

AVALIAÇÃO DE MÉTODOS DE CLASSIFICAÇÃO BASEADOS EM REGRAS DE ASSOCIAÇÃO PARA DETECÇÃO DE MALWARES ANDROID



UFAM



Vanderson Rocha
Diego Kreutz
Jonas Pontes
Eduardo Feitosa

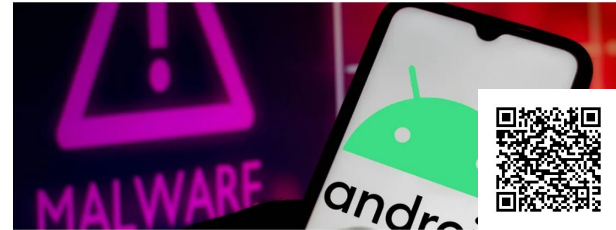
Universidade Federal do Amazonas (UFAM)
Universidade Federal do Pampa (UNIPAMPA)

DETECÇÃO DE MALWARES



BRATA: malware para Android rouba seus dados e reseta o celular

25/01/2022 às 17:50 • 1 min de leitura

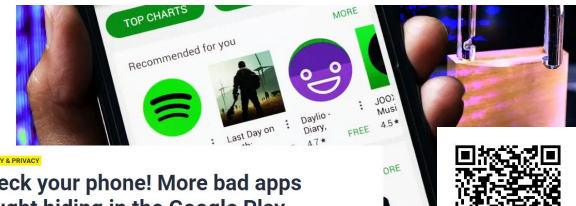


Crypto malware in patched wallets targeting Android and iOS devices

ESET Research uncovers a sophisticated scheme that distributes trojanized Android and iOS cryptocurrency wallets

Lukas Stefanko

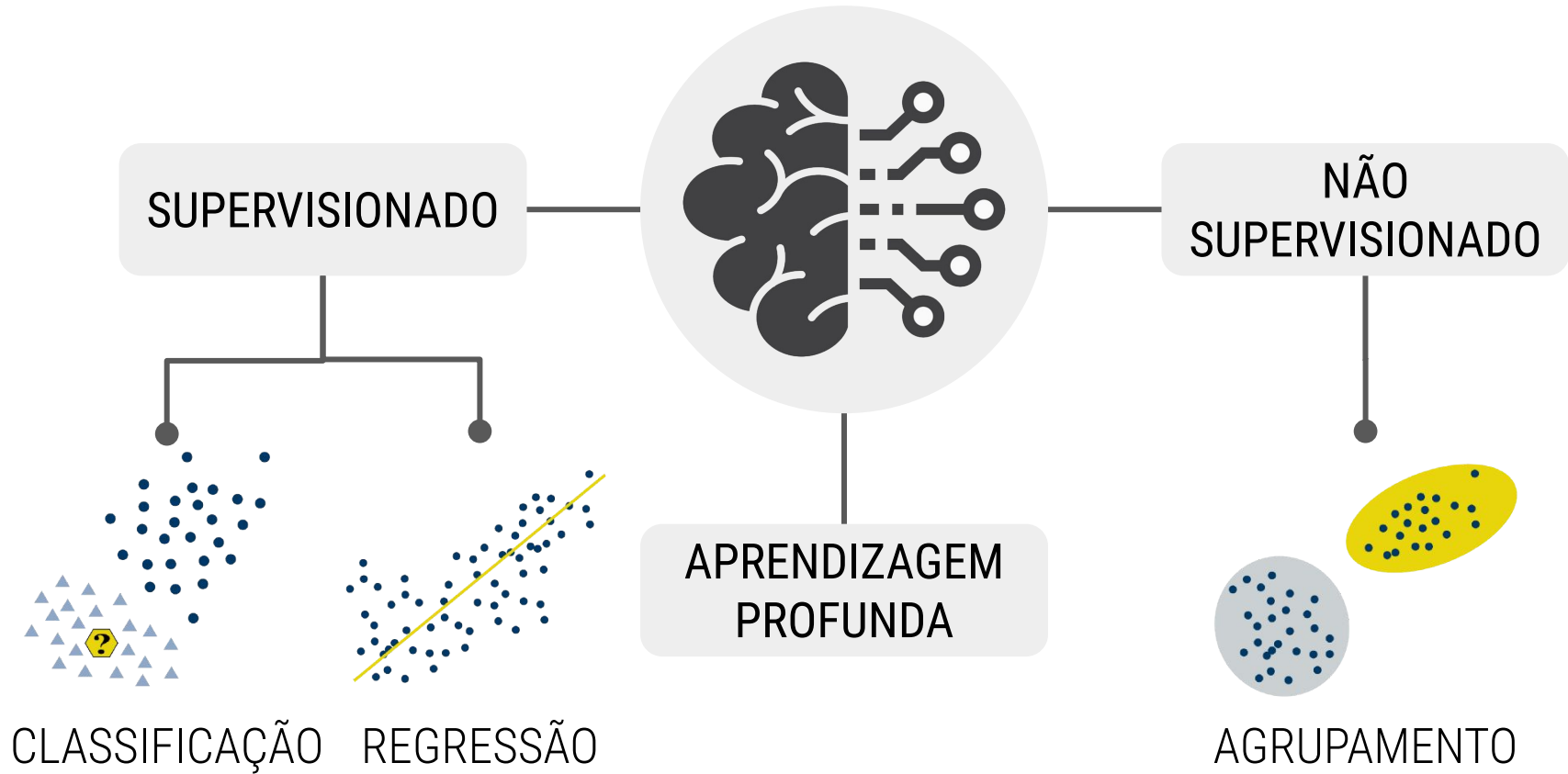
24 Mar 2022 - 01:30PM



Check your phone! More bad apps caught hiding in the Google Play Store

BY CHARLIE FRIPP KOMANDO • MARCH 21, 2022

APRENDIZADO DE MÁQUINA



CLASSIFICAÇÃO BASEADA EM REGRAS DE ASSOCIAÇÃO

- Gestão de Relacionamento
- Diagnóstico Médico
- Combate à Fraudes

OBJETIVO

Análise exploratória de diferentes métodos de classificação.

Algoritmos Clássicos	Método
Apriori	CBA
FP-Growth	CMAR
ECLAT	EQAR

ROTEIRO

- **Mineração de Regras de Associação**
- Método de Classificação EQAR
- Metodologia
- Resultados
- Considerações Finais

DEFINIÇÃO

Padrões, Correlações ou Associações

Se (Antecedente) \Rightarrow Então (Consequente)

Se $\{ACCESS_NETWORK_STATE, READ_PHONE_STATE, SEND_SMS\}$
 \Rightarrow Então Malicioso

PRINCIPAIS MÉTRICAS

- Suporte
- Confiança

Leite	Pão	Ovo
1	0	1
1	1	0
1	1	1
1	1	1
0	0	1

PRINCIPAIS MÉTRICAS

- Suporte

Frequência (Leite \cup Pão) / N = 3/5 = 0.6 (60%)

- Confiança

Leite	Pão	Ovo
1	0	1
1	1	0
1	1	1
1	1	1
0	0	1

PRINCIPAIS MÉTRICAS

- Suporte

Frequência (Leite \cup Pão) / N = 3/5 = 0.6 (60%)

- Confiança

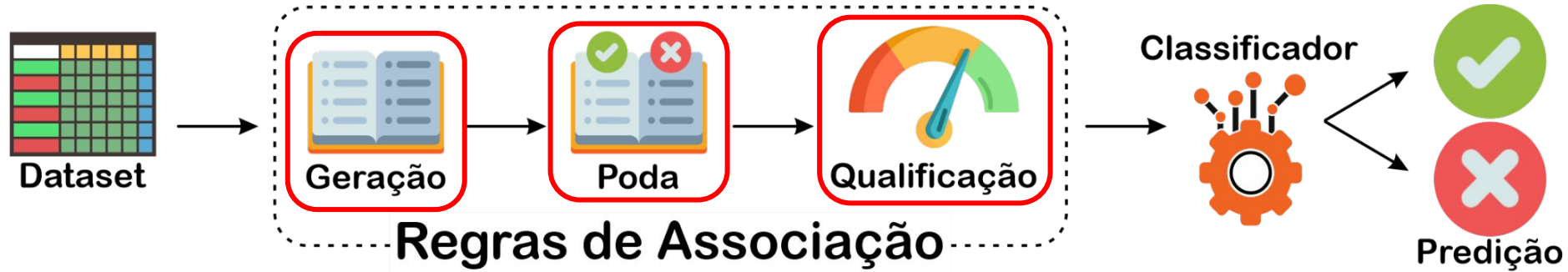
Frequência (Pão) / Frequência (Leite) = 3/4 = 0.75 (75%)

Leite	Pão	Ovo
1	0	1
1	1	0
1	1	1
1	1	1
0	0	1

ROTEIRO

- Mineração de Regras de Associação
- **Método de Classificação EQAR**
- Metodologia
- Resultados
- Considerações Finais

ETAPAS



GERAÇÃO

Algoritmo 1: Geração e Poda de Regras de Associação

Entrada:

D : *dataset* de treino

minSup: suporte mínimo

minConf: confiança mínima

```
1 datasetmalignas ←  $t \in D \wedge t[class] = 1$ ; // amostras malignas  
2 datasetbenignas ←  $t \in D \wedge t[class] = 0$ ; // amostras benignas
```

```
3 regrasmalignas ← gerar_regras(datasetmalignas, minSup, minConf)  
4 regrasbenignas ← gerar_regras(datasetbenignas, minSup, minConf)
```

```
5 regrasdif ← regrasmalignas \ regrasbenignas
```

```
6 regrassub ←  $X \in regras_{dif} \wedge X \not\subseteq Y, \forall Y \in regras$ 
```

```
7 regras ←  $X \in regras_{sub} \wedge X \not\subseteq Y, \forall Y$ 
```

```
8 retorna regras
```



$X \Rightarrow Y$

PODA

Algoritmo 1: Geração e Poda de Regras de Associação

Entrada:

D : *dataset* de treino

minSup: suporte mínimo

minConf: confiança mínima

```
1 datasetmalignas ←  $t \in D \wedge t[class] = 1$ ; //  $C_1 \cup C_2 \cup C_3 \Rightarrow$  malignas
2 datasetbenignas ←  $t \in D \wedge t[class] = 0$ ; // benignas
3 regrasmalignas ← gerar_regras(datasetmalignas, minSup, minConf)
4 regrasbenignas ← gerar_regras(datasetbenignas, minSup, minConf)
5 regrasdif ← regrasmalignas \ regrasbenignas
6 regrassub ←  $X \in regras_{dif} \wedge X \not\subseteq Y, \forall Y \in regras_{benignas}$ 
7 regras ←  $X \in regras_{sub} \wedge X \not\subseteq Y, \forall Y \in regras_{sub}$ 
8 retorna regras
```

$C_1 \cup C_8 \Rightarrow Classe$

QUALIDADE DE REGRAS

- Cobertura
 - Maximizar Amostras Positivas
- Consistência
 - Minimizar Amostras Negativas

	Malignas	Benignas
Amostras Cobertas Pela Regra	p	n
Total de Amostras	P	N

QUALIFICAÇÃO DE REGRAS

Algoritmo 2: Qualificação de Regras de Associação

Entrada:

D: *dataset* de treino

regras: conjunto de regras (Algoritmo 1)

métrica: função de métrica de qualidade a ser utilizada

```
1 P ← quantidade de malignas em D
2 N ← quantidade de benignas em D
3 regrasqualidade ← ∅; // lista (regra, qualidade)
4 para cada regra  $r \in$  regras faça
5   | p, n ← encontrar_cobertura(r)
6   | q ← métrica(p, n, P, N)
7   | regrasqualidade ← regrasqualidade ∪ (r, q)
8 retorna regrasqualidade
```


ROTEIRO

- Mineração de Regras de Associação
- Método de Classificação EQAR
- **Metodologia**
- Resultados
- Considerações Finais

MÉTODOS DE CLASSIFICAÇÃO

Método	Geração de Regras	Classificação das Regras	Poda	Predição
CBA	Apriori	Confiança, Suporte Geradas Primeiro	Cobertura (Método M1)	Máxima Probabilidade
CMAR	FP-Growth	Confiança, Suporte Cardinalidade	Cobertura (Método M1), Regras Redundantes	Múltiplas Regras
CPAR	PRM	Acurácia de Laplace	Seleção dos K Melhores	Múltiplas Regras
EQAR	ECLAT	Confiança, Suporte	Qualidade de Regras, Diferença de Conjuntos	Correspondência Exata

MÉTODOS DE CLASSIFICAÇÃO

Método	Geração de Regras	Classificação das Regras	Poda	Predição
CBA	Apriori	Confiança, Suporte Geradas Primeiro	Cobertura (Método M1)	Máxima Probabilidade
CMAR	FP-Growth	Confiança, Suporte Cardinalidade	Cobertura (Método M1), Regras Redundantes	Múltiplas Regras
CPAR	PRM	Acurácia de Laplace	Seleção dos K Melhores	Múltiplas Regras
EQAR	ECLAT	Confiança, Suporte	Qualidade de Regras, Diferença de Conjuntos	Correspondência Exata

MÉTODOS DE CLASSIFICAÇÃO

Método	Geração de Regras	Classificação das Regras	Poda	Predição
CBA	Apriori	Confiança, Suporte Geradas Primeiro	Cobertura (Método M1)	Máxima Probabilidade
CMAR	FP-Growth	Confiança, Suporte Cardinalidade	Cobertura (Método M1), Regras Redundantes	Múltiplas Regras
CPAR	PRM	Acurácia de Laplace	Seleção dos K Melhores	Múltiplas Regras
EQAR	ECLAT	Confiança, Suporte	Qualidade de Regras, Diferença de Conjuntos	Correspondência Exata

CONJUNTOS DE DADOS

Dataset	Características	Amostras		
		Malignas	Benignas	Total
KronoDroid Emulador	383	28745	35246	63991
KronoDroid Dispositivo Real	383	41382	36755	78137
DREBIN-215	215	5560	9476	15036

CONJUNTOS DE DADOS

Dataset	Características	Amostras		
		Malignas	Benignas	Total
KronoDroid Emulador	383	28745	35246	63991
KronoDroid Dispositivo Real	383	41382	36755	78137
DREBIN-215	215	5560	9476	15036

CONJUNTOS DE DADOS

Dataset	Características	Amostras		
		Malignas	Benignas	Total
KronoDroid Emulador	383	28745	35246	63991
KronoDroid Dispositivo Real	383	41382	36755	78137
DREBIN-215	215	5560	9476	15036

CONJUNTOS DE DADOS

Dataset	Características	Amostras		
		Malignas	Benignas	Total
KronoDroid Emulador	383	28745	35246	63991
KronoDroid Dispositivo Real	383	41382	36755	78137
DREBIN-215	215	5560	9476	15036

PARÂMETROS

- Tamanho Máximo da Regra: 5
- Confiança Mínima: 95%
- Validação
 - Cruzada
 - Estratificada
 - 5 Partições

PARÂMETROS

- RF
 - Árvores na Floresta: 100
 - Divisão Entre Classes: Raiz Quadrada
- SVM:
 - Kernel: Radial Basis Function (RBF)
 - Heurística de Redução

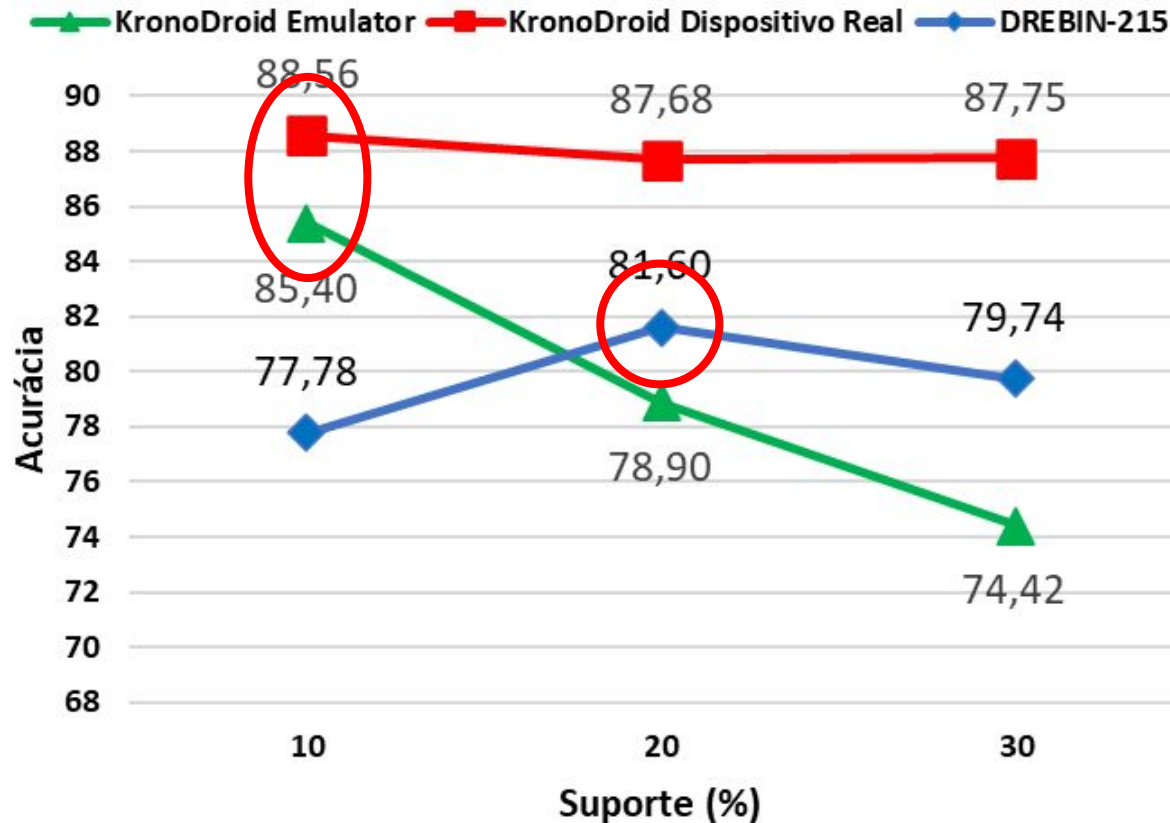
MÉTRICAS

- Acurácia
- Precisão
- Recall
- F1 Score
- MCC

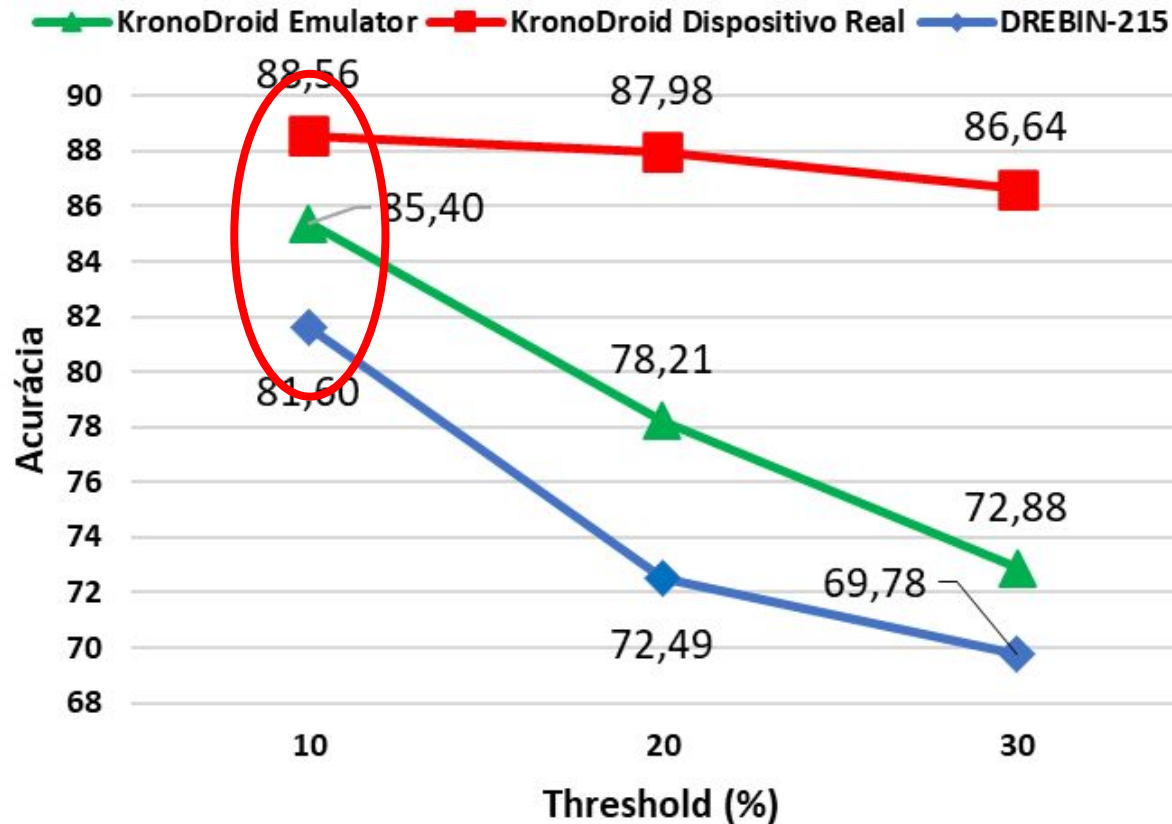
ROTEIRO

- Mineração de Regras de Associação
- Método de Classificação EQAR
- Metodologia
- **Resultados**
- Considerações Finais

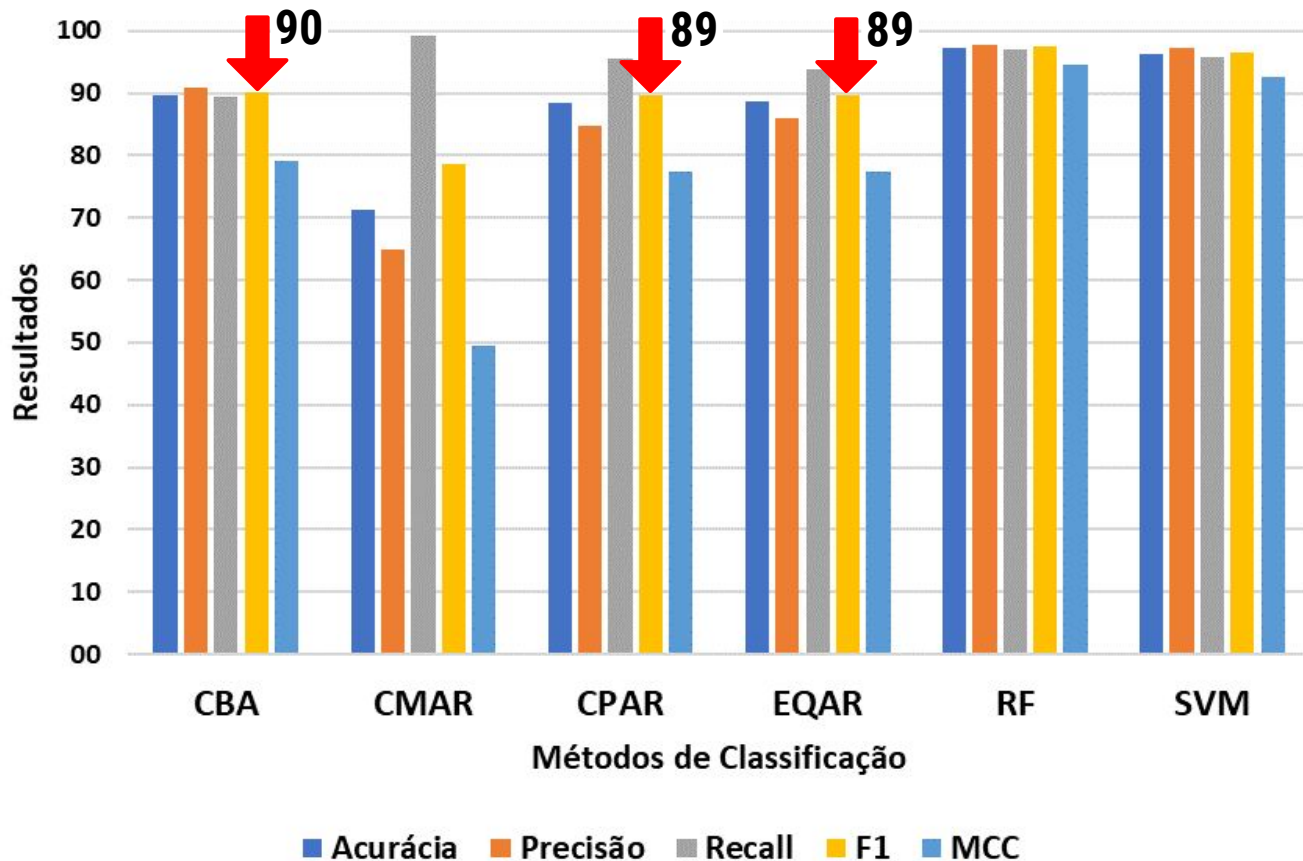
SUPOORTE



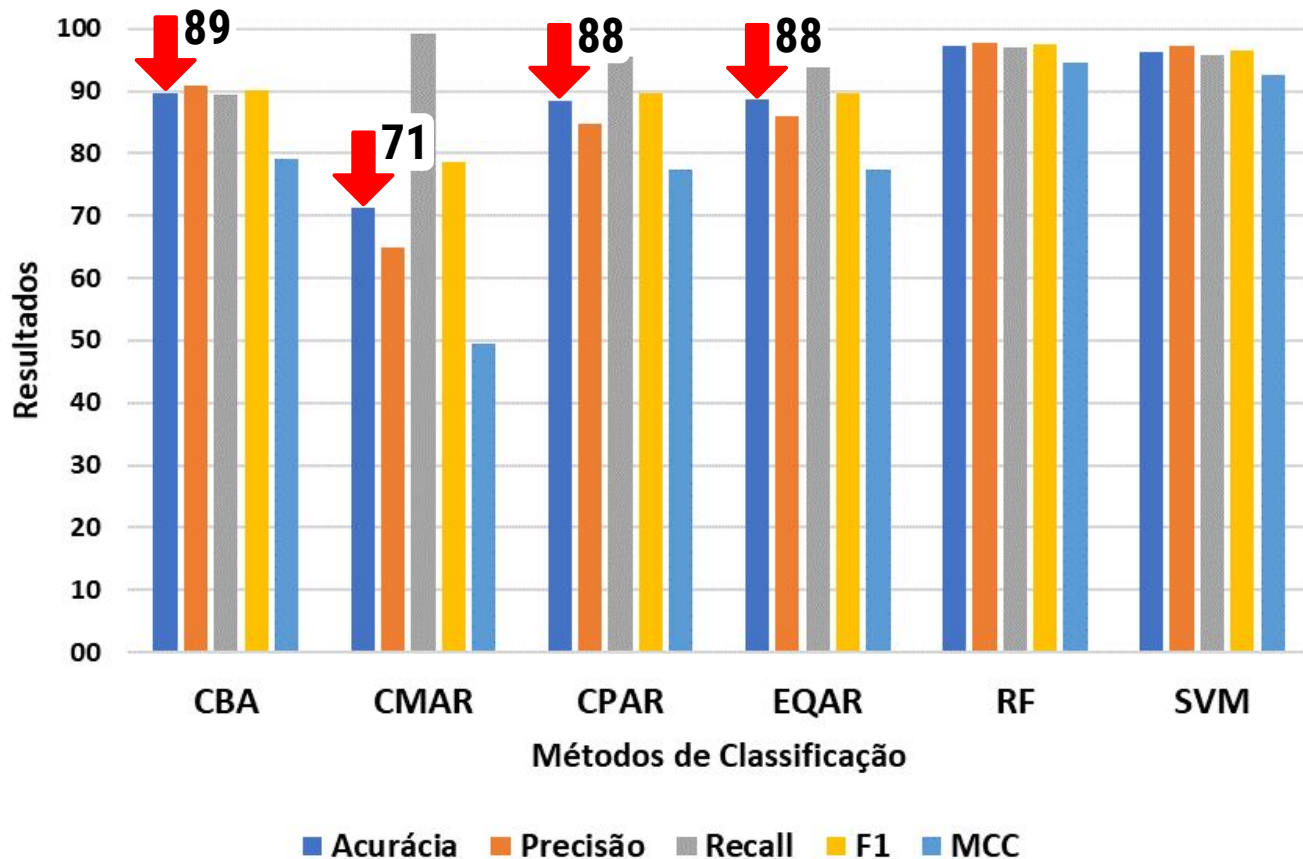
THRESHOLD



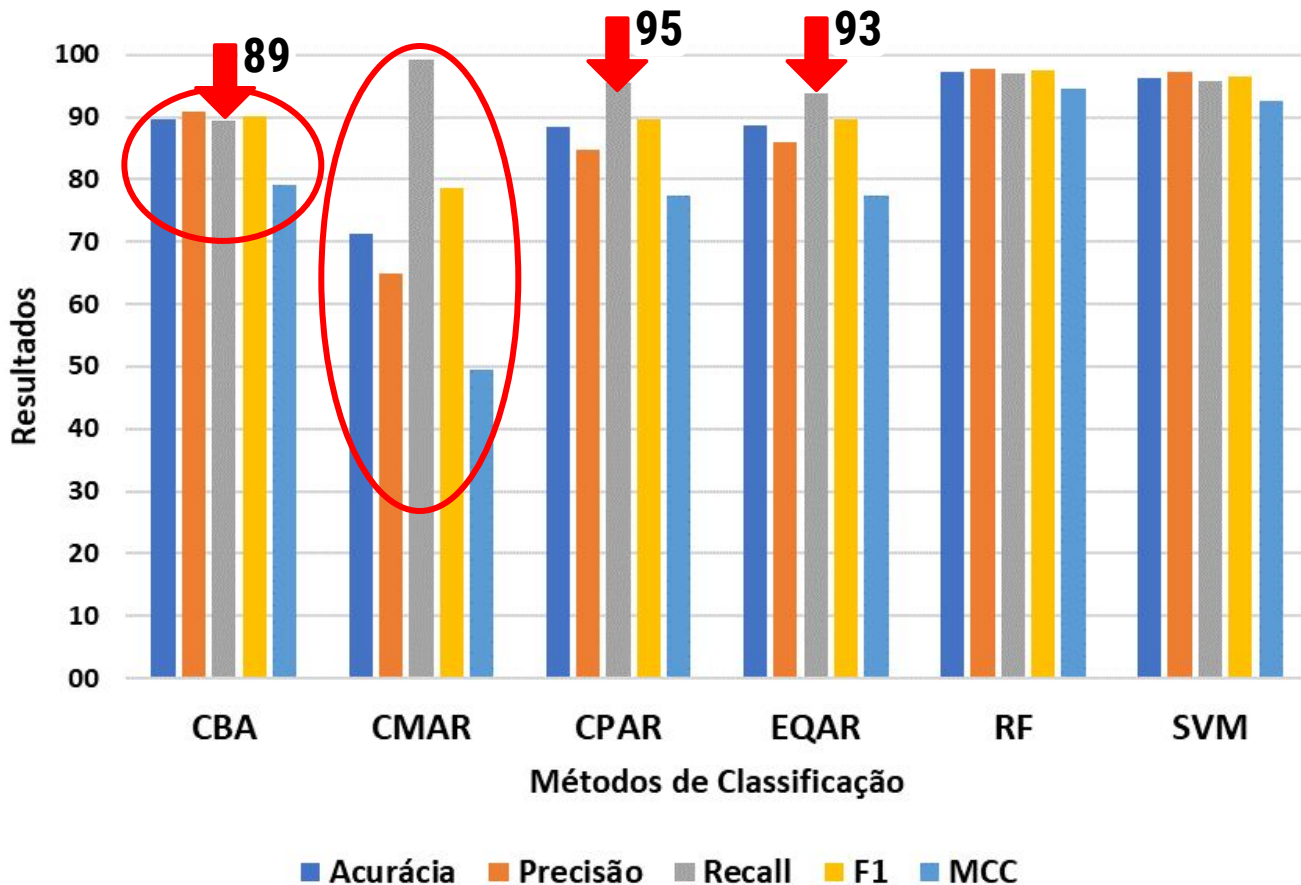
KRONODROID - DISPOSITIVO REAL



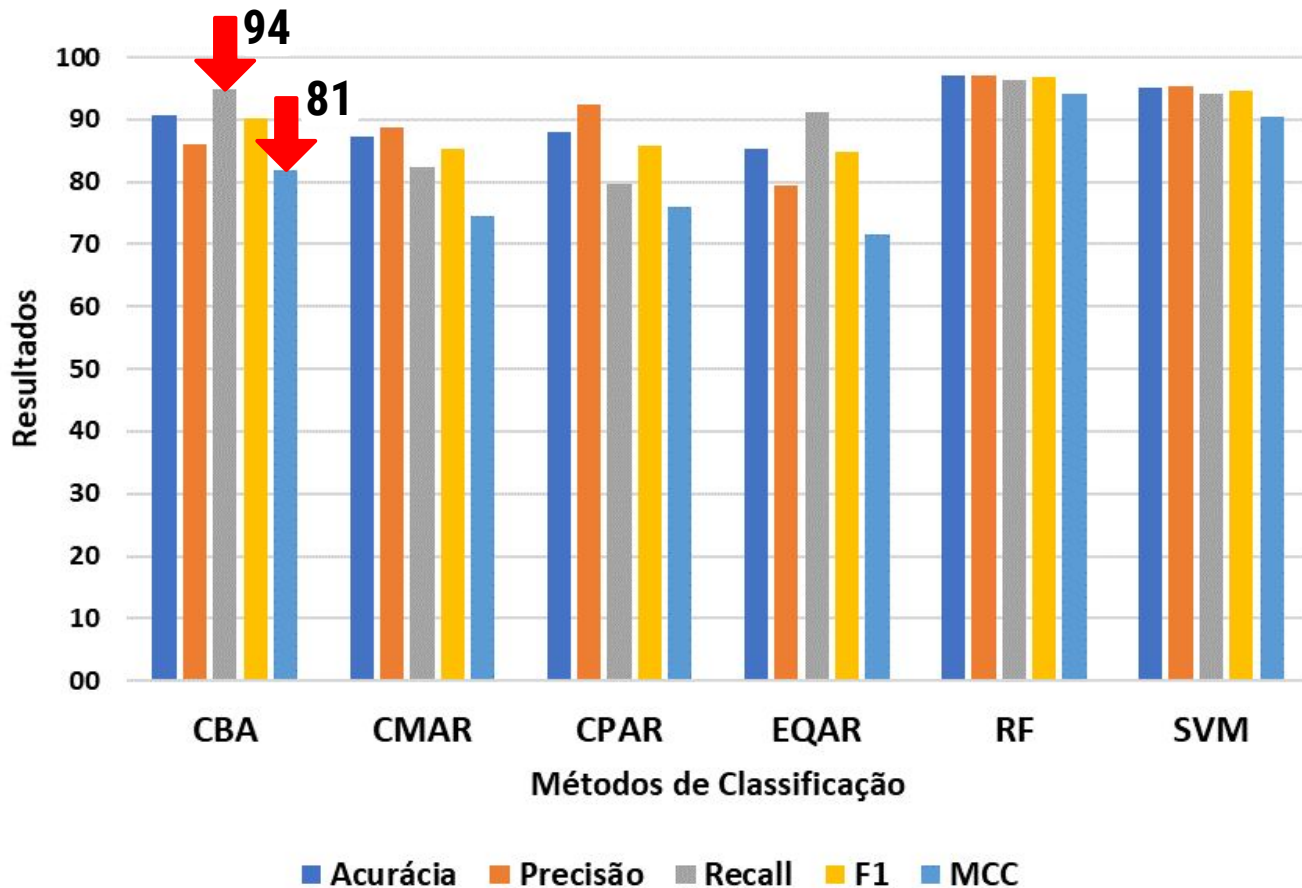
KRONODROID - DISPOSITIVO REAL



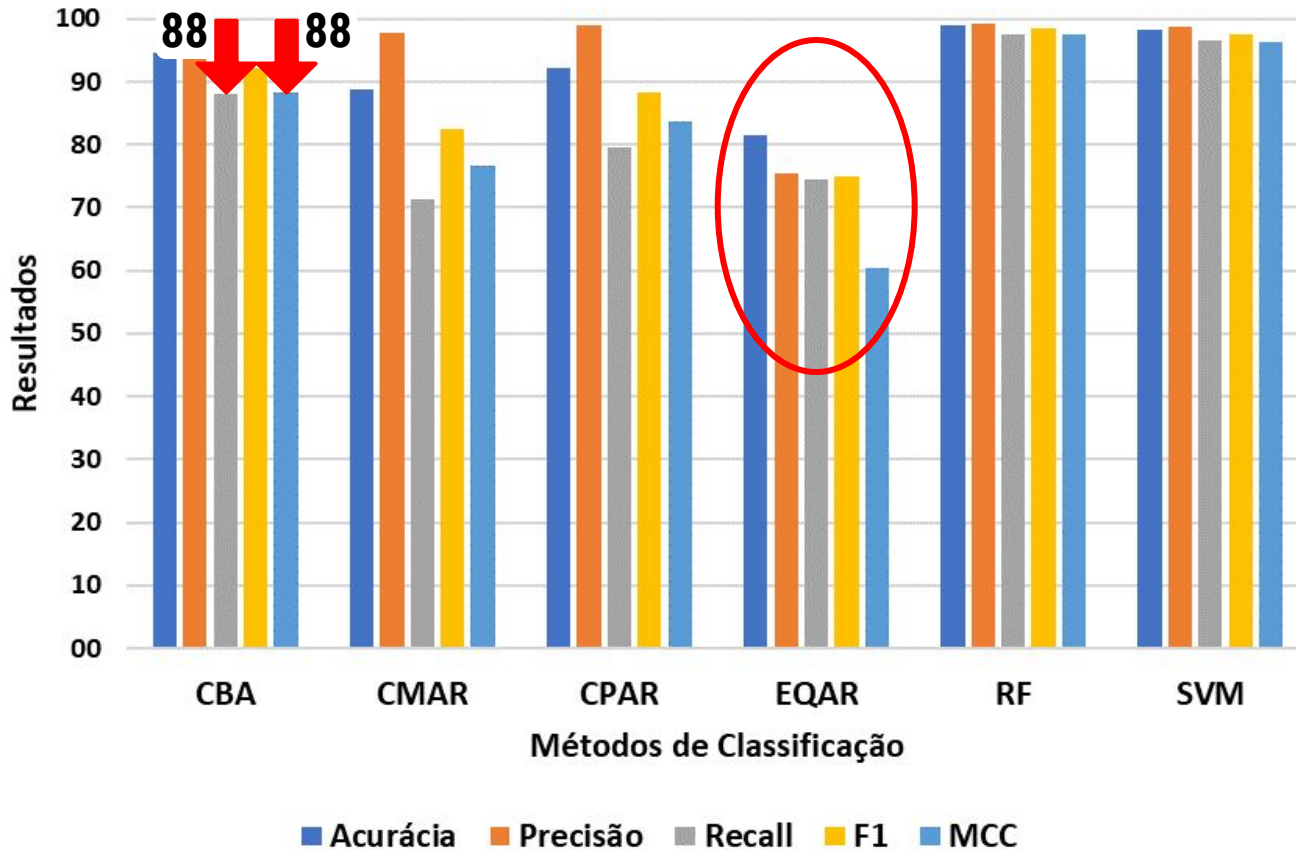
KRONODROID - DISPOSITIVO REAL



KRONODROID - EMULADOR



DREBIN-215



ROTEIRO

- Mineração de Regras de Associação
- Método de Classificação EQAR
- Metodologia
- Resultados
- **Considerações Finais**

TRABALHOS FUTUROS

- Expandir e Aprimorar o EQAR
- Reduzir FP e FN
- Avaliar
 - Novos Conjuntos de Dados
 - Outros Modelos de Aprendizagem de Máquina

OBRIGADO

Perguntas?

vanderson@ufam.edu.br
ppgi.ufam.edu.br