

Detecção de Malwares Android: reprodução da seleção de características do SigPID

Joner Assolin¹, Vanderson Rocha², Guilherme Silveira¹, Gustavo Rodrigues¹,
Eduardo Feitosa², Karina Casola¹, Diego Kreutz¹

¹ Universidade Federal do Pampa (Unipampa)

² Universidade Federal do Amazonas (UFAM)

{NomeSobrenome}.aluno, diegokreutz}@unipampa.edu.br

{vanderson, efeitosa}@icomp.ufam.edu.br

Resumo. Para atacar o desafio de escalabilidade na detecção de malwares em Android, há trabalhos que propõem a utilização de um número reduzido de permissões, como é o caso do SigPID. Neste trabalho, apresentamos a reprodução dos 3 níveis de seleção de permissões e avaliação dos principais métodos de aprendizagem do SigPID, utilizando um conjunto de dados publicamente disponível. Nosso estudo inicial indica que o número de permissões impacta o tempo de treinamento e execução, bem como a acurácia dos modelos. Entretanto, o tempo de execução pode não ser significativo a ponto de justificar um número menor de permissões para detecção de malwares.

1. Introdução

Dentre os métodos para detecção de *malwares* em aplicativos Android, os que utilizam aprendizado de máquina vêm ganhando destaque [Wu et al., 2021]. Independente de focarem suas abordagens na análise estática, dinâmica ou híbrida, esses trabalhos utilizam as permissões do Android para o desenvolvimento de modelos de detecção de *malwares* [Alsoghyer and Almomani, 2020]. Entretanto, utilizar todas as 247 permissões das APIs do Android, disponíveis para treinamento dos modelos de aprendizado de máquina, pode representar um desafio de escalabilidade [Yildiz and Doğru, 2019, Li et al., 2018] e impactar negativamente no tempo de execução das soluções.

Com o objetivo de mitigar o problema da escalabilidade, há trabalhos (e.g., [Li et al., 2018, Yildiz and Doğru, 2019]) que investigaram o impacto da redução dos números de permissões utilizadas para o treino dos modelos. Como resultado, verificaram que, mesmo utilizando um número menor de permissões, o tempo de execução pode ser reduzido sem comprometer de forma significativa a acurácia do modelo (i.e., aumento da escalabilidade sem comprometer o desempenho da classificação).

Neste trabalho avaliamos e discutimos a reprodutibilidade e o desempenho do trabalho de [Li et al., 2018, Sun et al., 2016] (SigPID - *Significant Permission Identification for Android Malware Detection*), que pode ser considerado um dos mais relevantes e mais bem citados (mais de 320 citações *Google Scholar Citations* - GSC, em setembro de 2021) sobre escalabilidade de modelos de detecção de *malwares* Android. Como o conjunto de dados empregado no trabalho original não está disponível, utilizamos um conjunto de dados de conhecimento público, que contém 113 permissões. Numa primeira etapa, reproduzimos a redução de características utilizada no SigPID, que consiste em três níveis

seleção de permissões. A aplicação dos três níveis de seleção resultou em 27 permissões mais significativas. Em seguida, avaliamos o desempenho do modelo em comparação com diferentes conjuntos de permissões, a saber: (a) 113 permissões (*baseline*) contidas no conjunto de dados; (b) 22 permissões identificadas no trabalho original do SigPID; (c) 32 permissões mais recorrentes em trabalhos de detecção de *malwares* Android; e (d) 22 permissões classificadas como perigosas pela Google ¹.

Como contribuição do trabalho, podemos destacar: (a) a implementação da estratégia de seleção de características empregada pelo SigPID em um conjunto de dados público; (b) a identificação de um subconjunto essencial de permissões (permissões significativas) que pode ser utilizado para identificar efetivamente *malwares* no Android; (c) um comparativo com o trabalho original, que identifica 22 permissões como significativas e também a comparação com outros conjuntos de dados com diferentes quantidades de permissões; e (d) a análise de aspectos de reprodutibilidade do trabalho original.

O trabalho está organizado como segue. Na Seção 2 apresentamos os requisitos para a reprodução do SigPID e detalhamos a metodologia de seleção das características. Nas Seções 3 e 4 apresentamos os resultados e as considerações finais.

2. Reprodução dos experimentos

2.1. Detalhamento do Ambiente

Para o desenvolvimento e avaliação dos experimentos, utilizamos um notebook com processador Intel Celeron 1007U (1.5GHz, Dual Core, 2MB L2), 4GB DDR3 1.600MHz, disco rígido de 320GB (SATA - 5.400rpm), Windows 10 Home Single Language, compilação 19042.1110. Para a implementação e avaliação dos modelos, utilizamos as ferramentas Jupyter Notebook (IPython 7.12.0, Python 3.7.6 (default, jan. 8 2020) e o Google Chrome Versão 91.0.4472.124 (Versão oficial) 64 bits. Com exceção do algoritmo *Functional Tree*, versão 1.0.4, implementado com a ferramenta Weka versão 3.9.5, os demais foram implementados utilizando a versão 0.22.1 da biblioteca Scikit-learn.

Para análise e uso do conjunto de dados, utilizamos uma divisão estratificada pseudo-aleatória (*test_size*) de 70%/30% [James et al., 2013], a partir dos dados iniciais, sendo 70% utilizado para treinos e 30% para testes. As divisões estratificadas são desejáveis em casos de conjuntos de dados desbalanceados, como é o caso do escolhido.

Para garantir a reprodutibilidade do experimento, definimos arbitrariamente a semente aleatória como 1 para `train_test_split`, de forma a controlar a seleção dos dados de treino e teste. Já os hiperparâmetros, variáveis que controlam o próprio processo de treinamento, foram seguidos conforme o padrão da biblioteca Scikit-learn.

2.2. Conjunto de dados

Como mencionado, o conjunto de dados original do trabalho está indisponível. Para a reprodução e comparação do SigPID com diferentes conjuntos de permissões, selecionamos o conjunto de dados Drebin_215, disponível publicamente no FigShare ², um subconjunto do *Drebin project* [Arp et al., 2014], por ser de acesso público e possuir permissões

¹Ao total, a Google define 30 permissões como perigosas. Destas, 22 estão presentes no conjunto de dados analisado.

²https://figshare.com/articles/dataset/Android_malware_dataset_for_machine_learning_2/5854653

do Android como características. O Drebin_215 possui 215 características extraídas de 15.036 aplicativos (5.560 malignos e 9.476 benignos), sendo que 113 características são permissões.

2.3. Seleção das Permissões através do MLDP

O MLDP (*Multi-Level Data Pruning*) é um método multinível (de 3 níveis) para seleção de características relevantes para uso na construção do modelo. A ideia por trás do método é diminuir o número de permissões (selecionando as mais significativas) e, conseqüentemente, o tempo de execução dos modelos. O MLDP assume, como parâmetro de seleção, uma taxa de detecção de *malware* acurácia e precisão de no mínimo 90%, sendo essa considerada uma taxa muito boa.

O método opera nos seguintes três níveis de seleção: (1) classificação de permissão com taxa negativa (*Permission Ranking Negative Rate* ou PRNR), onde o objetivo é selecionar as permissões que geram as maiores métricas (*e.g.*, acurácia, recall, F1-score); (2) classificação de permissão baseada em suporte (*Support Based Permission Ranking* ou SPR), onde o objetivo é selecionar o menor número de permissões que alcançam uma taxa maior que 90% de acurácia; e (3) mineração de permissões com regras de associação (*Permission Mining with Association Rules* ou PMAR), para eliminar permissões que geram redundância ao modelo. Cada um dos três níveis é detalhado na versão estendida do trabalho, disponível em [Assolin et al., 2021]. Para reprodutibilidade, os conjuntos de dados e códigos estão disponíveis online no GitHub³. Ao final, chegamos a seleção de 108, 30 e 27 permissões nos níveis 1, 2 e 3 de seleção do MLDP, respectivamente.

3. Resultados

Para a reprodução do SigPID, utilizamos o algoritmo *Support Vector Machine* (SVM), conforme proposto pelos autores no trabalho original [Li et al., 2018]. A Tabela 1 apresenta os resultados da execução do SVM para os diferentes conjuntos de dados, incluindo métricas para a avaliação do desempenho da detecção e o tempo de execução.

Tabela 1. Métricas de avaliação SVM

Conjunto de Dados	Quantidade de Permissões	Precisão	Recall	FPR	F1_Score	Acurácia	Tempo Execução (s)
Nível 1 (PRNR)	108	93,62	88,01	3,52	90,73	93,35	5,44
Nível 2 (SPR)	30	90,03	82,25	5,35	85,96	90,07	2,41
Nível 3 (PMAR)	27	90,13	82,07	5,28	85,91	90,05	2,26
Baseline	113	93,49	87,83	3,59	90,57	93,24	5,84
Perigosas Google	22	86,76	71,88	6,44	78,62	85,55	2,34
Recorrentes	32	88,54	81,53	6,19	84,89	89,27	3,17
SigPID	22	91,77	74,88	3,94	82,47	88,23	2,62

Como podemos observar, os conjuntos de dados que utilizam o mesmo número de permissões obtiveram resultados diferentes. As 22 permissões identificadas no SigPID se destacam em todas as métricas de avaliação quando comparadas com as 22 permissões consideradas perigosas pela Google, ou seja, utilizar permissões perigosas não leva a um

³<https://github.com/Malware-Hunter/SigPID>

melhor resultado qualitativo. Isto indica que a escolha das permissões possui, de fato, um impacto no desempenho na detecção de *malwares*.

Quando reduzimos o número de permissões de 113 (*baseline*) para 108 com o nível 1 do MLDP, alcançamos taxas de precisão e acurácia mais altas, acima de 93%. Apesar de haver aumento no *recall*, o F1-Score e a taxa de falsos positivos (FPR) permaneceu mais baixa em relação ao *baseline*, assim como o tempo de execução também diminuiu. Quando reduzimos o número de permissões de 113 para 30 com o nível 2 do MLDP, mantiveram-se a acurácia e precisão na faixa de 90%, porém com F1-Score abaixo de 90%. Além disso, obtivemos um ganho de mais de 2 segundos em relação ao tempo de execução.

Além do SVM, avaliamos também outros três algoritmos (*Random Forest*, *Decision Tree* e *Functional Trees* [Gama, 2004]) e comparamos os resultados com os conjuntos de dados discriminados na Tabela 2.

Tabela 2. Conjuntos de dados

Nº de Permissões	Conjunto de dados	Observação
113	<i>Baseline</i>	Contidas no conjunto de dados Drebin_215
22	SigPID	Identificadas no trabalho original do SigPID
32	Recorrentes	Identificadas por meio da interseção de permissões identificadas em outros trabalhos
22	Perigosas	Permissões na lista de perigosas da Google e estavam contidas na <i>baseline</i> deste trabalho
27	MLDP	Identificadas após aplicação dos três níveis de seleção sobre a <i>baseline</i>

A Tabela 3 sintetiza os resultados dos algoritmos *Random Forest*, *Decision Tree* e *Functional Trees*. As métricas apresentadas são a acurácia, que indica o desempenho geral do modelo, e F1-Score, a média harmônica entre o *recall* e a precisão.

Tabela 3. Métricas de avaliação dos conjuntos de dados

Conjunto de dados	Decision tree		Random forest		Functional Trees	
	F1_Score	Acurácia	F1_Score	Acurácia	F1_Score	Acurácia
113 Baseline	92,28	94,32	94,44	95,94	97,30	97,27
22 SigPID	88,46	92,00	89,30	92,57	92,90	93,05
27 MLDP	91,09	93,50	92,62	94,64	95,40	95,43
32 Recorrentes	90,89	93,44	92,46	94,57	95,60	95,63
22 Perigosas Google	81,62	87,74	85,31	88,85	89,10	89,02

Quatro dos cinco conjuntos de dados obtiveram acurácia acima dos 90% (*Baseline*, SigPID, MLDP e Recorrentes). Destes, o *Baseline* performou melhor que os demais. Quando utilizamos o *Decision Tree*, atingimos 94,32% de acurácia e 92,28% de F1-Score. Com os algoritmos *Random Forest* e *Functional Trees*, atingimos 95,94% de acurácia e 94,44% de F1-Score e acurácia e F1-Score acima dos 97%, respectivamente.

Utilizando 22 permissões do SigPID, a acurácia se manteve na faixa de 92% com *Decision Tree* e *Random Forest*. Já F1-Score se manteve abaixo de 90%. O *Functional Trees*, com o melhor desempenho, atingiu 93,05% de acurácia e 92,90% de F1-Score. A

Tabela 3 ainda mostra que os conjuntos de dados MLDP e Recorrentes obtiveram resultados próximos, ficando ambos com acurácia acima de 95% para o *Functional Trees*.

Na Figura 1 apresentamos os dados sobre o tempo de execução de cada algoritmo. Como podemos observar, o *baseline* tem o maior tempo de execução, chegando a 11,31 segundos com algoritmo *Functional Trees*. Entretanto, é interessante observar que o problema ocorre apenas para o *baseline*, já que o algoritmo executa em menos tempo que alguns dos demais para conjuntos de dados menores.

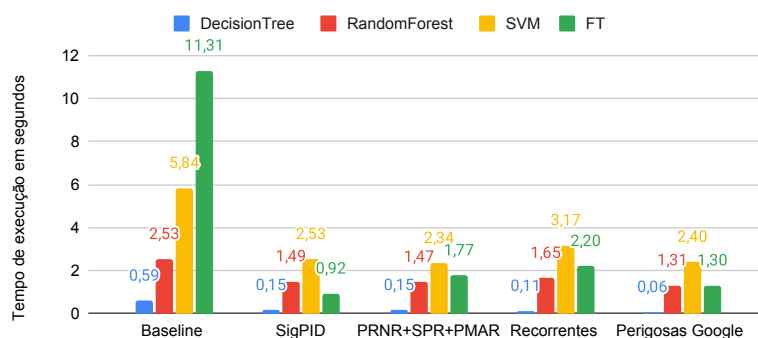


Figura 1. Tempo de execução para diferentes conjuntos de dados

O *Functional Trees* é um algoritmo baseado em árvores de classificação que podem ter funções de regressão logística nos nós internos ou folhas. Por ser capaz de lidar com variáveis binárias e multi-classe, atributos numéricos e valores ausentes, seu tempo de execução pode aumentar de acordo com o número de características de entrada [Gama, 2004]. Como a estrutura do algoritmo é uma generalização das *Multivariate Trees* [Gama, 2001], sua complexidade é, $\mathcal{O}(n^2)$, o que explica o comportamento apresentado no gráfico da Figura 1. De fato, o *Functional Trees* reduziu significativamente o tempo de execução para conjuntos de dados menores (e.g., 1,77s para as 27 permissões do MLDP).

Ao analisarmos os conjuntos de dados que utilizam 27 e 32 permissões, observamos que o tempo de execução e as métricas estão muito próximas para a *Functional Trees*. Nesse caso, podemos dizer que os modelos são equivalentes, isto é, não há diferença em utilizar as 32 permissões recorrentes ou as 27 permissões do MLDP. Isto indica que a identificação das permissões recorrentes, em trabalhos existentes na literatura, leva a resultados tão bons quanto os resultados do método de múltiplos níveis de seleção do SigPID, o que confirma uma das nossas hipóteses, isto é, permissões recorrentes podem ter um impacto positivo sobre os modelos de detecção de *malwares*.

4. Considerações Finais

Aplicando os 3 níveis de seleção do SigPID nas 113 permissões do conjunto de dados Drebin_215, conseguimos reduzir em 76% o número de permissões a serem analisadas para detecção de *malwares* Android, mantendo acurácia acima de 90% com SVM, 93,50% com *Decision Tree*, 94,64% com *Random Forest* e 95% com o *Functional Trees*. Além disso, conseguimos também reduzir o tempo de execução dos modelos, porém, ao custo de um leve aumento na taxa de falsos positivos.

Além das 27 permissões resultantes dos 3 níveis de seleção, utilizamos também conjuntos de 113, 22 e 32 permissões. Percebemos que o conjunto de dados que utiliza

todas permissões (113) foi o que melhor performou (e.g., 97% de acurácia), porém, ao preço de um tempo de execução significativamente maior. Um caso interessante ocorreu ao compararmos os conjuntos com 22 permissões, sendo um oriundo do trabalho original do SigPID e 22 permissões classificadas como perigosas pela Google. O SigPID chegou a 93% de acurácia enquanto que as perigosas mantiveram acurácia abaixo de 90%, o que indica que o fato da permissão ser classificada como perigosa não a torna necessariamente relevante para detecção de *malwares*. Outra observação interessante é o fato de os conjuntos de dados com 32 permissões mais recorrentes e das 27 identificadas atingirem resultados muito próximos ao aplicar os 3 níveis de seleção. Isto indica que escolher as permissões de acordo com a recorrência pode ser um caminho a ser investigado.

Como trabalhos futuros, podemos elencar: (a) testes com conjuntos de dados maiores; (b) testes com conjuntos de dados atuais; (c) avaliação dos níveis de seleção para outras características (e.g., *intents* e chamadas de API); (d) otimização de hiperparâmetros; (e) testar os modelos em *smartphones* modernos; e (f) mensurar o tempo de execução dos modelos em CPUs modernas projetadas para acelerar a computação de algoritmos de aprendizado de máquina.

Agradecimentos

Esta pesquisa foi financiada, conforme previsto nos Arts. 21 e 22 do decreto nº 10.521/2020, nos termos da Lei Federal nº 8.387/1991, através de convênio nº003/2021, firmado entre ICOMP/UFAM, Flextronics da Amazônia Ltda e Motorola Mobility Comércio de Produtos Eletrônicos Ltda.

Referências

- Alsoghyer, S. and Almomani, I. (2020). On the effectiveness of application permissions for android ransomware detection. In *2020 6th Conference on Data Science and Machine Learning Applications (CDMA)*, pages 94–99.
- Arp, D., Spreitzenbarth, M., Hubner, M., Gascon, H., Rieck, K., and Siemens, C. (2014). Drebin: Effective and explainable detection of android malware in your pocket. In *Ndss*, volume 14, pages 23–26.
- Assolin, J., Rocha, V., Kreutz, D., Siqueira, G., Rodrigues, G., Feitosa, E., and Casola, K. (2021). Detecção de Malwares Android: reprodução da seleção de características do SigPID. https://arxiv.kreutz.xyz/wrseg2021_sigpid_vel.pdf.
- Gama, J. (2001). Functional trees for classification. In *Proceedings 2001 IEEE International Conference on Data Mining*, pages 147–154. IEEE.
- Gama, J. (2004). Functional trees. *Machine learning*, 55(3):219–250.
- James, G., Witten, D., Hastie, T., and Tibshirani, R. (2013). *An introduction to statistical learning*, volume 112. Springer.
- Li, J., Sun, L., Yan, Q., Li, Z., Srisa-an, W., and Ye, H. (2018). Significant permission identification for machine-learning-based android malware detection. *IEEE Transactions on Industrial Informatics*, 14(7):3216–3225.
- Sun, L., Li, Z., Yan, Q., Srisa-an, W., and Pan, Y. (2016). Sigpid: significant permission identification for android malware detection. In *2016 11th international conference on malicious and unwanted software (MALWARE)*, pages 1–8. IEEE.
- Wu, Q., Zhu, X., and Liu, B. (2021). A survey of android malware static detection technology based on machine learning. *Mobile Information Systems*, 2021.
- Yildiz, O. and Doğru, I. A. (2019). Permission-based android malware detection system using feature selection with genetic algorithm. *International Journal of Software Engineering and Knowledge Engineering*, 29(02):245–262.