

Uma Análise de Métodos de Seleção de Características aplicados à Detecção de Malwares Android

XXII SIMPÓSIO BRASILEIRO DE SEGURANÇA DA INFORMAÇÃO E DE SISTEMAS COMPUTACIONAIS



UFAM

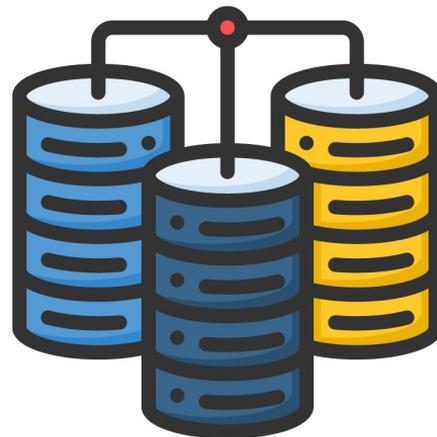


Taina Soares, Diego Kreutz,
Vanderson Rocha, Estevão Costa,
Luiza Leão, Jonas Pontes, Joner Assolin,
Gustavo Rodrigues, Eduardo Feitosa

DETECÇÃO DE MALWARES ANDROID E SELEÇÃO DE CARACTERÍSTICAS

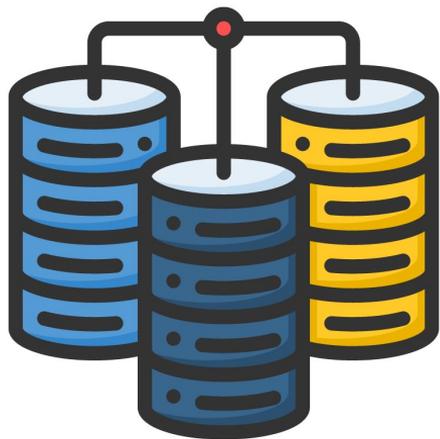


modelos de
aprendizado de máquina



datasets de alta
dimensionalidade

DETECÇÃO DE MALWARES ANDROID E SELEÇÃO DE CARACTERÍSTICAS



datasets de alta dimensionalidade

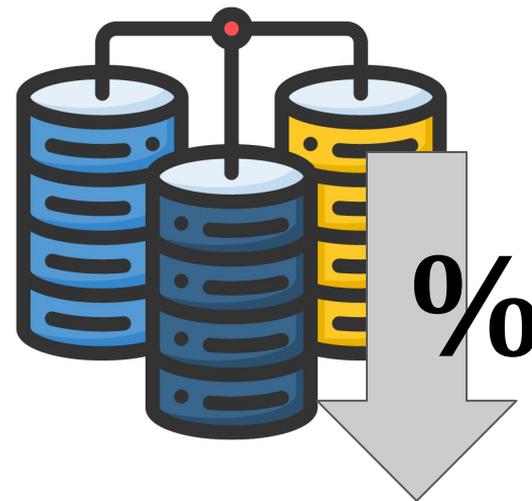
Exemplos:

- DefenseDroid > 11k features
- Drebin > 1m features e amostras

DETECÇÃO DE MALWARES ANDROID E SELEÇÃO DE CARACTERÍSTICAS



modelos de
aprendizado de máquina

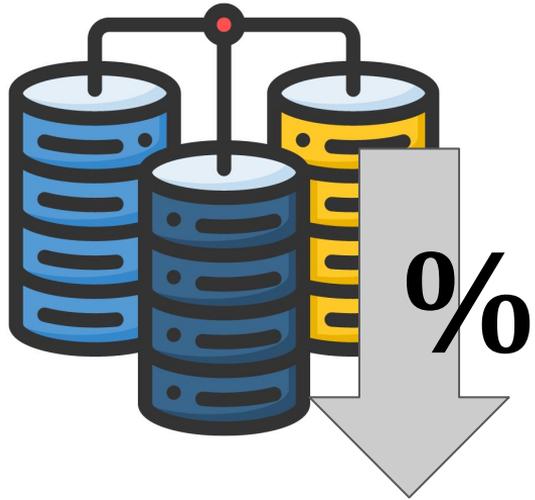


datasets reduzidos

MÉTODOS DE SELEÇÃO DE CARACTERÍSTICAS

- Específicos para algum tipo de característica
 - SigAPI - chamadas de API
 - SigPID - permissões
- Redução de tempo computacional
- Modelos com maior capacidade de generalização

MÉTODOS DE SELEÇÃO DE CARACTERÍSTICAS - REDUÇÃO



Exemplo: SigAPI ~75% de redução

- Drebin (subconjunto)
 - 73 API Calls
- MalGenome (subconjunto)
 - 69 API Calls

ROTEIRO

1. Trabalhos relacionados
2. Motivação
3. Métodos selecionados e implementados
4. Metodologia
5. Resultados e discussão
6. Considerações finais
7. Trabalhos futuros

TRABALHOS RELACIONADOS

- Comparação de técnicas clássicas de seleção estatística
 - Redução
 - Qualidade de detecção

TRABALHOS RELACIONADOS

- Comparação de técnicas clássicas de seleção estatística

- Redução
- Qualidade de detecção

Exemplos:

- chi-quadrado
- ganho de informação
- regras de associação

TRABALHOS RELACIONADOS

- Comparação de técnicas clássicas de seleção estatística
 - Redução
 - Qualidade de detecção

TRABALHOS RELACIONADOS

- Propostas de métodos elaborados de seleção
 - SigPID
 - SigAPI
 - ALR
 - RFG

EXEMPLO: FUNCIONAMENTO DO SIGPID

- Seleção de dados multinível
- 3 níveis de seleção:
 - PRNR (classificação de permissão com taxa negativa)
 - SPR (classificação de permissão baseada em suporte)
 - PMAR (mineração de permissões com regras de associação)
- Dataset reduzido (características selecionadas)

EXEMPLO: FUNCIONAMENTO DO SIGPID

- Drebin_215 (113 permissões)
 - PRNR (108)
 - SPR (30)
 - PMAR (27)

MOTIVAÇÃO DO ESTUDO

- Comparação de métodos sofisticados
- Métodos originalmente avaliados para datasets específicos

MOTIVAÇÃO DO ESTUDO

Método	Dataset
SigPID	Próprio (*)
SigAPI	Próprio (*)
RFG	Próprio (*)
ALR	Android Malware Dataset

MOTIVAÇÃO DO ESTUDO

Método	Dataset
SigPID	Próprio (*)
SigAPI	Próprio (*)
RFG	Próprio (*)
ALR	Android Malware Dataset

- Drebin
- Google Play
- APKPure

MOTIVAÇÃO DO ESTUDO

Qual o comportamento para diferentes datasets?

MOTIVAÇÃO DO ESTUDO

Qual o comportamento para diferentes datasets?

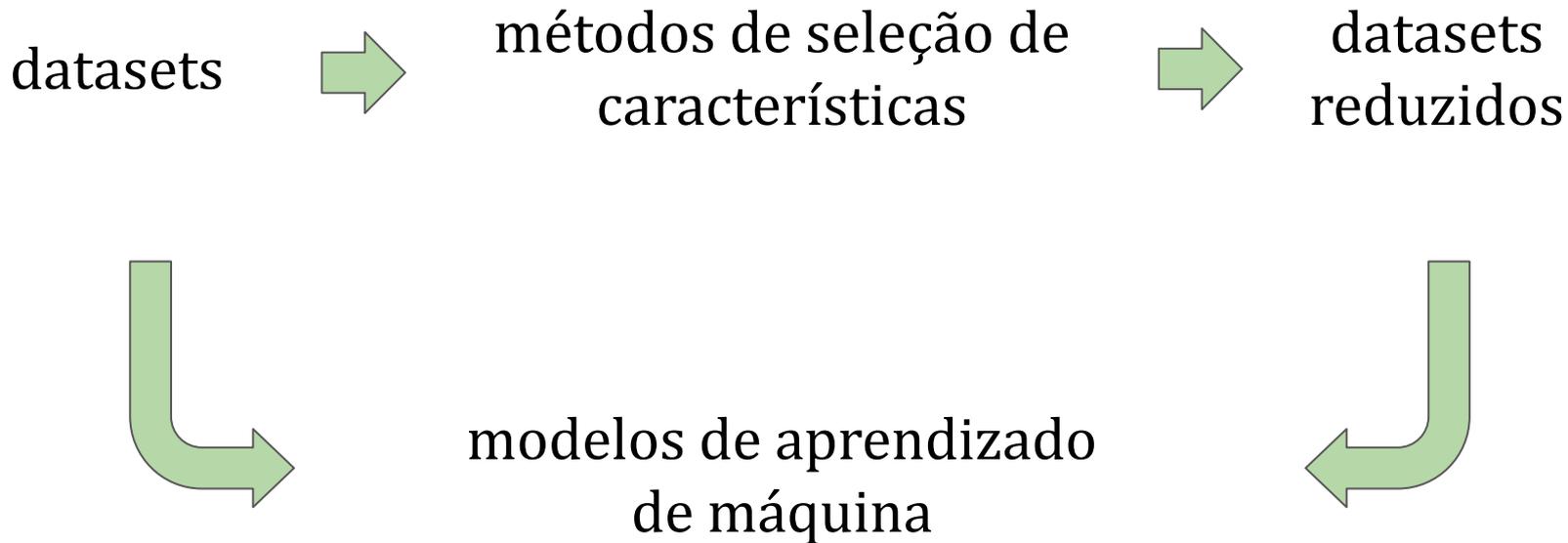
- Dimensionalidade
- Esparsidade
- Balanceamento (B:M)

MÉTODOS SELECIONADOS E IMPLEMENTADOS

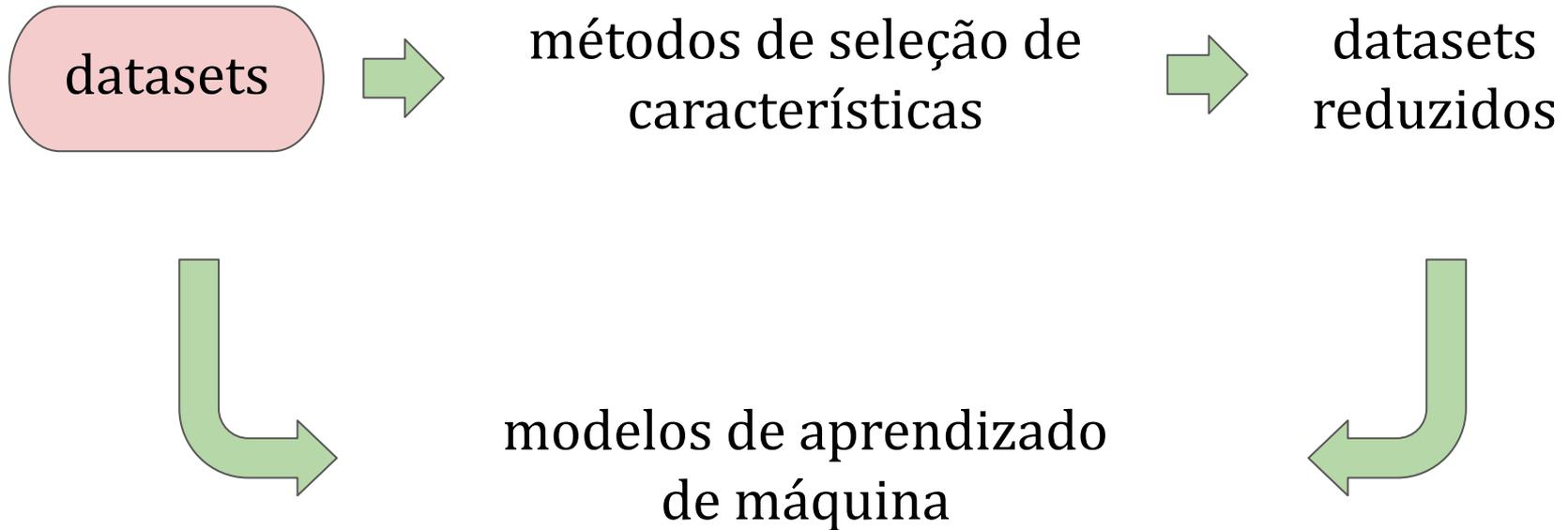
- Reprodutibilidade
- Bons resultados
- Complexidade

Método	Característica
SigPID	<i>Permissões</i>
ALR	<i>Permissões</i>
SigAPI	<i>Chamadas de API</i>
RFG	<i>Chamadas de API</i>

METODOLOGIA



METODOLOGIA - DATASETS



METODOLOGIA - DATASETS

Nome	Número de amostras	B : M
<i>md46k</i>	46.670	15,4 : 1
<i>androcrawl</i>	96.732	8,5 : 1
<i>drebin_215</i>	15.031	1,7 : 1

METODOLOGIA - DATASETS

datasets

- datasets de chamadas de API
- datasets de permissões

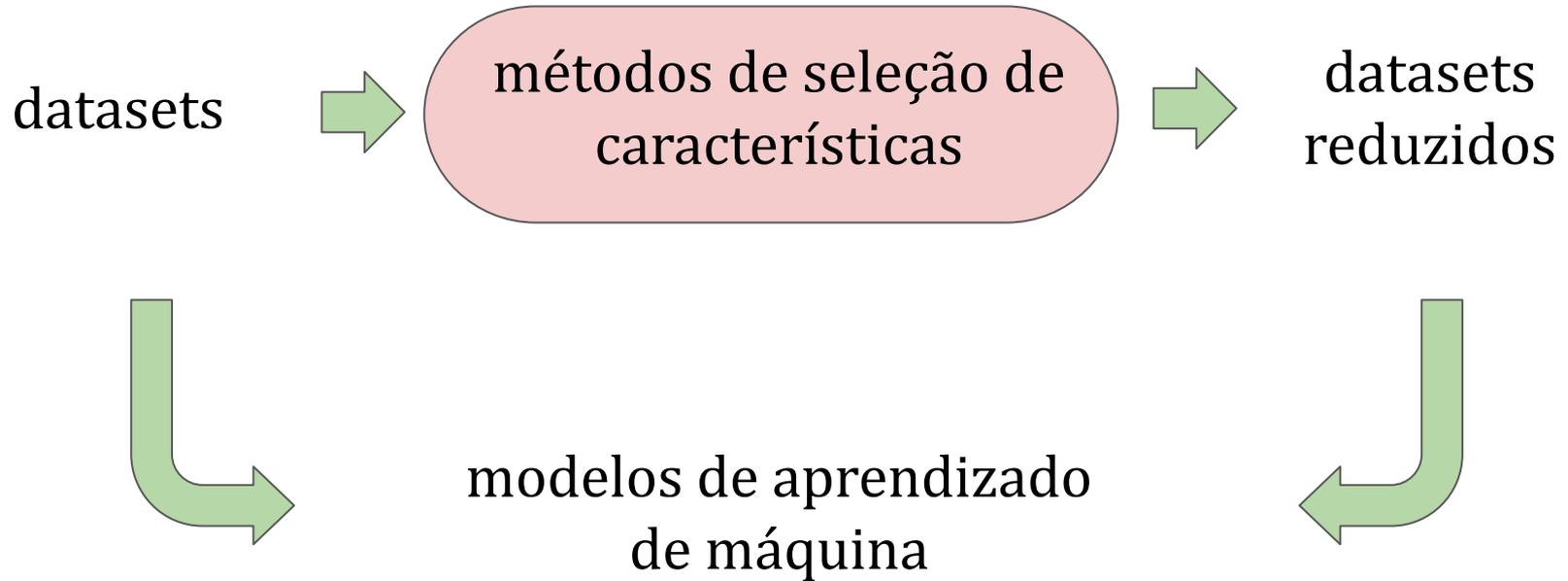
METODOLOGIA - DATASETS DE CHAMADAS DE API

Nome	Número de características
<i>md46k</i>	1.524
<i>androcrawl</i>	24
<i>drebin_215</i>	73

METODOLOGIA - DATASETS DE PERMISSÕES

Nome	Número de características
<i>md46k</i>	1.316
<i>androcrawl</i>	49
<i>drebin_215</i>	113

METODOLOGIA - MÉTODOS



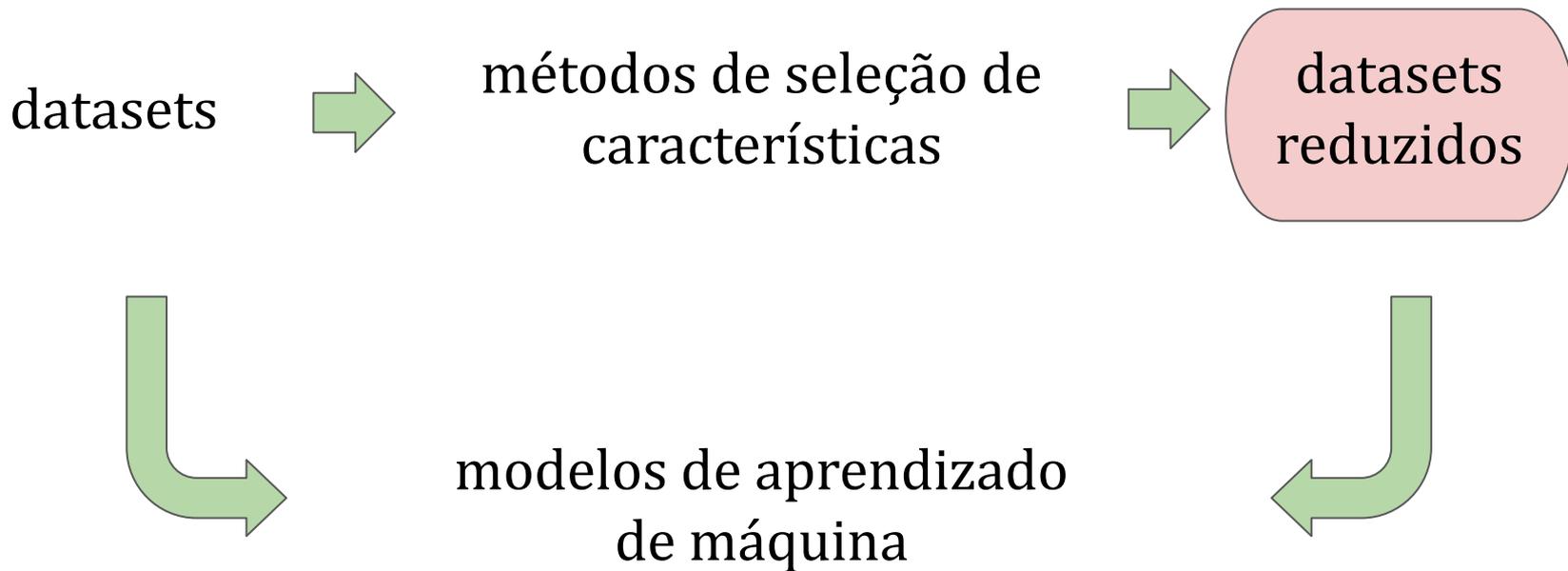
METODOLOGIA - MÉTODOS

métodos de seleção de características

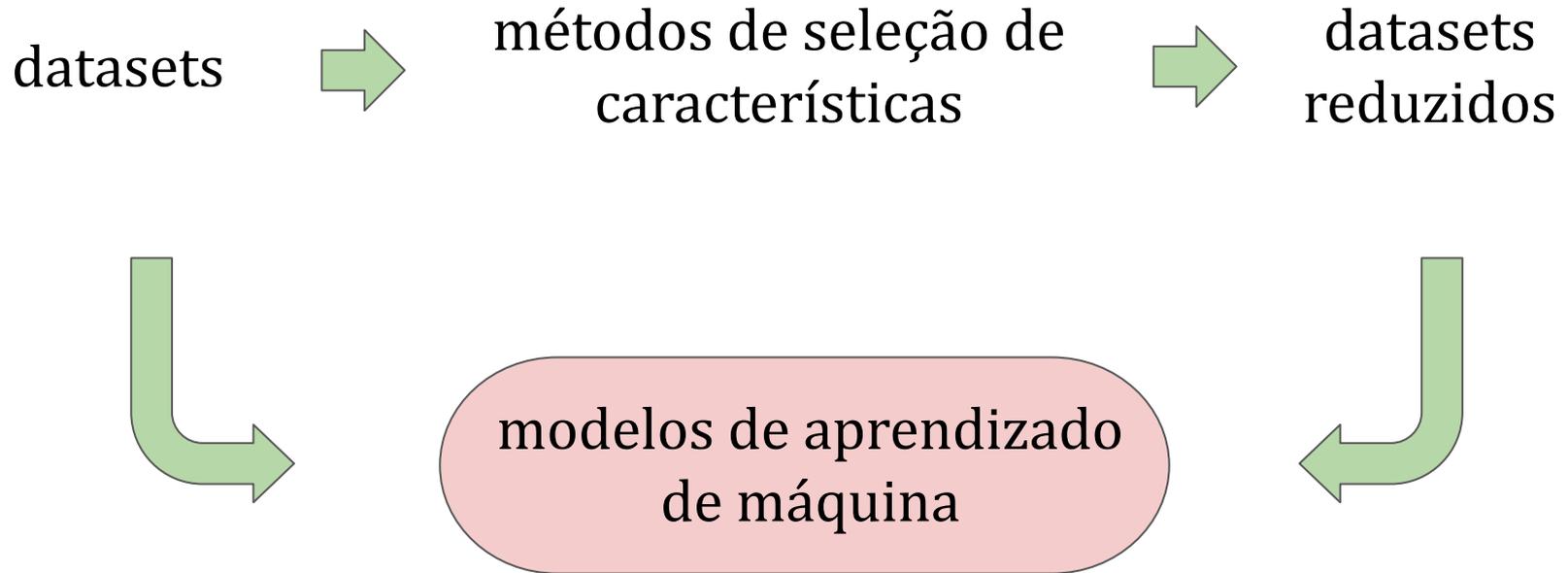
datasets de chamadas de API → SigAPI e RFG

datasets de permissões → SigPID e ALR

METODOLOGIA - DATASETS REDUZIDOS



METODOLOGIA - MODELOS



METODOLOGIA - MODELOS

modelos de aprendizado
de máquina

- Random Forest (RF)
- Support Vector Machine (SVM)

METODOLOGIA - MODELOS

RF:

- 100 árvores
 - Raiz quadrada da quantidade de características para procurar a melhor divisão entre classes
- ☐ Valores padrões da scikit learn

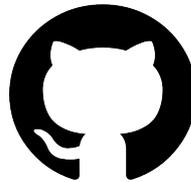
METODOLOGIA - MODELOS

SVM:

- Kernel Radial Basis Function (RBF) e heurística de redução para diminuir tempo de treino

- Valores padrões da scikit learn

REPOSITÓRIO FEATURE SELECTION



RESULTADOS - SELEÇÃO DE CHAMADAS DE API (SIGAPI E RFG)

Dataset	Redução (%)		Tempo (s)		Desempenho (ROC_AUC)		
	SigAPI	RFG	SigAPI	RFG	SigAPI	RFG	Sem redução
<i>drebin_215</i>	79,45	71,23	540	203	94,40	91,71	97,55
<i>md46k</i>	99,15	8,07	180803	4322	70,23	79,95	79,95
<i>androcrawl</i>	95,83	12,5	314	391	81,56	90,94	91,18

RESULTADOS - SELEÇÃO DE CHAMADAS DE API (SIGAPI E RFG)

Dataset	Redução (%)		Tempo (s)		Desempenho (ROC_AUC)		
	SigAPI	RFG	SigAPI	RFG	SigAPI	RFG	Sem redução
<i>drebin_215</i>	79,45	71,23	540	203	94,40	91,71	97,55
<i>md46k</i>	99,15	8,07	180803	4322	70,23	79,95	79,95
<i>androcrawl</i>	95,83	12,5	314	391	81,56	90,94	91,18

Mais demorado -> utilização de sete técnicas

RESULTADOS - SELEÇÃO DE CHAMADAS DE API (SIGAPI E RFG)

Dataset	Redução (%)		Tempo (s)		Desempenho (ROC_AUC)		
	SigAPI	RFG	SigAPI	RFG	SigAPI	RFG	Sem redução
<i>drebin_215</i>	79,45	71,23	540	203	94,40	91,71	97,55
<i>md46k</i>	99,15	8,07	180803	4322	70,23	79,95	79,95
<i>androcrawl</i>	95,83	12,5	314	391	81,56	90,94	91,18

Redução maior em todos os casos

RESULTADOS - SELEÇÃO DE CHAMADAS DE API (SIGAPI E RFG)

Dataset	Redução (%)		Tempo (s)		Desempenho (ROC_AUC)		
	SigAPI	RFG	SigAPI	RFG	SigAPI	RFG	Sem redução
<i>drebin_215</i>	79,45	71,23	540	203	94,40	91,71	97,55
<i>md46k</i>	99,15	8,07	180803	4322	70,23	79,95	79,95
<i>androcrawl</i>	95,83	12,5	314	391	81,56	90,94	91,18

Melhor resultado para *drebin_215_api_calls*

-> mais balanceado e conjunto reduzido de features

RESULTADOS - SELEÇÃO DE CHAMADAS DE API (SIGAPI E RFG)

Dataset	Redução (%)		Tempo (s)		Desempenho (ROC_AUC)		
	SigAPI	RFG	SigAPI	RFG	SigAPI	RFG	Sem redução
<i>drebin_215</i>	79,45	71,23	540	203	94,40	91,71	97,55
<i>md46k</i>	99,15	8,07	180803	4322	70,23	79,95	79,95
<i>androcrawl</i>	95,83	12,5	314	391	81,56	90,94	91,18

Grande variação

RESULTADOS - SELEÇÃO DE CHAMADAS DE API (SIGAPI E RFG)

Dataset	Redução (%)		Tempo (s)		Desempenho (ROC_AUC)		
	SigAPI	RFG	SigAPI	RFG	SigAPI	RFG	Sem redução
<i>drebin_215</i>	79,45	71,23	540	203	94,40	91,71	97,55
<i>md46k</i>	99,15	8,07	180803	4322	70,23	79,95	79,95
<i>androcrawl</i>	95,83	12,5	314	391	81,56	90,94	91,18

Melhores resultados em geral -> estabilidade

RESULTADOS - SELEÇÃO DE PERMISSÕES (SIGPID E ALR)

Dataset	Redução (%)		Tempo (s)		Desempenho (ROC_AUC)		
	SigPID	ALR	SigPID	ALR	SigPID	ALR	Sem redução
<i>drebin_215</i>	71,68	69,03	48	2	94,05	89,18	95,86
<i>md46k</i>	99,77	46,05	7106	746	50,00	52,81	62,06
<i>androcrawl</i>	87,76	97,96	176	3	50,00	50,00	49,84

RESULTADOS - SELEÇÃO DE PERMISSÕES (SIGPID E ALR)

Dataset	Redução (%)		Tempo (s)		Desempenho (ROC_AUC)		
	SigPID	ALR	SigPID	ALR	SigPID	ALR	Sem redução
<i>drebin_215</i>	71,68	69,03	48	2	94,05	89,18	95,86
<i>md46k</i>	99,77	46,05	7106	746	50,00	52,81	62,06
<i>androcrawl</i>	87,76	97,96	176	3	50,00	50,00	49,84

Mais demorado

RESULTADOS - SELEÇÃO DE PERMISSÕES (SIGPID E ALR)

Dataset	Redução (%)		Tempo (s)		Desempenho (ROC_AUC)		
	SigPID	ALR	SigPID	ALR	SigPID	ALR	Sem redução
<i>drebin_215</i>	71,68	69,03	48	2	94,05	89,18	95,86
<i>md46k</i>	99,77	46,05	7106	746	50,00	52,81	62,06
<i>androcrawl</i>	87,76	97,96	176	3	50,00	50,00	49,84

Ótimo resultado para *drebin_215_permissions*
-> leva em consideração balanceamento

RESULTADOS - SELEÇÃO DE PERMISSÕES (SIGPID E ALR)

Dataset	Redução (%)		Tempo (s)		Desempenho (ROC_AUC)		
	SigPID	ALR	SigPID	ALR	SigPID	ALR	Sem redução
<i>drebin_215</i>	71,68	69,03	48	2	94,05	89,18	95,86
<i>md46k</i>	99,77	46,05	7106	746	50,00	52,81	62,06
<i>androcrawl</i>	87,76	97,96	176	3	50,00	50,00	49,84

Reduções pouco seletivas em datasets
esparsos e desbalanceados

CONSIDERAÇÕES FINAIS

- Relação dos métodos com particularidades dos datasets
 - Resultados oscilam para diferentes datasets
 - SigAPI, SigPID e ALR -> datasets balanceados
 - RFG mostrou mais estabilidade
 - Etapa robusta de validação

TRABALHOS FUTUROS

- Avaliação dos métodos com uma gama ainda maior e mais diversificada de datasets
 - conjuntos de dados maiores e sem pré-processamento

TRABALHOS FUTUROS

- Avaliação dos métodos com uma gama ainda maior e mais diversificada de datasets
 - conjuntos de dados maiores e sem pré-processamento
- Avaliação de outros métodos de seleção de características

TRABALHOS FUTUROS

- Avaliação dos métodos com uma gama ainda maior e mais diversificada de datasets
 - conjuntos de dados maiores e sem pré-processamento
- Avaliação de outros métodos de seleção de características
- Identificação de parâmetros de otimização que podem impactar os métodos de seleção

TRABALHOS FUTUROS

- Avaliação dos métodos com uma gama ainda maior e mais diversificada de datasets
 - conjuntos de dados maiores e sem pré-processamento
- Avaliação de outros métodos de seleção de características
- Identificação de parâmetros de otimização que podem impactar os métodos de seleção
- Avaliação de resultados a partir da combinação de métodos

TRABALHOS FUTUROS

- Avaliação dos métodos com uma gama ainda maior e mais diversificada de datasets
 - conjuntos de dados maiores e sem pré-processamento
- Avaliação de outros métodos de seleção de características
- Identificação de parâmetros de otimização que podem impactar os métodos de seleção
- Avaliação de resultados a partir da combinação de métodos

OBRIGADA PELA ATENÇÃO!



UFAM



Uma Análise de Métodos de Seleção de Características aplicados à Detecção de Malwares Android



tainasoaes.aluno@unipampa.edu.br

