

ARL

Grammatical inference and language frameworks for LANGSEC

Kerry N. Wood Richard Harang U.S. Army Research Laboratory 21 MAY 2015

The Nation's Premier Laboratory for Land Forces

UNCLASSIFIED





- 1. My background is primarily in networking, modeling, and optimization. (I am <u>not an</u> expert in formal language theory.)
- 2. A natural beginning for our work was a solid literature review. This paper and presentation are a direct result of trying to summarize that work.
- 3. This is not comprehensive. We focused on those frameworks and techniques we thought might be useful.

We came into this knowing that inferring grammars is difficult. However, rewriting all the existing code to use LANGSEC is *probably* harder.



Our contribution is mostly a survey of useful methods.



- A brief introduction and review of grammars and LANGSEC. (OR) What we think LANGSEC is / is supposed to be.
- Models of learning: definitions are important.
- Models of teaching: is data available / how is it presented?
- Pattern Languages / EFS : our two favorite formalisms.
- Some initial work : where "initial" means "small amount."
- Conclusions and future work : hope springs eternal.



We all remember Chomsky classes and the recognizer types.



CS101 – formal production rules for 'strings' that have an exact mapping to different models of computation

Grammar	Power to recognize
Regular	State machine
Context-free	Stack machine
Context-sensitive	Bounded automata
Unrestricted	Turing machine

Grammatical inference:

- You are given the output of the grammar (strings).
- You have to identify the generating structure.
- In essence, it is *decompiling*. Just harder.





Trouble begins when we mean to define a recognizer for S, but actually define one for A (and don't know it).

UNCLASSIFIED



Why not just parse?

"LangSec" – language theoretic security

- Use only simple grammars with provable properties
- Validate all inputs before any processing; reject failures aggressively



Right now: really only helpful if you "bake it in"

Significant restrictions on message complexity

Deterministic CFG or weaker (or some non-Chomsky class recognizable in polynomial time

Learning: access to examples.



Learning from *positive examples:*

U.S. ARMY RDECOM®

- All strings in the grammar are shown to the learner $(S_+ = \forall s \in \mathbb{L})$
- Positive and negative examples (aka "complete presentation"):
 - Members and non-members are labeled and shown to the learner.
 - $(S_+ : \forall s \in \mathbb{L})$ and $(S_- : \forall s \in \Sigma^* \mathbb{L})$
- Learning with an oracle (aka "membership queries"):
 - The learner presents strings, the oracle returns membership labels.
- Learning with a teacher (aka "equivalence queries"):
 - The learner presents guesses as to the hypothesis grammar.
 - Teacher verifies correctness, or returns an example string in the language, but not in the learner-suggested grammar.
- **Compressed / simple examples:**
 - The minimum set of examples that allow for identification of the grammar.



Models of learning.



In general: a learner \mathcal{A} takes some data D, and outputs a (possible) generating grammar \mathbb{G} .

Learning in the limit (Gold-style)

U.S. ARMY RDECOM®

- A learner is presented with examples from the entire grammar.
- At each iteration, it guesses a hypothesis grammar.
- At a certain iteration, the guesses never change.

Fixed-time identification

- Learner will see strings in *D* one at a time; determines beforehand how many steps it will take to produce a correct $\mathcal{A}(D)$.

Finite identification

As above, but learner decides after seeing each string whether to stop

Probably Approximately Correct (PAC) learning

- Given access to examples, a finite-length grammar, and finite-length examples; with probability $(1 - \delta)$, a learner will output a hypothesis that is ϵ - good.





Pattern languages: replace variables with substrings.



- Constants and strings: $\Sigma, \Sigma^*, \Sigma^+$
- Variables: $(x_1, x_2, ...) \in \mathcal{X}$
- Patterns : $\mathcal{P} = \pi_1, \pi_2, \dots = \Sigma^* \cup \mathcal{X}$
- Substitution: θ

•
$$\pi\theta = \left[\frac{s_1}{x_1}, \frac{s_2}{x_2}, \dots, \frac{s_k}{x_k}\right]$$
 where $|s_i|! = 0$

A pattern π_i is <u>regular</u> if each variable appears at most once.







The good:

U.S.ARMY

They look just like regular expressions.

Angluin showed they are learnable.

U.S. ARMY RDECOM®

- Lange and Wiehagen have a simple, but *inconsistent* algorithm.

The bad:

- In general, intractable to determine membership.
 - NP − complete
- In general, they do not map to Chomsky.
 - Regular patterns encode regular languages.
- In general, relationships are hard to determine
 - Not closed under union.
- Lots of "academic" sub-classes and cases.
- PAC learnability is bad.

Examples: $\overline{\mathcal{X}} = \{x_1, x_2, ...\}$ xaybza $\chi\chi$



Elementary formal systems (EFS) are a

logic programming analog.

ARL

$$\{a^{n}b^{n}c^{n} \mid n \ge 1\} :$$

$$\Gamma = \begin{cases} p(x_{1}, x_{2}, x_{3}) \leftarrow q(x_{1}, x_{2}, x_{3}) \\ q(ax_{1}, bx_{2}, cx_{3}) \leftarrow q(x_{1}, x_{2}, x_{3}) \\ q(a, b, c) \leftarrow \end{cases}$$

$$\{a^{n}b^{n} \mid n \ge 1\} :$$

$$\Gamma = \begin{cases} p(ax_{1}b) \leftarrow p(x_{1}) \\ p(ab) \end{cases}$$

EFS' with specific properties map INCLASSIFIED EFS' with specific properties map INCLASSIFIED A CLAUSE / EFS ($A \leftarrow B_1, \dots, B_m$) is... **A Clause / EFS (** $A \leftarrow B_1, \dots, B_m$) is... **Variable bounded** $- v(A) \supseteq v(B_1) \cup \dots \cup v(B_m)$ **Recursively enumerable Inclusion Incluston Inclusion Inclusion Inclusion Incluston Inclusion Inclusion Inclusion Incluston I**

- π is regular for all $\pi \in \Gamma$

Right / left linear

- pattern of the head is xw for some $w \in \Sigma^+$

- pattern of the head is wx for some $w \in \Sigma^+$

Regular

Context-free

Pattern languages are a single-clause EFS:
 Γ = {p(π) ←}

UNCLASSIFIED

Inferability results are *slightly* more encouraging.



U.S. ARMY RDECOM®

 $- p(\pi_1, \dots, \pi_n) \leftarrow q_1(\tau_1, \dots, \tau_{t_1}), q_2(\tau_{t_1}, \dots, \tau_{t_2}), \dots, q_l(\tau_{t_{l-1}+1}, \dots, \tau_{t_l})$

- If for each $j = 1, ..., t_l$ pattern τ_j is a substring of some pattern π_i .

LB-H-EFS(m,k):

- − $|A\theta| \ge |B_1\theta| + \dots + |B_m\theta|$ (length bounded) <u>and</u> hereditary
- $\leq m$ clauses
- $\leq k$ variable occurrences in the head of each clause

Language in Chomsky Hierarchy	EFS	Gold Inferable		(Hereditary) ^b polynomial-PAC Inferable	
	# of clauses \rightarrow	∞a	N ^a	∞	(<i>m</i> , <i>k</i>)
recursively enumerable	Variable bounded	NO	NO	NO	
context-sensitive	Length bounded	NO	YES	NO	NO ^c
context-free	regular	NO	YES	NO	YES
regular	right/left linear	NO	YES	NO	YES



Many problems involve differentiating these classes.



UNCLASSIFIED

Q: If we look at these as grammar inference problems, can we develop any intuition / bounds / definition of success?

Security problem	Subset	Learning model
Supervised learning	SUA \ S	PAC with positive and negative examples
Anomaly detection, type 1	S	PAC with positive examples
Anomaly detection, type 2	A	PAC with noisy positive examples
Fuzzing	$\mathbb{A}\setminus\mathbb{S}$	Learning with an oracle
Unsupervised clustering	Subsets of S or A \ S	PAC with positive examples in a mixture setting



UNCLASSIFIED

Simple \rightarrow inferable.



Hypothesis: *most* data "in the wild" comes from simple grammars.

A protocol or code may claim to accept a large, complex grammar. Really, most messages are members of simpler sub-grammars (or unions of sub-grammars).

Sometimes classifying messages is a good place to start.

• Are most messages as complex as the overall grammar?

• Use program structure as a guide.

- Can we apply a grammar to those control flow data?
- Can we segment a large grammar via sub-grammars?

Weblog data from live system.

POST Requests (Anonymized URIs)	# of Samples
/7NWS/XBPD07F/HXZRT/6BR/PR9M/_x0Q70J	83199
/0YDO8J33LKSC/DKWPZJ/QCSY1BWNRIG65;R9LZ64B6GI= <u>x0</u> ! <u>x1</u> ! <u>x2</u>	2366
/5HEZP8EKVK3R9RGF9D_x0x1x2x3_	204950
GET Requests (Anonymized URIs)	# of Samples
GET Requests (Anonymized URIs) /KPGF/7OH0EHB1.HE9?id= x0 %F95X4L2% x1 x2 x3 x4 ' x5 x6 % x7 % x8 x2 x3 x11 ' x12 x6 % x14 % x15 = ' x16 '	# of Samples

Lange and Weihagen's (LW) algorithm was applied to weblog data.

- Log entries are batched by endpoint and "inferred."
- A straightforward, initial example.

U.S. ARMY RDECOM®

UNCLASSIFIED

UNCLASSIFIED		
EXAMPLE COIVI Fuzzing [?] Grammar Infer	ence	ARL
<corpus>/jpeg/full/images/*src:000619*</corpus>		
ffd8ffe000104a46494600010102001c001c0000ffdb004300281c1e231e19282321232d2b 28303c64413c37373c7b585d4964918099968f808c8aa0b4e6c3a0aadaa x0 8a8cc8ffcbd aeef5fffff9bc1ffffffaffe6fdfff8ffdb0043012b2d2d3c353c76 x1 4176f8a58ca5f 8f8f8f8f8f8f8f8f8f8f8f8f8f8f8f8f8f8f8	25	(504,117)
<pre>/jpeg/edges-only/images/*src:003554*</pre>		
ffd8ffdd0004_x0_00b400e8ffdb00438028101efffaffe6fdfff8ffdb0043012b2d2d3c35 3c64414176f8a58ca5f8f8d500010000f8f8f8f0f8f8f84a64653d79906fc317d3caf8f8f8 f8f8f8f8f8f801001c0000ffdb004300281c1e231e19282336232d2bef303c64413c37373c 7b585d49649180999600018c8aa0b4ad8a8cc8ffe4daee097affff9bc1ffffffaffd1fdff f8ffdb0043012b2d2d3cf8f8f8f8e3f8f0e9f900f8f84a6a653d79906fdf17d3b7f8f8f8f8 f8f8f8f8f8ffca001108_x2_0_x3_x4_x4_2003012200021101031101ffda00060001180 12200021101031101ffda00060001030003000100000000004000000001ffda000801010 101010101010 x1	12	(260,0)
<pre>/jpeg/edges-only/images/*src:000512*</pre>	1.0	(170,80)
ffd8ffe000104a x0_x3_x4_x3_x1_0_x2_		

• LW applied as in previous slide.

UNCLASSIFIED

Zalewski's AFL does the test case generation / organizing.



Future work.



- Explore ongoing work on binary analysis / malware autogeneration to help us identify grammars.
 - Instrument binaries via fuzzing.
 - Analyze binaries via lifting and decompiling tools.
 - Learn "low hanging fruit" grammars.
 - If most of the grammar is ignored most of the time, maybe the "simple stuff" is a good place to start.
- Input normalizer / firewall for known applications.
 - Once you know (a) the grammar, enforce the grammar.
- EFS *seems* intuitive. Can we use it to help define protocols and message structures that adhere to specific classes?



Questions?



The Nation's Premier Laboratory for Land Forces

UNCLASSIFIED





Backup Slides

The Nation's Premier Laboratory for Land Forces

UNCLASSIFIED





For example...

Can we program provably interoperable stacks?

- Does your protocol need to be written in BNF? Then the answer is NO
 - BNF is context-free
 - Equivalence between two context free grammars is *undecideable* (unless you know your grammar is deterministic)
 - In other words: if you have two *different* implementations of the same (context free) protocol, then you can *never* prove that they accept the same set of strings (unless you prove that both are deterministic)
 - See Sassaman et al. 2013: x.509 certificate forgeries

Can we write a behavioral malware detector?

- Are you trying to detect the behavior before it actually happens? Then the answer is NO
- Program X: "Halt if this embedded program does Y, otherwise loop forever"; does program X halt?

If I have a bunch of good traffic and a bunch of bad traffic, can I learn a classifier that will generalize?

- -NO (as long as factoring integers is as hard as we think it is)
- Valiant and Kearns: learning even a regular grammar in polynomial time that (with high probability) will have bounded error is not possible
- If you can learn DFAs in the PAC framework, you can factor Blum integers

Can I at least validate that incoming messages are well-formed first?

- PROBABLY NOT FAST ENOUGH (since most network protocols are stronger than context-free (see Davidson et al. 2009)
- Parsing CSGs can be anywhere up to NP (Satta, 1992)





"Finite Thickness" and "finite elasticity"

- Finite Thickness: No (non-empty) set of strings occurs in infinitely many languages in the class under consideration
- Finite Elasticity: There is no infinite set of strings that is consistent with every infinite sequence of languages in the class

The existence of "tell-tales" in the class

- A string or set of strings unique to a single member of the class
- In a security context: "signature"

U.S. ARMY RDECOM®

Finite language, positive presentation, no noise:

- Memorize everything you've ever seen
- Doesn't scale





Having an oracle that knows the grammar and can provide complex responses helps

- To learn DFAs in polynomial time with arbitrary data <u>requires</u> membership and equivalence queries (Angluin, '88)
 - Extends to a small subset of CFGs (those that can be written compactly in CNF; algorithm is polynomial in size of set of symbols)

A training set designed to teach the learner your grammar helps

- CFGs are PAC-learnable from "simple" presentations

Some non-Chomsky grammars are identifiable under various conditions

- Notable example: Pattern languages can be identified from positive examples...
 - ...but are NP-complete to recognize

- Finding minimal DFAs from positive and negative examples is NP-hard
 - Angluin '78; Gold '78

U.S.ARMY

U.S. ARMY RDECOM®

- Approximating a minimal target DFA from positive and negative examples to within any finite polynomial factor is NP-hard
 - Pitt and Warmuth, '93
- Representation-independent result: DFAs are cryptographically hard to learn from labeled examples
 - Existence of a weak learner for a DFA implies a polynomial advantage in recovering the LSB of the plaintext of an RSA-encrypted message
 - Kearns and Valiant, '94
- To learn DFAs in polynomial time with arbitrary data <u>requires</u> membership and equivalence queries
 - Extends to a small subset of CFGs (those that can be written compactly in CNF; algorithm is polynomial in size of set of symbols)
 - Angluin, '88

UNCLASSIFIED

- DFAs are learnable from "simple" examples, but those examples must be selected based on the DFA and learning algorithm
- Very, very simple CFGs appear to be empirically learnable from data...
 - -...With a lot of effort, a little luck, and no theoretical guarantees
 - See, e.g., Omphalos competition