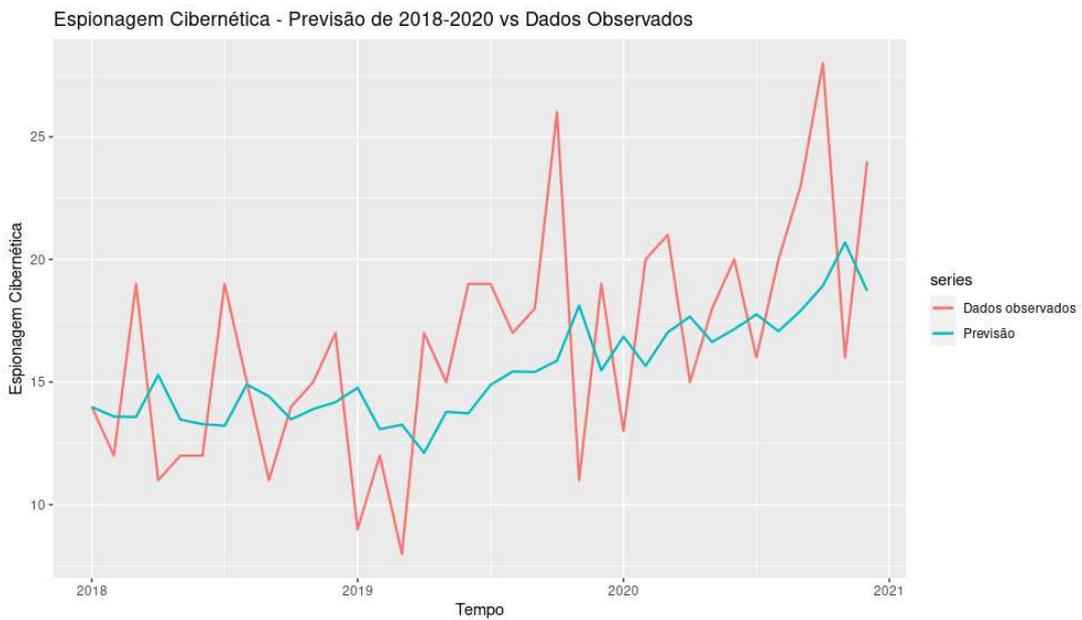


Figura 54 – Espionagem Cibernética - Previsão de 2018 a 2020 usando todo o período de estimativa de 2011 a 2017. Fonte: Do autor.

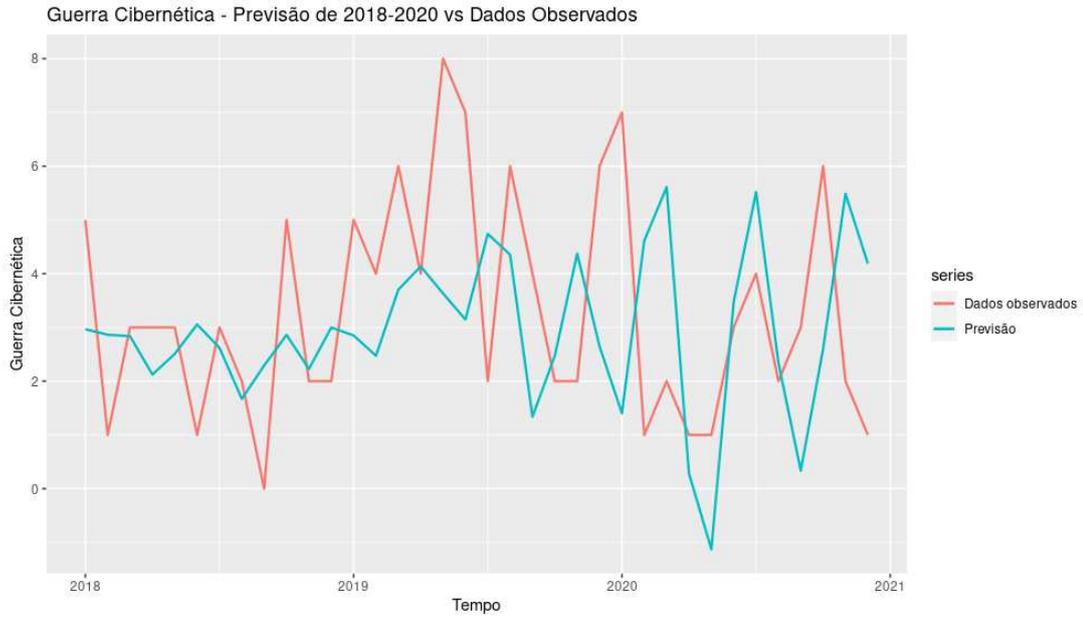


Figura 55 – Guerra Cibernética - Previsão de 2018 a 2020 usando todo o período de estimativa de 2011 a 2017. Fonte: Do autor.

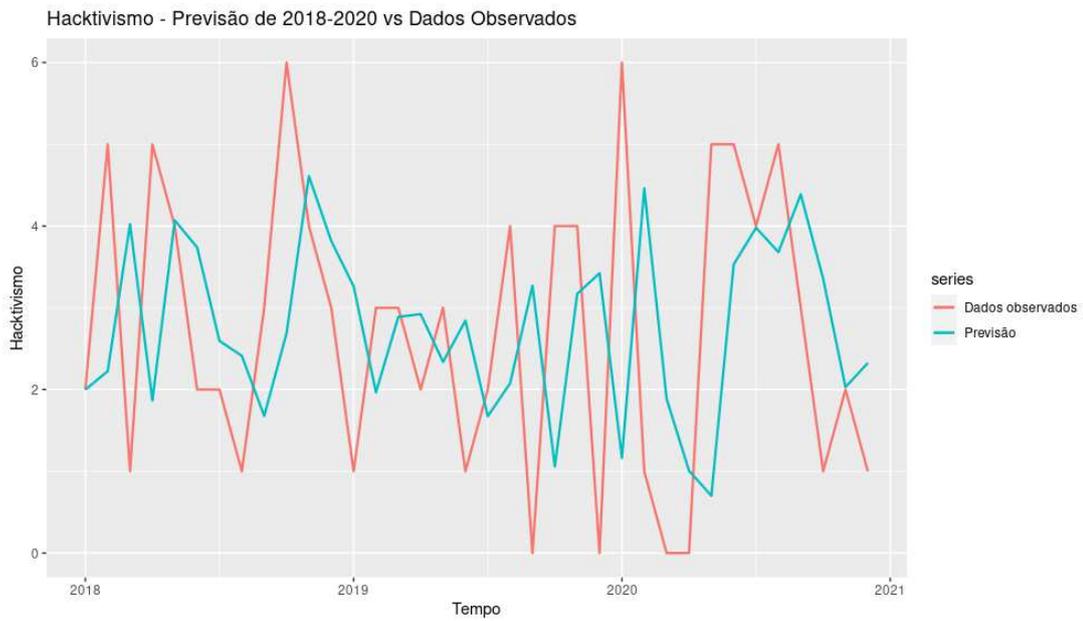


Figura 56 – Hacktivismo - Previsão de 2018 a 2020 usando todo o período de estimativa de 2011 a 2017. Fonte: Do autor.

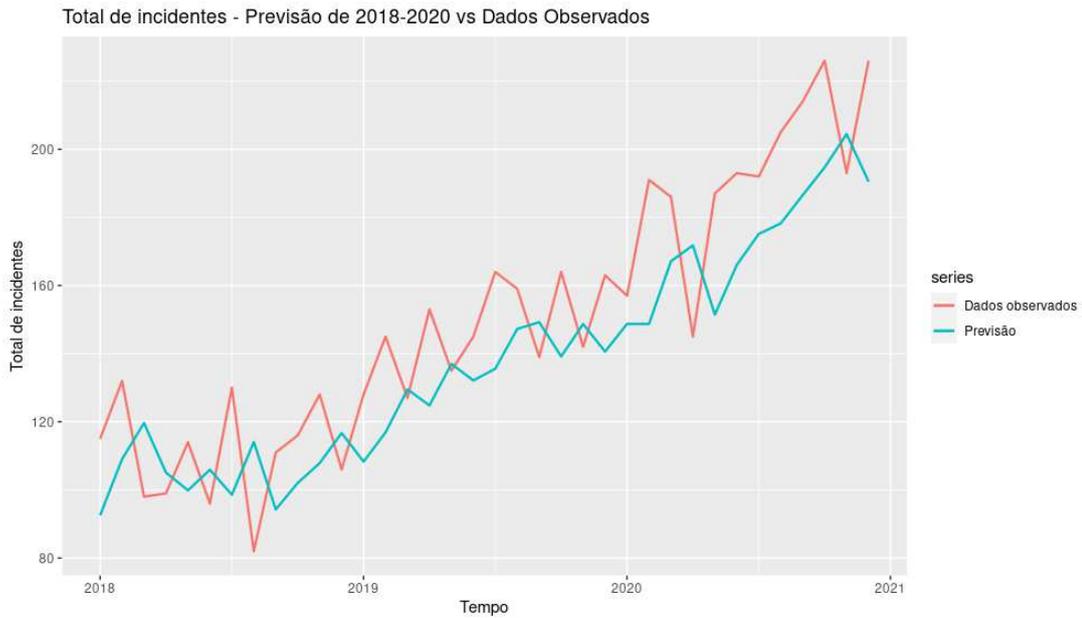


Figura 57 – Total de incidentes - Previsão de 2018 a 2020 usando todo o período de estimativa de 2011 a 2017. Fonte: Do autor.

#### 4.2.2 Previsão de 2018 a 2020 com o período de estimativa de 2014 a 2017

Nas Figuras de Crimes Cibernéticos 58 e Total de incidentes 62, percebe-se que os valores ajustados da previsão acompanham os dados observados e ficam bem próximos dos dados reais. Na Figura de Espionagem Cibernética 59, a curva de dados ajustados se encontra no valor médio dos picos da série. As previsões para o Hacking na Figura 61 parecem seguir bem as tendências dos dados observados, mas não capturando as peculiaridades dos extremos dos dados. Pela Figura 60 da Guerra Cibernética, parece que há o pior ajuste entre os valores ajustados da previsão e a curva real dos dados, pois nem parece acompanhar a tendência da mesma.

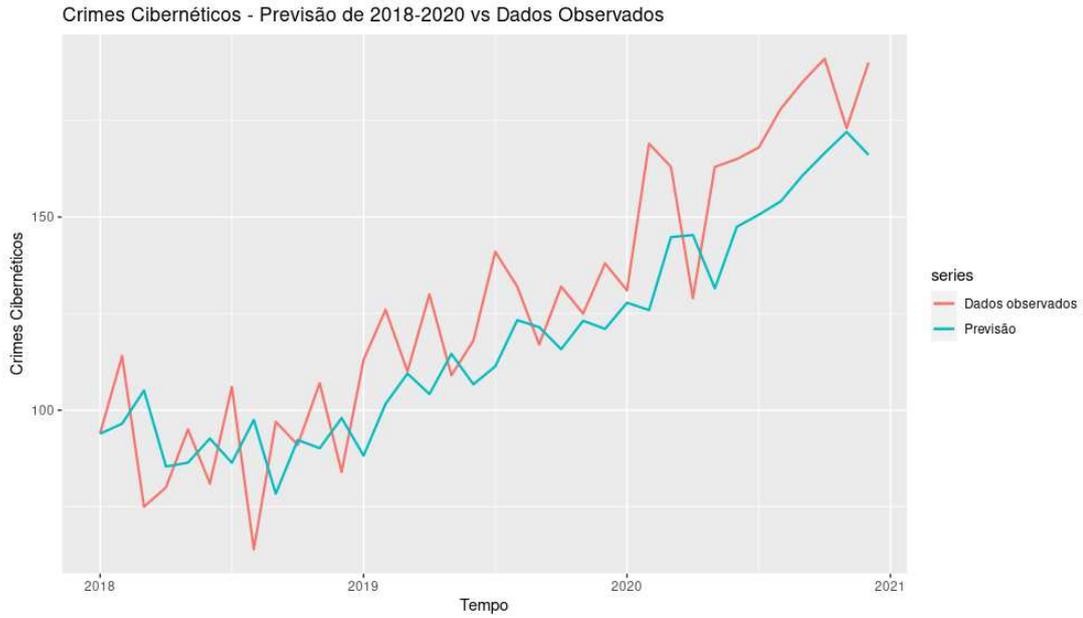


Figura 58 – Crimes Cibernéticos - Previsão de 2018 a 2020 usando todo o período de estimativa de 2014 a 2017. Fonte: Do autor.

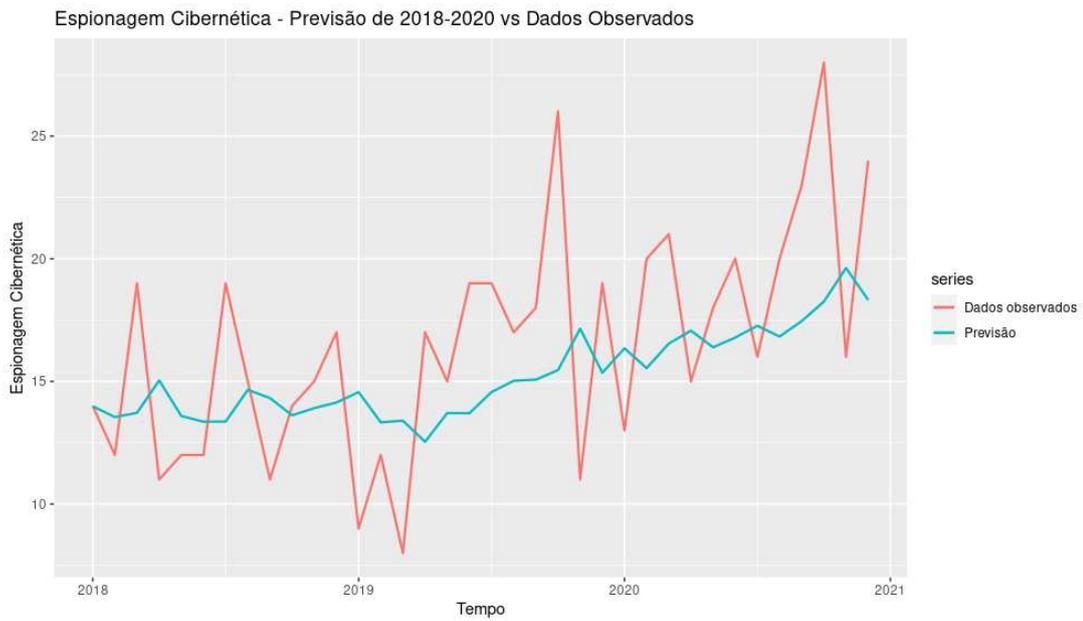


Figura 59 – Espionagem Cibernética - Previsão de 2018 a 2020 usando todo o período de estimativa de 2014 a 2017. Fonte: Do autor.

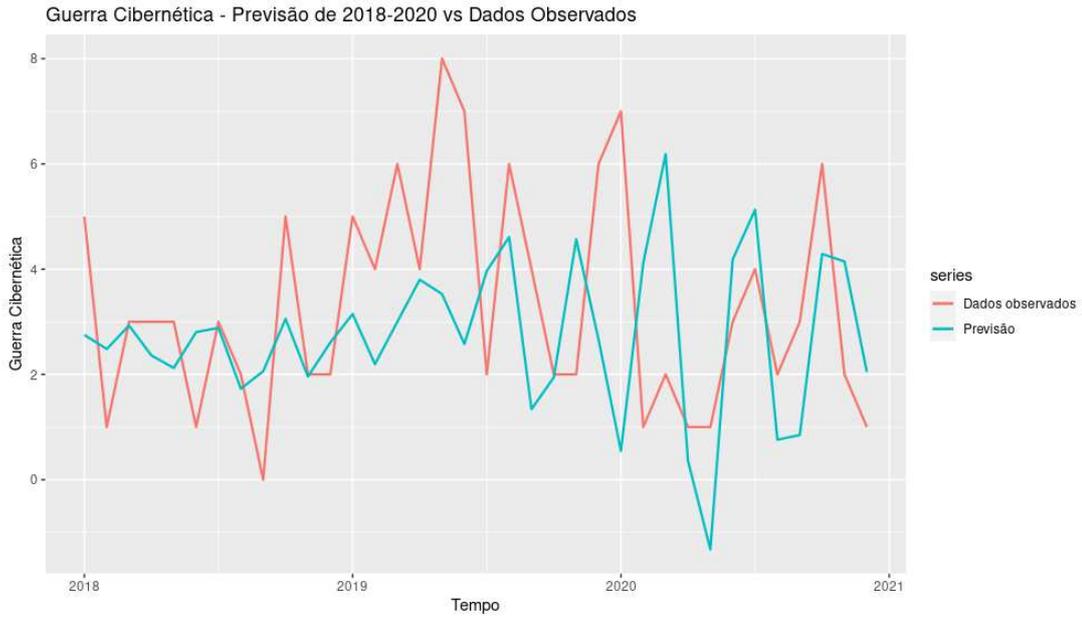


Figura 60 – Guerra Cibernética - Previsão de 2018 a 2020 usando todo o período de estimativa de 2014 a 2017. Fonte: Do autor.

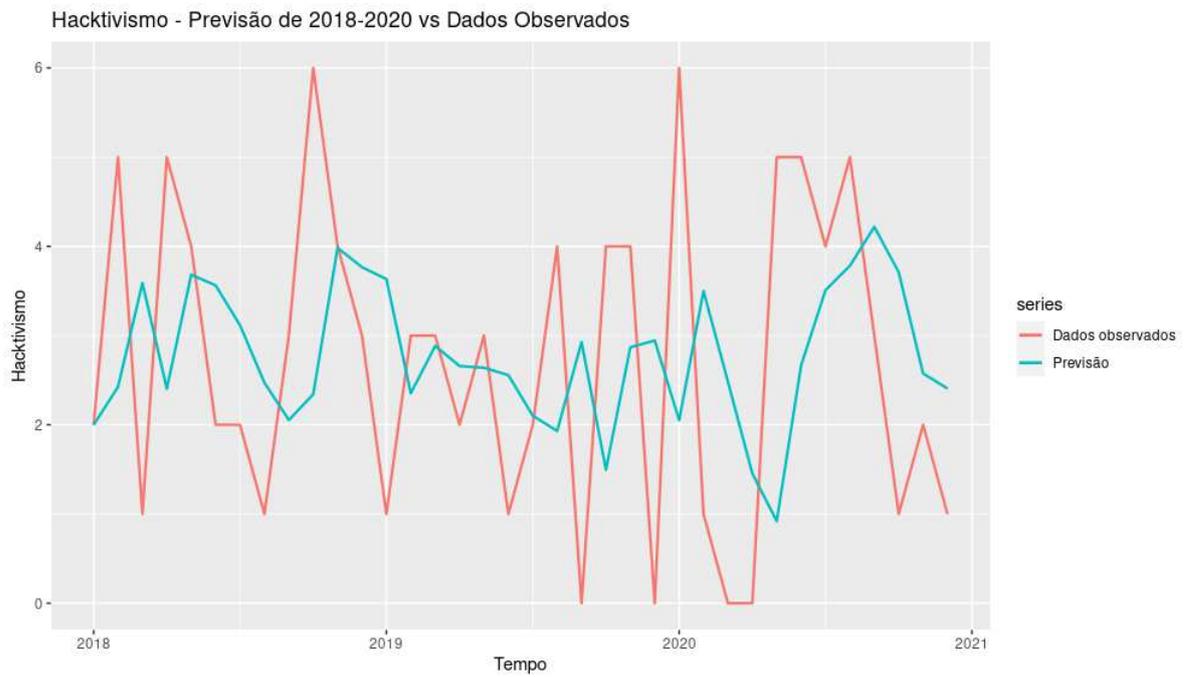


Figura 61 – Hacktivismo - Previsão de 2018 a 2020 usando todo o período de estimativa de 2014 a 2017. Fonte: Do autor.

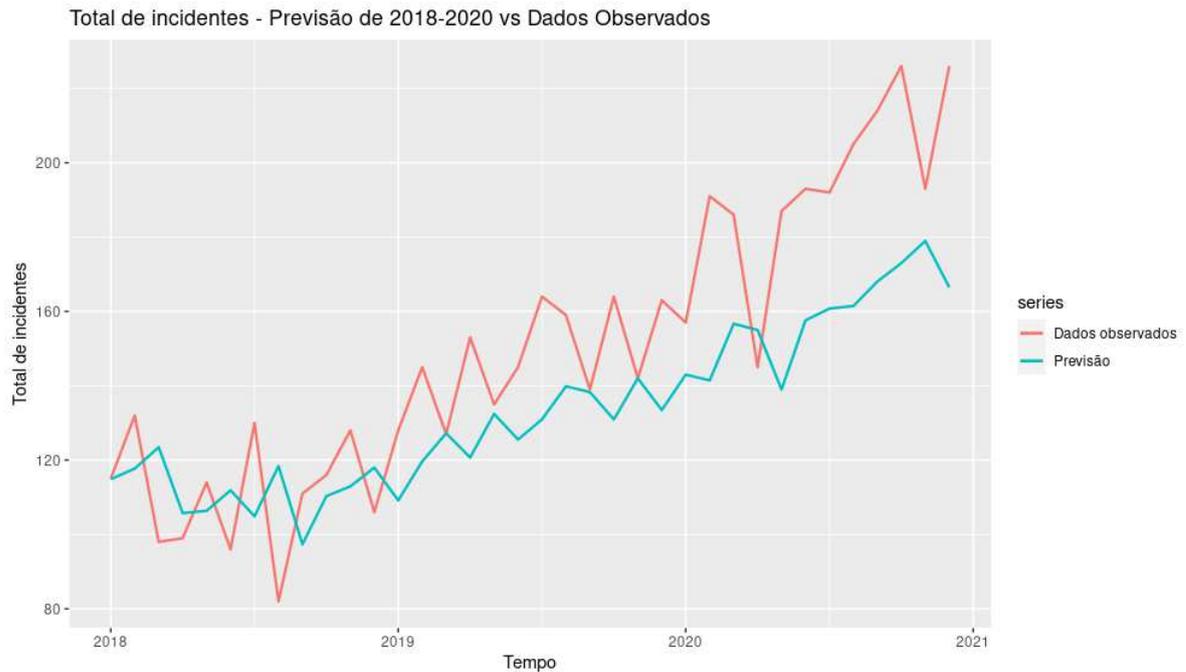


Figura 62 – Total de incidentes - Previsão de 2018 a 2020 usando todo o período de estimativa de 2014 a 2017. Fonte: Do autor.

#### 4.2.3 Comparativo de desempenho entre os dois períodos de estimativa

Como é mostrado na Tabela 21, quando observado o Erro Médio (ME) das séries, as de Espionagem Cibernética, Guerra Cibernética e Hacking acabaram se comportando melhor que as outras. Agora para o Erro Absoluto Médio (MAE), tem-se que as séries de Hacking se comportaram melhor por terem valores da métrica mais baixos do que as outras séries, enquanto as séries de Total de Incidentes tiveram a pior performance. Já para a última medida, o Erro Absoluto Percentual Médio (MAPE), para séries em qual se teve algum mês sem incidentes, não é possível estimar a métrica pois quando se tem zeros na série esse valor tende ao infinito. Nas séries que tem essa medida é observado independente da quantidade de tempo de estimativa do modelo o seu valor se mantém aproximadamente o mesmo (menos de 2% de variação), tendo Espionagem Cibernética o pior desempenho.

No geral, independente do volume da base de estimativa, os valores dos erros tiveram pouca variação, indicando que não necessariamente conhecer um período maior gerará previsões de maior desempenho.

Tabela 21 – Comparativo entre medidas de erros dos diferentes períodos de estimativa

	ME	MAE	MAPE
Total de incidentes (2011-2017)	12.440	20.219	13.733
Total de incidentes (2014-2017)	17.010	22.938	14.675
Crimes Cibernéticos (2011-2017)	12.839	18.312	14.926
Crimes Cibernéticos (2014-2017)	9,650	16.443	13.550
Espionagem Cibernética (2011-2017)	1.103	3.561	22.387
Espionagem Cibernética (2014-2017)	1.319	3.574	22.167
Guerra Cibernética (2011-2017)	0.372	2.005	<i>Inf</i>
Guerra Cibernética (2014-2017)	0.573	1.867	<i>Inf</i>
Hacktivismo (2011-2017)	-0.031	1.716	<i>Inf</i>
Hacktivismo (2014-2017)	-0.046	1.658	<i>Inf</i>

## 5 Conclusão

Com o aumento dos ataques cibernéticos, analisar as tendências futuras desses se tornou importante. O objetivo do trabalho foi tentar usar os modelos ARIMA para modelar séries temporais e fazer previsões em uma base de dados dos últimos 10 anos para categorias de incidentes de segurança (*Crimes Cibernéticos, Espionagem Cibernética, Guerra Cibernética e Hacktivismo*) e para o total somado dessas categorias.

Em resumo, as séries temporais construídas a partir dos dados da base *Hackmageddon* foram separadas por tipo de incidente, analisadas e modeladas utilizando modelos ARIMA. Usando o R como ferramenta para a modelagem, foram realizados testes estatísticos e construídos gráficos para auxiliar na definição do modelo, previsão de incidentes futuros e análise de desempenho entre os modelos. A fim de analisar o desempenho dos modelos, foram construídos gráficos comparativos entre os valores das previsões alcançadas e os dados observados e, em seguida, calculadas as medidas de erro para cada período de estimativa e tipo de incidente.

Ao longo do trabalho, foi possível observar uma previsão com desempenho médio de 85% para Total de Incidentes e Crimes Cibernéticos e de 78% para Espionagem Cibernética. Para as demais séries que tinham períodos sem incidentes, ou seja, quando existia zeros ao longo da série original, não foi possível calcular essa métrica de erro percentual, pois a fórmula da mesma tende ao infinito quando isso ocorre. Entretanto, independente do período usado para estimativa, as métricas de erro tiveram pouca variação e para certos tipos de incidente um período menor de estimação teve menos erro, enquanto para outras, teve-se o contrário, podendo-se dizer que não necessariamente um maior ou menor período de estimativa trará previsões de maior desempenho para modelos ARIMA. Com base nisso, para trabalhos futuros, a fim de obter previsões mais precisas, é preciso fazer a modelagem e posterior previsão usando outros modelos de séries temporais e também podendo utilizar técnicas de aprendizado de máquina.

# Referências

- ANTONAKAKIS, M.; APRIL, T.; BAILEY, M.; BURSZTEIN, E.; COCHRAN, J.; DURUMERIC, Z.; HALDERMAN, J. A.; MENSCHER, D.; SEAMAN, C.; SULLIVAN, N.; THOMAS, K.; ZHOU, Y.; BURSZTEIN, M. A. T. A. M. B. M. B. E.; INVERNIZZI, B. J. J. C. Z. D. A. H. L.; MA, B. M. K. D. K. C. L. Z.; MENSCHER, J. M. D.; THOMAS, B. C. S. N. S. K.; ZHOU, B. Y. This paper is included in the Proceedings of the 26th USENIX Security Symposium Understanding the Mirai Botnet Understanding the Mirai Botnet. 2017. Disponível em: <<https://www.usenix.org/conference/usenixsecurity17/technical-sessions/presentation/antonakakis>>. Citado na página 15.
- Cert.br. *Cartilha de Segurança para Internet*. 2020. Disponível em: <<https://cartilha.cert.br/>>. Citado na página 18.
- CONDON, E.; HE, A.; CUKIER, M. Analysis of computer security incident data using time series models. *Proceedings - International Symposium on Software Reliability Engineering, ISSRE*, p. 77–86, 2008. ISSN 10719458. Citado 4 vezes nas páginas 15, 16, 21 e 22.
- CRYER, J. D.; CHAN, K.-S. *Springer Texts in Statistics Time Series Analysis With Applications in R Second Edition*. [S.l.], 2008. Citado na página 19.
- DICKEY, D. A.; FULLER, W. A. Distribution of Estimators of Autoregressive Time Series with a Unit Root. *Journal of the American Statistical Association*, v. 74, n. 366, p. 427–431, 6 1979. Citado na página 26.
- EDWARDS, B.; HOFMEYR, S.; FORREST, S. Hype and heavy tails: A closer look at data breaches. *Journal of Cybersecurity*, v. 2, n. 1, p. 3–14, 2016. ISSN 20572093. Citado 3 vezes nas páginas 15, 21 e 22.
- EHRENFELD, J. M. WannaCry, Cybersecurity and Health Information Technology: A Time to Act. *Journal of Medical Systems*, Journal of Medical Systems, v. 41, n. 7, p. 10916, 2017. ISSN 0148-5598. Citado na página 15.
- FONTES, E. L. G. *Segurança da informação*. [S.l.]: Editora Saraiva, 2017. Citado na página 15.
- GUALBERTO, E. S.; SOUSA, R. T. D.; DEUS, F. E. G. D.; DUQUE, C. G. *Proposição de uma Ontologia de Apoio à Gestão de Riscos de Segurança da Informação*. [S.l.], 2013. v. 6, n. 1, 30–43 p. Citado na página 18.
- HYNDMAN, R. J.; ATHANASOPOULOS, G. *Forecasting Principles and Practices*. 2013. Citado 2 vezes nas páginas 18 e 19.
- HYNDMAN, R. J.; ATHANASOPOULOS, G. *Forecasting: principles and practice*. 2nd edition. ed. Melbourne, Australia.: [s.n.], 2018. Citado 7 vezes nas páginas 7, 19, 20, 21, 27, 28 e 29.
- JARQUE, C. M.; BERA, A. K. Efficient tests for normality, homoscedasticity and serial independence of regression residuals. *Economics Letters*, v. 6, n. 3, p. 255–259, 1980. ISSN 0165-1765. Disponível em: <<https://www.sciencedirect.com/science/article/pii/0165176580900245>>. Citado na página 29.

Jonathan Stempel; Jim Finkle. *Yahoo 2013 Account Security Update FAQs*. 2017. Disponível em: <<https://help.yahoo.com/kb/account/SLN28451.html>>. Citado na página 15.

KESAN, J. P.; ZHANG, L. *Analysis of Cyber Incident Categories Based on Losses*. [S.l.], 2020. Disponível em: <<https://doi.org/0>>. Citado na página 22.

KUYPERS, M. A.; MAILLART, T.; PATÉ-CORNELL, E. An Empirical Analysis of Cyber Security Incidents at a Large Organization. 2017. Disponível em: <[https://fsi.stanford.edu/sites/default/files/kuypersweis\\_v7.pdf](https://fsi.stanford.edu/sites/default/files/kuypersweis_v7.pdf)>. Citado 4 vezes nas páginas 15, 16, 21 e 22.

LIMA, R. *Avaliando previsões com o R*. 2015. Disponível em: <<https://analisemacro.com.br/economia/macroeconometria/avaliando-previsoes-com-o-r/>>. Citado na página 30.

LIU, M.; LIU, Y.; ZHANG, J.; BAILEY, M.; KARIR, M.; SARABI, A. Predicting Cyber Security Incidents Using Feature-Based Characterization of Network-Level Malicious Activities. p. 3–9, 2015. Citado na página 16.

LIU, Y.; SARABI, A.; ZHANG, J.; NAGHIZADEH, P.; KARIR, M.; BAILEY, M. Cloudy with a Chance of Breach : Forecasting Cyber Security Incidents This paper is included in the Proceedings of the. *24th USENIX Security Symposium (USENIX Security 15)*, 2015. Citado 3 vezes nas páginas 16, 21 e 22.

MIRANDA, I. P. D. Comparação de diferentes Métodos de Previsão em Séries Temporais com valores discrepantes. 2014. Citado na página 31.

MUSHTAQ, R. *Testing times series data for stationarity*. [S.l.], 2011. Disponível em: <<http://ssrn.com/abstract=1911068>>. Citado na página 26.

PASSERI, P. *Hackmageddon.com*. 2018. Disponível em: <<http://hackmageddon.com/>>. Citado 4 vezes nas páginas 7, 16, 24 e 25.

PRC. *Privacy Rights Clearinghouse*. 2019. Disponível em: <<http://www.privacyrights.org/data-breach>>. Citado na página 23.

REIS, M. M. *Estatística Aplicada à Administração I*. [S.l.], 2016. Citado na página 18.

REUTERS. *Yahoo says all three billion accounts hacked in 2013 data theft*. 2017. Disponível em: <<https://www.reuters.com/article/us-yahoo-cyber/yahoo-says-all-three-billion-accounts-hacked-in-2013-data-theft-idUSKCN1C8201>>. Citado na página 15.

SAX, C.; EDDELBUETTEL, D. *Seasonal Adjustment by X-13ARIMA-SEATS in R*. [S.l.], 2018. Disponível em: <<http://www.seasonal.website>>. Citado na página 27.

SOLMS, R. V.; NIEKERK, J. V. From information security to cyber security. *Computers and Security*, Elsevier Ltd, v. 38, p. 97–102, 2013. ISSN 01674048. Citado na página 17.

STALLINGS, W. *Criptografia e Segurança de Redes - 6ª Ed.* 2014. 2014. Citado na página 17.

VAIDYA, R. Statistical Release: Cyber Security Breaches Survey 2018. p. 1–58, 2018. Disponível em: <[https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment\\_data/file/702074/Cyber\\_Security\\_Breaches\\_Survey\\_2018\\_-\\_Main\\_Report.pdf](https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/702074/Cyber_Security_Breaches_Survey_2018_-_Main_Report.pdf)>. Citado na página 15.

VAIDYA, R. *Cyber Security Breaches Survey 2020*. [S.l.], 2020. Citado na página 15.

VERIZON. *VERIS Community Database*. 2019. Disponível em: <<https://github.com/vz-risk/VCDB>>. Citado na página 24.

WASC. *Web Hacking Incident Database*. 2019. Disponível em: <<http://projects.webappsec.org/w/page/13246995/Web-Hacking-Incident-Database>>. Citado na página 24.

XU, M.; SCHWEITZER, K. M.; BATEMAN, R. M.; XU, S. Modeling and Predicting Cyber Hacking Breaches. *IEEE Transactions on Information Forensics and Security*, Institute of Electrical and Electronics Engineers Inc., v. 13, n. 11, p. 2856–2871, 11 2018. ISSN 15566013. Citado na página 22.