Beijing Forest Studio 北京理工大学信息系统及安全对抗实验中心



HinDroid: An Intelligent Android Malware Detection System Based on Structured Heterogeneous Information Network

7018在11月11日
Refection System Based on Structured
Heterogeneous Information Network

内容提要



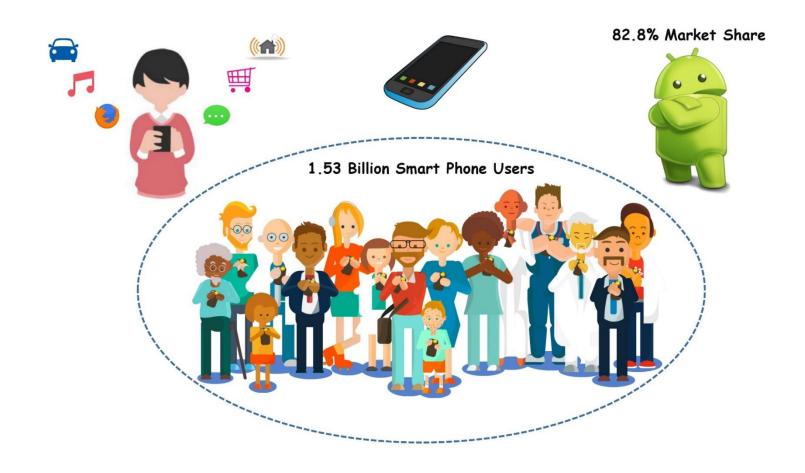
- 背景简介
- 基本概念
- 框架原理
- 实验分析
- 优劣分析
- 应用总结
- 参考文献



• 预期收获

- 1.了解一种最新恶意软件智能检测方法
- 2.理解异构信息网络、SVM、多核学习等基础知识
- 3. 学习信息网络在数据挖掘问题中的运用







































Send SMS message

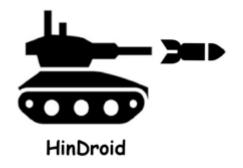




Download unwanted app

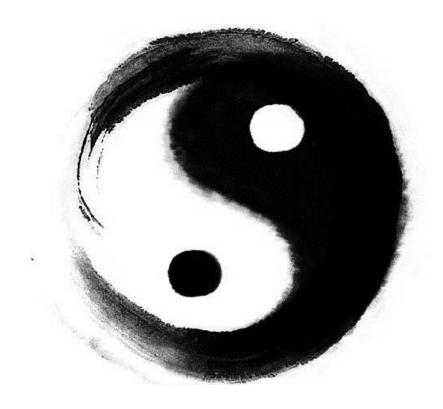


Abdulhayoglu, Melih, et al. "HinDroid: An Intelligent Android Malware Detection System Based on Structured Heterogeneous Information Network." *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* ACM, 2017:1507-1515.



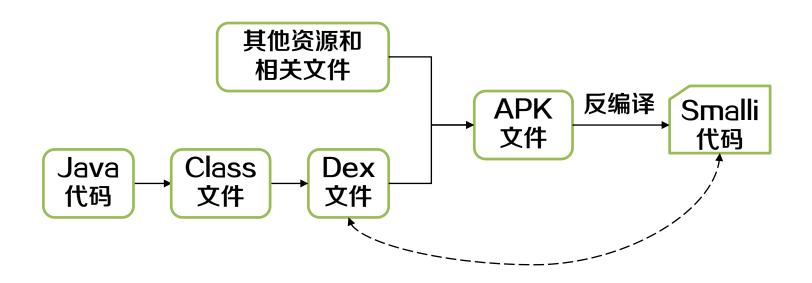








- Smali语言
 - Dalvik虚拟机的寄存器语言
 - Dex文件反编译后得到Smalli代码



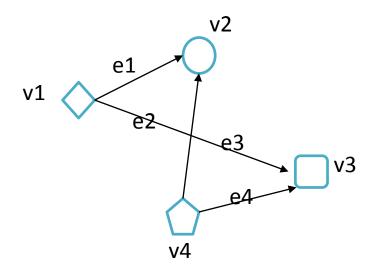


· Smalli代码实例

```
□ Java2Smali -- 作者: 隐心
                                                                                                                  \times
 文件 关于
                                                                   GO (F5)
                                                                     1 .class public LTest;
 1 public class Test {
                                                                     2 .super Ljava/lang/Object;
                                                                      3 .source "Test.java"
       public static void main(String[] args) {
           System.out.println("Hello World!");
                                                                       # direct methods
                                                                      .method public constructor <init>()V
                                                                           .registers 1
                                                                           .prologue
                                                                           .line 1
                                                                           invoke-direct {p0}, Ljava/lang/Object; -><init>(
                                                                           return-void
                                                                    14 .end method
                                                                    16 .method public static main([Ljava/lang/String;)V
                                                                           .registers 3
                                                                    18
                                                                           .prologue
                                                                    19
                                                                           .line 5
                                                                           sget-object v0, Ljava/lang/System; ->out:Ljava/i
                                                                           const-string v1, "Hello World!"
                                                                    23
                                                                    24
                                                                           invoke-virtual {v0, v1}, Ljava/io/PrintStream;-
                                                                    25
                                                                    26
                                                                           .line 6
                                                                           return-void
                                                                    28 .end method
Java代码
                                                                ..: Smali代码
```

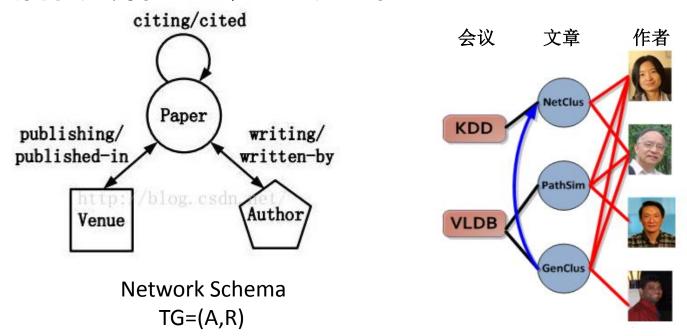


- 信息网络
 - 信息网络可以用一个有向图G = (V, E) 来表示





- 异构信息网络(HIN)
 - 如果|A| > 1 或者 |R| > 1 ,则该信息网络为异构信息网络,或简称为异构网络,否则为同构网络。



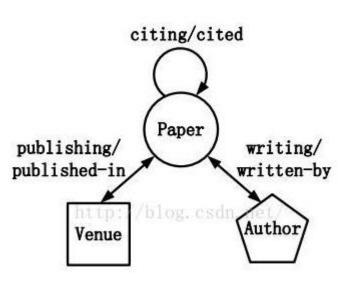
A= {会议、文章、作者}

R= { publishing \(\) published-in \(\) writing \(\) written-by \(\) citing \(\) cityed \(\)

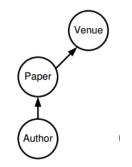


元路径

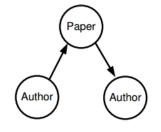
- 元路径P定义在网络模式TG=(A, R)之上
- 链接两类对象的一条路径



Network Schema TG=(A,R)



元路径: APV



元路径: APA



交換矩阵

- 元路径中实体邻接矩阵乘积

元路径:
$$P = (A_1 A_2 A_3 \cdots A_l)$$

交换矩阵:
$$M = W_{A_1A_2}W_{A_2A_3}\cdots W_{A_{l-1}A_l}$$

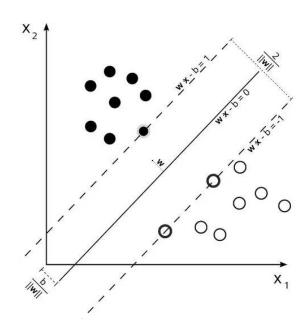
例如:
$$P=(AVA)$$
 \rightarrow $M=W_{AV}W_{VA}$

$$M = W_{AV} * W_{AV}^{T} = \begin{bmatrix} 5 & 120 & 4 & 5 & 0 \\ 120 & 2900 & 100 & 120 & 0 \\ 4 & 100 & 5 & 4 & 1 \\ 5 & 120 & 4 & 5 & 0 \\ 0 & 0 & 1 & 0 & 2 \end{bmatrix}$$



- SVM分类原理
 - 特征空间里寻找超平面,以最小的错分率把正负样本分开

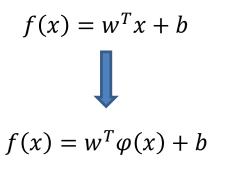
$$D = \{(x_1, y_1), (x_2, y_2), \dots (x_m, y_m)\}$$
$$y_i \in (-1, +1)$$
$$f(x) = w^T x + b$$

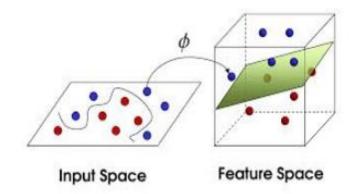




核技巧

- 将输入空间映射到高维特征空间



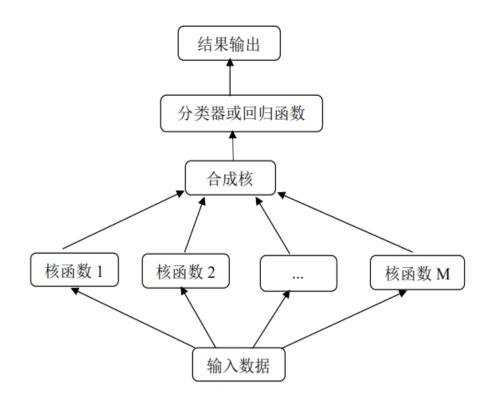


$$\max_{\alpha} \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i=1}^{m} \sum_{j=1}^{m} \alpha_i \alpha_j \ y_i \ y_j \varphi(x_j)^T \varphi(x_i)$$

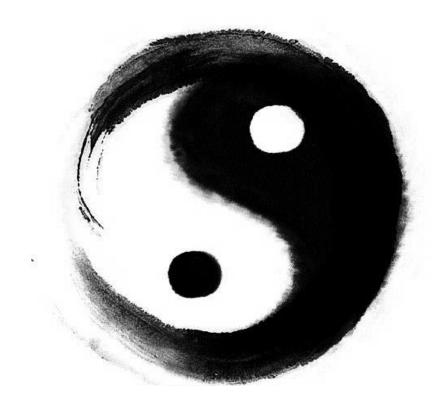
$$s.t. \sum_{i=1}^{m} \alpha_i y_i = 0$$
核函数!



- 多核学习
 - 针对不同的特征,采用不同核函数
 - 充分发挥了各个基本核的不同特征映射能力







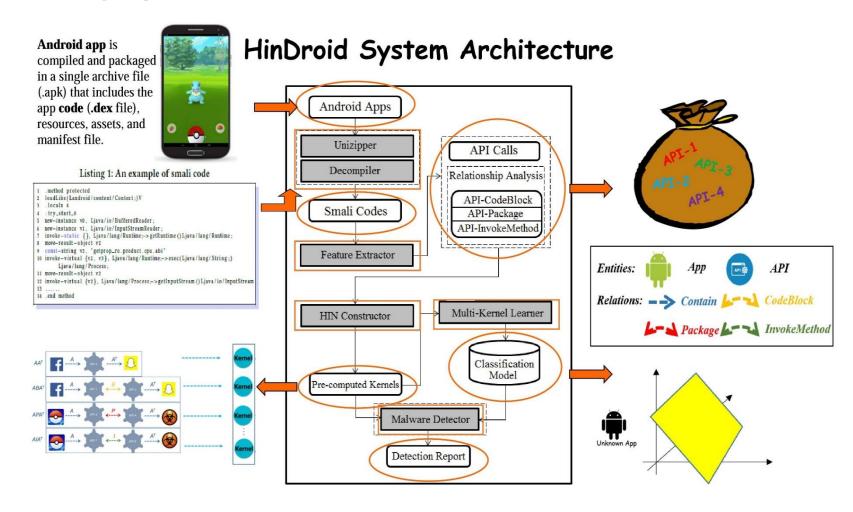


Т	对Android应用文件进行分析,判断是否为恶意软件
I	APK文件
Р	1.预处理 2.特征构建 3.机器学习模型分类
0	APK文件是否为恶意

Р	基于常规特征检测方式容易被代码改写等方式绕过
С	具备大量有标签APK样本
D	构建表征能力强的特征
L	2017 KDD 顶级会议



整体框架图

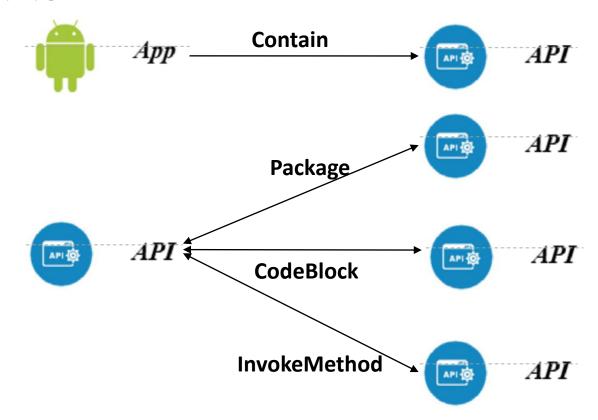




- 特征构建
 - 2种实体四种关系
 - 异构信息网络
- · 基于HIN的多核学习模型构建

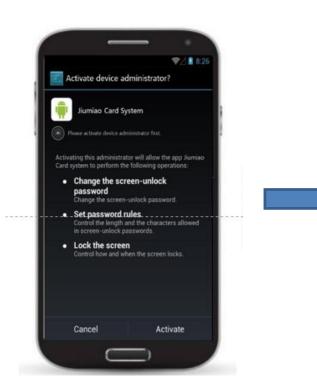


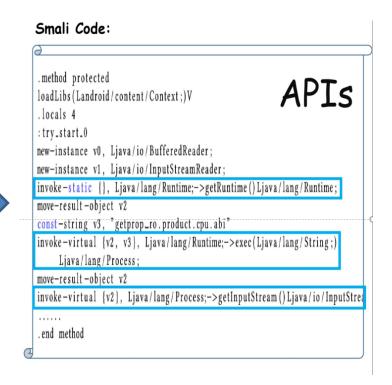
- 2种实体
 - APP(微信、QQ、微博等)
 - API(文件读写、网络连接等关键API)
- 四种关系





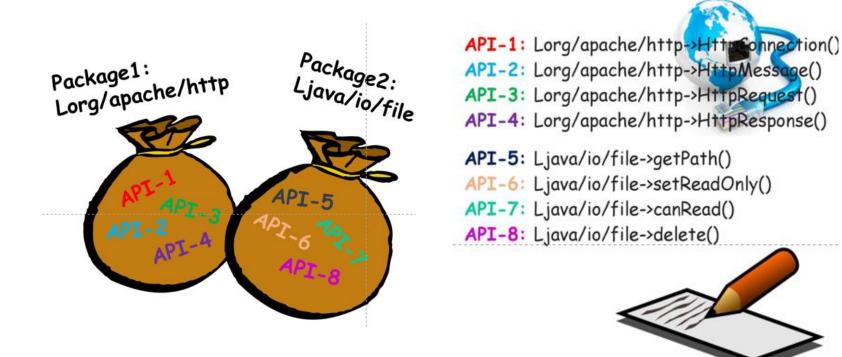
关系1 -> Contain关系(APP-API)







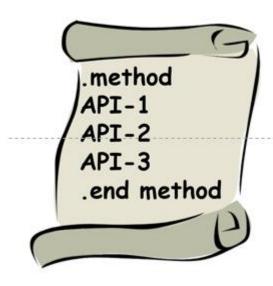
关系2 -> Package关系(API-API)





关系3 -> CodeBlock关系(API-API)

CodeBlock



Smali Code:

```
.method protected
                                                        APIS
loadLibs (Landroid/content/Context;)V
locals 4
:trv_start_0
new-instance v0, Ljava/io/BufferedReader;
new-instance v1, Ljava/io/InputStreamReader;
invoke-static {}, Ljava/lang/Runtime;->getRuntime()Ljava/lang/Runtime;
move-result-object v2
const-string v3, "getprop_ro.product.cpu.abi"
invoke-virtual {v2, v3}, Ljava/lang/Runtime;->exec(Ljava/lang/String;)
     Liava/lang/Process:
move-result-object v2
invoke-virtual {v2}, Ljava/lang/Process;->getInputStream()Ljava/io/InputStream
. . . . . .
.end method
```



关系4 -> InvokeMethod关系(API-API)

Smali Code:

```
.method protected
loadLibs(Landroid/content/Context;)V
.locals 4
:try.start.0
new-instance v0, Ljava/io/BufferedReader;
new-instance v1, Ljava/io/InputStreamReader;
invoke-static {}, Ljava/lang/Runtime;->getRuntime()Ljava/lang/Runtime;
move-result-object v2
const-string v3, *getprop_ro.product.cpu.abi*
invoke-virtual {v2, v3}, Ljava/lang/Runtime;->exec(Ljava/lang/String;)
Ljava/lang/Process;
move-result-object v2
invoke-virtual {v2}, Ljava/lang/Process;->getInputStream()Ljava/io/InputStream....
.end method
```

1.invoke-static: 调用实例的静态方法

2.invoke-virtual:调用实例的虚方法

3.invoke-direct: 调用实例的直接方法

4.invoke-super: 调用实例的父类方法

5.invoke-interface: 调用实例接口方法

算法原理

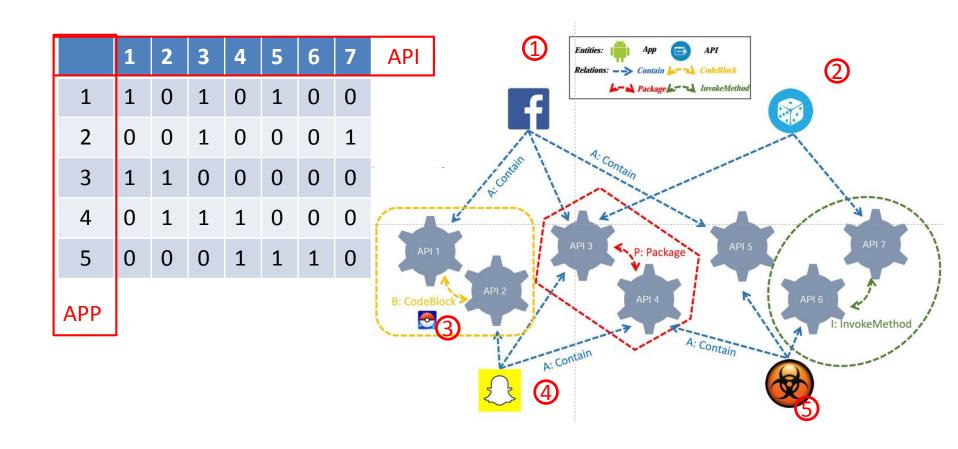


• 四个邻接矩阵

矩阵	元素	含义
А	a_{ij}	APP_i 包含 API_j , a_{ij} =1 否则, a_{ij} =0
В	b_{ij}	API_i 和 API_j 在同一个函数中, b_{ij} =1 否则, b_{ij} =0
Р	p_{ij}	API_i 和 API_j ,有相同的包名 p_{ij} =1 否则, p_{ij} =0
l	i_{ij}	API_{ij} 和 API_{ij} 被调用的方式相同, i_{ij} =1 否则, i_{ij} =0

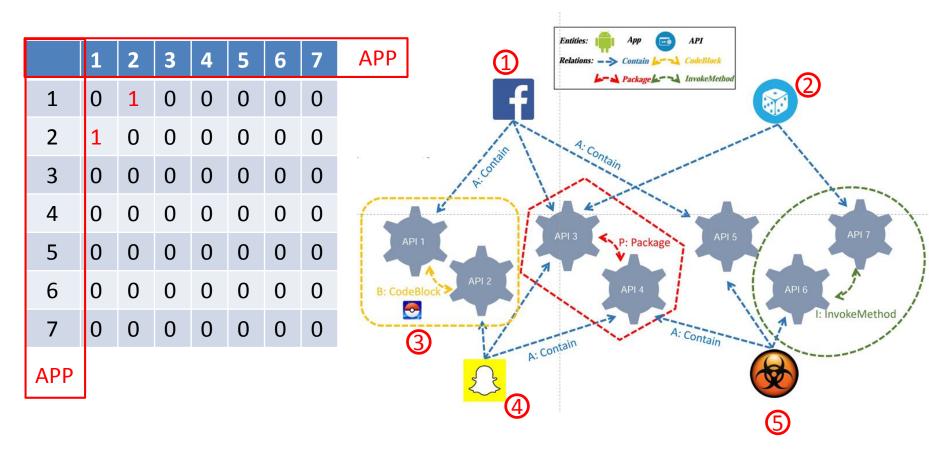


• A矩阵(5x7)





• B矩阵(7x7) API_i 和 API_j 在同一个函数里, p_{ij} =1 否则, p_{ij} =0



I矩阵、P矩阵行列的大小都是7x7



交换矩阵构建(只选取对称元路径)

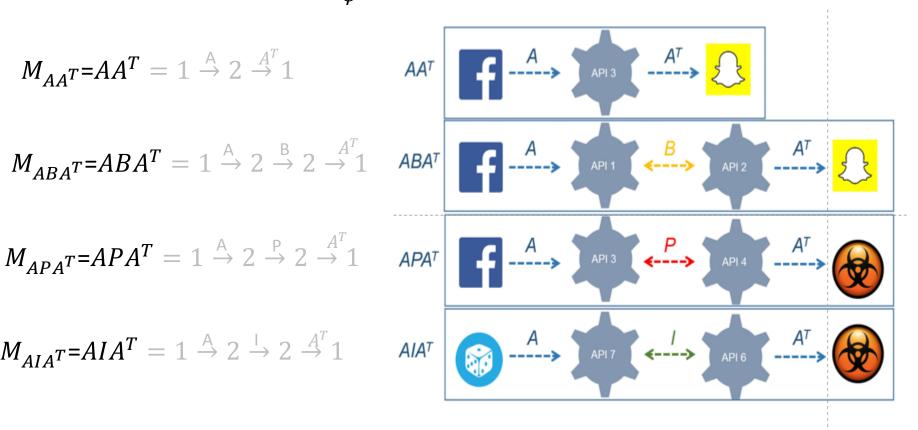
Commuting matrix: M_{ϕ}

$$M_{AA}T = AA^T = 1 \stackrel{A}{\rightarrow} 2 \stackrel{A^T}{\rightarrow} 1$$

$$M_{ABA}T$$
= $ABA^T=1\stackrel{\mathsf{A}}{ o}2\stackrel{\mathsf{B}}{ o}2\stackrel{A^T}{ o}1$

$$M_{APA}^{T} = APA^{T} = 1 \stackrel{A}{\rightarrow} 2 \stackrel{P}{\rightarrow} 2 \stackrel{A^{T}}{\rightarrow} 1$$

$$M_{AIA}$$
 = $AIA^T = 1 \stackrel{A}{\rightarrow} 2 \stackrel{I}{\rightarrow} 2 \stackrel{A^T}{\rightarrow} 1$





• Commuting matrix(M_{φ})含义

	1	2	3	4	5	APP
1						
2						
3						
4						
5						
APP						

$M_{\varphi}(i,j)$ 物理含义?

以
$$AA^T$$
 为例:

 $t(x_i)=[t_{i1},t_{i1}....,t_{in}]$ $\mathbf{t}(x_i)$ 表示特定元路径下一个APP的特征

在 AA^T 路径下: $\mathbf{t}(x_i)$ 表示两个APP之间 具有的相同的API

 $M_{\varphi}(i,j) = t(x_i) t(x_j)$: 具有的相同的API 为特征,计算两个APP之间的相似度

$$M_{\varphi}(i,j) = t(x_i) t(x_j)$$

表示在元路径表征的关系特征下,两个 APP之间的相似度。



- SVM分类
 - 假如有数据集合D: m个有标签的APP文件

$$D = \{(x_1, y_1), (x_2, y_2), \dots (x_m, y_m)\}$$
 $y_i \in (-1, +1)$ 其中, x_m 代表APP的id, y_m 代表APP标签

条件一



Commuting matrix M_{φ} (m x m)

条件二



$$f(x) = w^T \varphi(x) + b$$

目标



模型构建

- 假设存在一个核函数 $k(x_j, x_i) = \varphi(x_j)^T \varphi(x_i)$

$$\min_{w,b} \frac{1}{2} ||w||^2$$

s.t.
$$y_i^*(w^T \varphi(x_i) + b) \ge 1, i = 1,2,3 \dots m$$



通过拉格朗日乘子,转换为对偶问题

$$\max_{\alpha} \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i=1}^{m} \sum_{j=1}^{m} \alpha_i \alpha_j \ y_i \ y_j \varphi(x_j)^T \varphi(x_i)$$

 $M_{\varphi}(i,j) = k(x_i, x_i)$

$$s.t. \sum_{i=1}^{m} \alpha_i y_i = 0$$



二次规划算法

求解出 α_i



- 模型构建
 - 判别函数生成

$$f(x) = w^{T} \varphi(x) + b$$

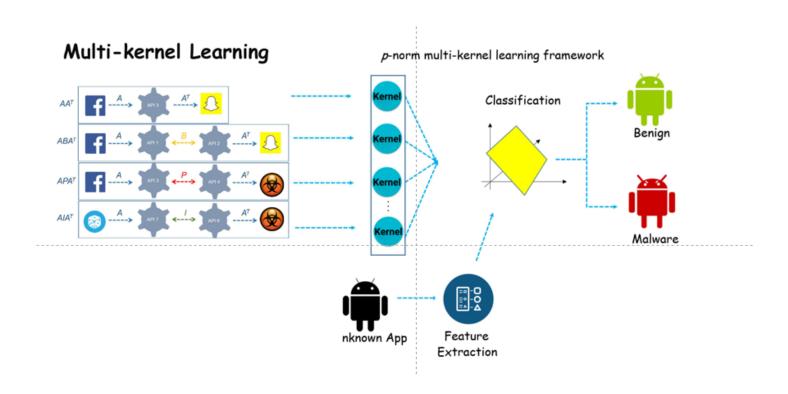
$$= \sum_{i=1}^{m} \alpha_{i} y_{i} \varphi(x_{i})^{T} \varphi(x) + b$$

$$= \sum_{i=1}^{m} \alpha_{i} y_{i} k(x, x_{i}) + b$$

对于新来的样本x,可以根据f(x)完成分类



· 基于HIN的多核学习模型



实验分析



• 评价指标

Table 2: Performance indices of Android malware detection

Indices	Description
TP	# of apps correctly classified as malicious
TN	# of apps correctly classified as benign
FP	# of apps mistakenly classified as malicious
FN	# of apps mistakenly classified as benign
Precision	TP/(TP + FP)
Recall	TP/(TP + FN)
ACC	(TP + TN)/(TP + TN + FP + FN)
<i>F</i> 1	2 * Precision * Recall/(Precision + Recall)

实验分析



• 单核和多核学习效果对比分析

Table 3: Detection performance evaluation

PID	Method	F1	β	ACC	TP	FP	TN	FN
1	AA^T	0.9529	0.1069	94.40%	283	19	189	19
2	ABA^T	0.9581	0.0900	95.00%	286	9	189	16
3	APA^T	0.9495	0.0858	94.20%	273	0	198	29
4	AIA^T	0.9183	0.0623	90.40%	270	16	182	32
5	$ABPB^TA^T$	0.9479	0.0670	94.00%	273	1	197	29
6	$APBP^TA^T$	0.9502	0.0565	94.20%	277	4	194	25
7	$ABIB^TA^T$	0.8683	0.0639	84.60%	254	29	169	48
8	$AIBI^TA^T$	0.8722	0.0639	85.00%	256	29	169	46
9	$\mathbf{APIP}^T \mathbf{A}^T$	0.8373	0.0445	81.20%	242	34	164	60
10	$AIPI^TA^T$	0.8761	0.0572	86.60%	237	2	196	65
11	$\mathbf{A}\mathbf{B}\mathbf{P}\mathbf{I}\mathbf{P}^T\mathbf{B}^T\mathbf{A}^T$	0.9184	0.0616	90.80%	259	3	195	43
12	$\mathbf{APBIB}^T \mathbf{P}^T \mathbf{A}^T$	0.8597	0.0617	84.60%	236	11	187	66
13	$ABIPI^TB^TA^T$	0.9284	0.0426	91.80%	266	5	193	36
14	$AIBPB^TI^TA^T$	0.8237	0.0426	82.60%	218	3	195	84
15	$AIPBP^TI^TA^T$	0.8597	0.0469	81.60%	215	5	193	87
16	$\mathbf{APIBI}^T \mathbf{P}^T \mathbf{A}^T$	0.8597	0.0458	84.60%	236	11	187	66
17	Combined-kernel (5)	0.9214		91.20%	258	0	198	44
18	Combined-kernel (16)	0.9740		96.80%	300	14	184	2
19	Multi-kernel (5)	0.9834		98.00%	297	5	193	5
20	Multi-kernel (16)	0.9884		98.60%	299	4	194	3

实验分析



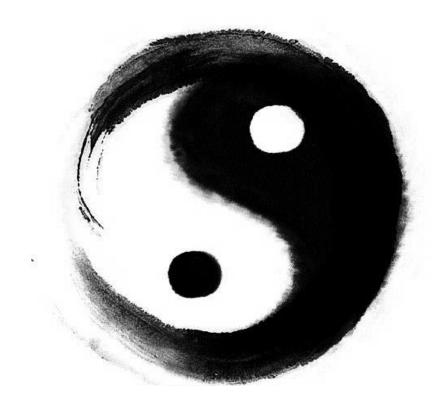
• 对比分析

Table 5: Comparisons with other mobile security products

Family	Sample #	Norton	Lookout	CM	HinDroid
Lotoor	78	75	74	76	78
RevMob	52	46	50	48	52
Malapp	33	29	32	30	33
Fakebank	31	29	30	29	30
Generisk	29	29	29	29	29
GhostPush	19	15	16	18	18
Fakegupdt	16	15	14	14	16
Danpay	21	19	20	20	21
HideIcon	12	11	9	8	12
Idownloader	11	10	9	9	10
Total	302	278	283	281	299
DetectionRate - 92.05% 93.71% 93.05% 99.01					99.01%

For the comparisons, we use all the latest versions of the mobile security products (i.e., Clean Master (CM): 2.08, Lookout: 10.9-7f33b3e, and Norton: 3.17.0.3205)





优劣分析



- 表征能力强。让恶意软件变种难以逃脱检测。
- 基于静态分析,无法处理加固或者具有动态加载功能的恶意软件。

应用总结



框架具有较好的扩充性。可运用与医疗或者生物信息等各个方面。

季考文献



- Abdulhayoglu, Melih, et al. "HinDroid: An Intelligent Android Malware Detection System Based on Structured Heterogeneous Information Network." ACM SIGKDD International Conference on Knowledge Discovery and Data Mining ACM, 2017:1507–1515.
- 2. https://blog.csdn.net/swy520/article/details/78962895
- 3. http://home.cse.ust.hk/~yqsong/papers/2017-HIN-Metagraph.pdf



谢谢!

大成若缺,其用不弊。大盈若冲,其用不穷。大直若屈。 大巧若拙。大辩若讷。静胜 躁,寒胜热。清静为天下正。

