

[13]2020-Towards a Deep Learning Model for Vulnerability Detection on Web Application Variants-ICSTW-PHP--56篇参考文献

• 总结

• 数据集----SARD, File, PHP, 随机, 7: 3

• 代码特征表示

• 通过VLD转换为Opcode

• 向量化

• 通过Vocabularycomposed of VLD opcodes来生产整数, 也即numeric vector

• $i*5+j+1$

• 表1

TABLE I: Vocabulary composed of VLD opcodes. The index of the opcode in position i, j is given by $i * 5 + j + 1$.

i \ j	0	1	2	3	4
0	NOF	NOF	ADD	SEB	MEET
1	DIV	MOD	SR	SR	CONCAT
2	BW_AND	BW_AND	BW_NOT	BW_NOT	BW_NOT
3	BOOL_XOR	IS_IDENTICAL	IS_NOT_IDENTICAL	IS_EQUAL	IS_NOT_EQUAL
4	IS_SMALLER	IS_SMALLER_OR_EQUAL	IS_EQUAL	QW_ASSIGN	ASSIGN_ADD
5	ASSIGN_SUB	ASSIGN_MUL	ASSIGN_DIV	ASSIGN_AND	ASSIGN_OR
6	ASSIGN_OR	ASSIGN_CONCAT	ASSIGN_BW_OR	ASSIGN_BW_AND	ASSIGN_BW_XOR
7	PRE_INC	POST_INC	POST_DEC	POST_DEC	ASSIGN
8	ASSIGN_REF	ECHO	GENERATOR_CREATE	JMP	JMPZ
9	IMPOZ	IMPOZ	JMPZ_EX	JMPZ_EX	CASE
10	CHECK_VAR	SEND_VAR_NO_REF_EX	MAKE_REF	BOOL	FAST_CONCAT
11	ROPE_INT	ROPE_ADD	ROPE_END	BEGIN_SILENCE	END_SILENCE
12	INT_FCALL_BY_NAME	DOLFCALL	INT_FCALL	RETURN	RECV
13	RECV_INT	SEND_VAL	SEND_VAL_EX	SEND_REF	NEW
14	INT_NS_FCALL_BY_NAME	FREE	INT_ARRAY	ADU_ARRAY_ELEMENT	INCLUDE_OR_EVAL
15	UNSET_VAR	UNSET_VAR	UNSET_VAR	PL_RESET_R	PL_FETCH_R
16	EXIT	UNSET_VAR	UNSET_VAR	PL_FETCH_R	PL_FETCH_W
17	FETCH_OBL_W	FETCH_OBL_W	FETCH_OBL_W	FETCH_OBL_W	FETCH_OBL_W
18	FETCH_OBL_W	FETCH_OBL_W	FETCH_OBL_W	FETCH_OBL_W	FETCH_OBL_W
19	FETCH_OBL_FUNC_ARG	FETCH_OBL_FUNC_ARG	FETCH_OBL_FUNC_ARG	FETCH_OBL_FUNC_ARG	FETCH_OBL_FUNC_ARG
20	FETCH_OBL_FUNC_ARG	FETCH_OBL_FUNC_ARG	FETCH_OBL_FUNC_ARG	FETCH_OBL_FUNC_ARG	FETCH_OBL_FUNC_ARG
21	EXT_NOP	TICKS	SEND_VAR_NO_REF	CATCH	THROW
22	FETCH_CLASS	CLOSE	RETURN_BY_REF	INT_STATIC_METHOD_CALL	INT_STATIC_METHOD_CALL
23	ISSET_ISEMPTY_VAR	ISSET_ISEMPTY_OBL	SEND_VAL_EX	SEND_VAR	INT_USER_CALL
24	UNKNOWN[0]	UNKNOWN[0]	UNKNOWN[0]	UNKNOWN[0]	UNKNOWN[0]
25	VERBP_RETURN_TYPE	PL_FETCH_W	PL_FETCH_W	PL_FREE	INT_DYNAMIC_CALL
26	DOLFCALL	DOLFCALL	DOLFCALL_BY_NAME	PRE_INC_OBI	PRE_DEC_OBI
27	POST_INC_OBI	POST_DEC_OBI	ASSIGN_OBI	OF_DWA	INSTANCEOF
28	DECLARE_CLASS	DECLARE_INHERITED_CLASS	RASL_ABSTRACT_ERROR	DECLARE_CONST	DECLARE_CONST
29	ADD_INTERFACE	VERIFY_INSTANCEOF	VERIFY_ABSTRACT_CLASS	ASSIGN_OBI	ISSET_ISEMPTY_PROP_OBI
30	HANDLE_EXCEPTION	USER_OPCODE	ASSERT_CHECK	JMP_SET	DECLARE_LAMBDA_FUNCTION
31	ADD_THROW	BINDTRAP	SEPARATE	FETCH_CLASS_NAME	JMP_SET_NA
32	DISCARD_EXCEPTION	YIELD	GENERATOR_RETURN	FAST_CALL	FAST_RET
33	REC_VARIADIC	SINGL_UNPACK	POW	ASSIGN_POW	BIND_GLOBAL
34	COALESCE	FETCH_STATIC_PROP_W	DECLARE_ANON_CLASS	DECLARE_ANON_INHERITED_CLASS	FETCH_STATIC_PROP_R
35	FETCH_STATIC_PROP_W	FETCH_STATIC_PROP_W	FETCH_STATIC_PROP_IS	FETCH_STATIC_PROP_FUNC_ARG	FETCH_STATIC_PROP_UNSET
36	UNSET_STATIC_PROP	ISSET_ISEMPTY_STATIC_PROP	FETCH_CLASSICAL_CONSTANT	BIND_LEXICAL	BIND_STATIC
37	FETCH_THIS	UNKNOWN[16]	ISSET_ISEMPTY_THIS	SWITCH_LONG	SWITCH_STRING
38	IN_ARRAY	COUNT	GET_CLASS	GET_CALLED_CLASS	GET_TYPE
39	PUNC_NCM_ARGS	PUNC_GET_ARGS	ISSET_ISEMPTY		

• 模型

• 总

• DL 模型由多层构成。每一层都接收前一层的输出作为输入, 并对输入应用一些转换。通过具有多个层, 模型可以基于更简单和更广泛的概念 (来自前几层) [24]学习更复杂和更有用的概念。

• Embeddding

• 将正整数, 也即索引值, 转换为固定大小的向量

• N LSTM+N Dropout

• 考虑每个操作码的上下文, 也即保存顺序性

• 减少过拟合

• Dense

• 1

• 学习the slice和label之间的关系, 将输入转换成一个具有相同形状的向量

• 2

• 对the slice进行分类并存在SQLi漏洞的概率, sigmoid函数[44]

• 优化器+调参

• 优化器

- 三种常用的基于 SGD 的优化器
 - ADADELTA[25]
 - 它使用动态学习率 (lr) 计算每个维度
 - RMSProp[26]
 - 它使用重新缩放的梯度生成更新
 - ADAM[27]
 - 它使用梯度的运行平均值生成更新。
- 调参
 - HS - HIDDEN SIZE,---学习能力相关
 - LSTM中HIDDEN SIZE的单元数
 - 拥有的单元越多，学习的内容也越多
 - δ - dropout rate,----减少过拟合
 - 它表示dropout rate，一个与dropout Layer相关的0到1之间的值，并对应于输入向量的每个条目变为0的概率。
 - NE - number of epochs---提高学习效果
 - 定义优化器更新参数的次数
 - 一般来说，该参数越高，模型学习效果越好。但是，我们需要平衡学习和泛化能力，不要让模型过拟合。

• 评价

- 在训练集中对每个模型进行了3次10倍交叉验证

• 结果

- 最佳模型达到了95%以上的**准确率**，**正确率**和**召回率**，且RMSProp优化器最佳
- 表 5

TABLE V: Results of the accuracy, precision and recall for the various configurations analysed.

Optimizer	HS	δ	NE	Accuracy	Precision	Recall
ADADELTA	80	0.30	160	0.9487	0.9837	0.9344
RMSProp	80	0.15	70	0.9535	0.9651	0.9614
ADAM	70	0.30	35	0.9413	0.9876	0.9189

HS - HIDDEN_SIZE, δ - dropout rate, NE - number of epochs

• 注意的点

- 数据集问题
 - 输入可能被净化，但仍会危及应用程序；-----过滤未全面，可以绕过。
 - 恶意输入可能会在对其他变量的赋值过程中，在代码片之间进行传播----恶意代码存在恶意传递，数据流分析
- 优化器超参数保留问题--默认值
 - 简化第一个实验并想获得对该问题的一些直觉
- 加入n LSTM layers的原因---[44]

- 在漏洞检测任务中，操作码在切片中的顺序是非常相关的。
- 产生输出向量，编码之前的操作码及其顺序的信息
- 加入Dropout layer的原因
 - 位置，在每个LSTM的后面
 - 可以进一步减少过拟合--概率--设置为零--噪声-防止学习到无关模式
 - 该层根据给定的概率(δ)，将其输入的一些条目随机设置为零。该方法的目标是在模型中引入噪声，以防止它记忆由LSTM层偶然学到的无关模式。
- Opcode的优势
 - 类似于C/C++的汇编操作码的操作码格式处理PHP代码片，被认为是一种中间语言。
 - 更接近于语言的内部结构---有助于解放分类任务
 - 由于切片可以属于任何web应用程序变体，使用中间语言来表示它们可以促进变体之间的分析横向性。
 - 这种中间语言尚未用于解决这个任务
- 嵌入层的作用
 - CNN 和 RNN 层之前通常有一个嵌入层，它将离散符号映射到连续向量，解决自然语言数据的稀疏性问题。此外，通常将这些组件的输出提供给前馈组件，该组件学习执行所需的任务，如分类[28]。
 - one-hot vectors and embedding vectors
 - embedding vectors
 - CBOW
 - skip-gram