

End-to-End Active Learning for Computer Security Experts

Anaël Beaugnon, Pierre Chifflier, Francis Bach

`anael.beaugnon@ssi.gouv.fr`



ANSSI, ENS Paris, INRIA

AICS 2018



Outline

- 1 Active Learning and Computer Security
- 2 ILAB: an End-to-End Active Learning System
- 3 Importance of Features



Outline

- 1 Active Learning and Computer Security
- 2 ILAB: an End-to-End Active Learning System
- 3 Importance of Features



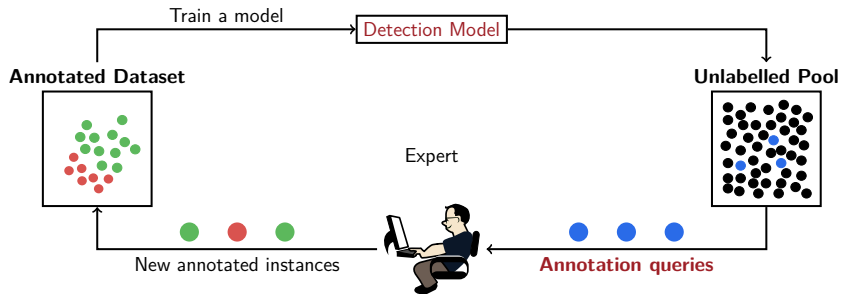
Lack of Representative Training Data !

- ✗ Public datasets \neq deployment environments
- ✗ Crowd-sourcing is not suited for Computer Security

In-situ labelling with Active Learning
Annotate data from the deployment environment



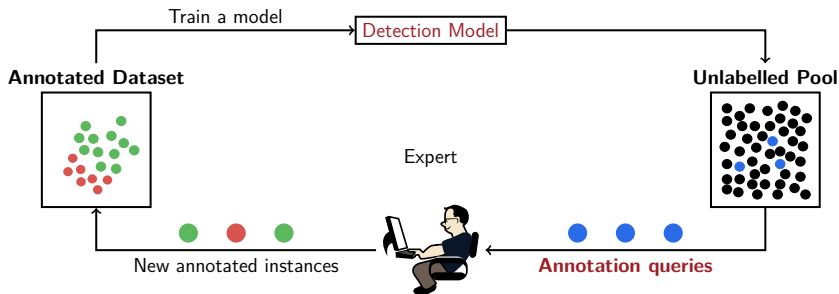
Active Learning





Active Learning

Don't forget the expert !



Active learning is not only a query strategy ...



Outline

- 1 Active Learning and Computer Security
- 2 ILAB: an End-to-End Active Learning System
- 3 Importance of Features



ILAB: an End-to-End Active Learning System

End-to-end active learning system
=
Active learning strategy + Annotation system

Active Learning Strategy

Selects cleverly the instances to be annotated

RAID'17 Beaugnon et al., ILAB: An Interactive Labelling Strategy for Intrusion Detection

Annotation System

Displays the annotation queries and gathers the answers



ILAB: an End-to-End Active Learning System

End-to-end active learning system
=
Active learning strategy + Annotation system

Active Learning Strategy

Selects cleverly the instances to be annotated

RAID'17 Beaugnon et al., ILAB: An Interactive Labelling Strategy for Intrusion Detection

Annotation System

Displays the annotation queries and gathers the answers

How to assess properly a whole active learning system ?



Anomaly detection from NetFlow data

Dataset

Num. flows	$1.2 \cdot 10^8$
Num. IP	463,913
Num. features	134

How does ILAB help experts annotate ?



How does ILAB help experts annotate ?

Uncertain

 >>>

Malicious

 >>>

Benign

Next Iteration

Annotation Queries

Family

slow_scan

4 / 5

Prev

Next

Display Families

Annotation Query

1 / 9

Prev

Next

Instance 374335

Annotation

Suggestion

slow_scan

Malicious Families

ICMP_scan
TCP_Syn_flooding
misconfiguration
obvious_scan
slow_scan

Add

Benign Families

DNS
SMTP
web

Add

Ok

Remove

Description

NetFlows

Features

Start	Duration	Proto	Src IP	Src port	Dst IP	Dst port	Flags	Num bytes	Num packets
08:22:23.341	8.835	TCP		43805		23S.	168	3



How does ILAB help experts annotate ?

Problem

Need for contextual informations to annotate

Solution: Problem-specific visualization

Description

NetFlows

Features

Start	Duration	Proto	Src IP	Src port	Dst IP	Dst port	Flags	Num bytes	Num packets
08:22:23.341	8.835	TCP		43805		23S.	168	3

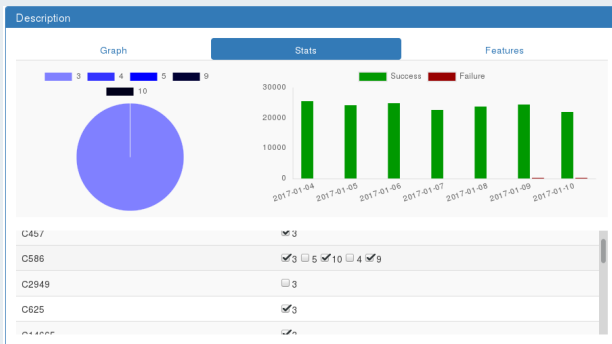


How does ILAB help experts annotate ?

Problem

Need for contextual informations to annotate

Solution: Problem-specific visualization



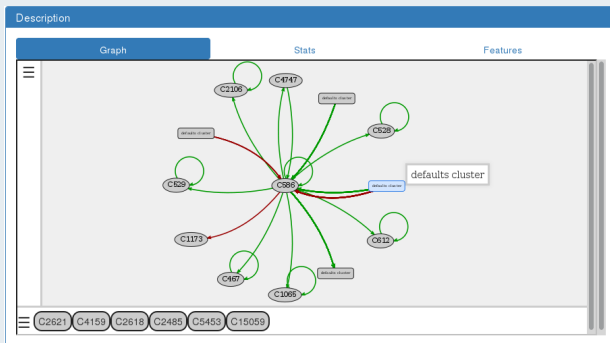


How does ILAB help experts annotate ?

Problem

Need for contextual informations to annotate

Solution: Problem-specific visualization





How does ILAB help experts annotate ?

Problem

- ▶ What families should I create ?
- ▶ How should I group the data ?

Families definitions may evolve across iterations ...

Solution: Family Editor

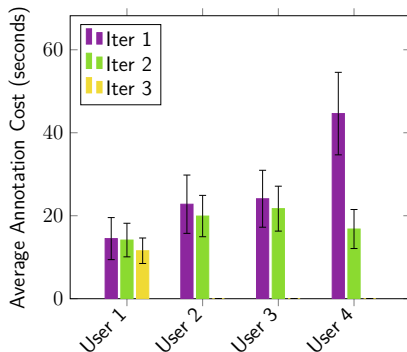
- ▶ Edit the name of a family
- ▶ Merge several families
- ▶ Change the label associated with a family



It is hard to annotate, but ILAB helps !

Solutions

- ▶ Problem-specific visualization
- ▶ Family Editor





Outline

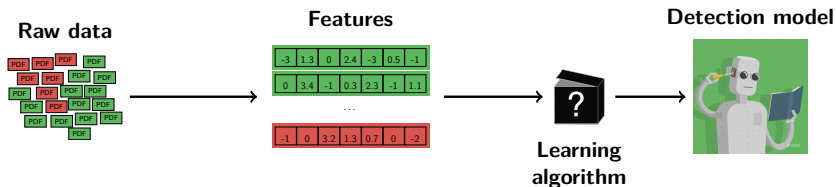
- 1 Active Learning and Computer Security
- 2 ILAB: an End-to-End Active Learning System
- 3 Importance of Features**



Importance of features

Problem

Features may be not expressive enough



Loss of information with feature extraction ...



Importance of features

Problem

Features may be not expressive enough

Features

Num. bytes sent/received:

- ▶ globally
- ▶ on port 80
- ▶ on port 53
- ▶ on port 25

User annotation

- ▶ anomalous
- ▶ traffic on port 1258



Importance of features

Problem

Features may be not expressive enough

Solutions

- ▶ Annotators must know the extracted features
- ▶ Make features evolve across iterations
 - ▶ manually
 - ▶ or even better, automatically

ECML'14 Boule, Towards automatic feature construction for supervised classification

DSAA'15 Kanter et al., Deep feature synthesis: towards automating data science endeavors

EURASIP'16 Šrndić et al., Hidost: a static machine learning based detector of malicious files



ILAB: an End-to-End Active Learning System

<https://github.com/ANSSI-FR/SecuML>

ILAB helps experts to annotate

- ▶ Problem-specific visualization
- ▶ Family editor

Features Expressiveness

- ▶ Annotators must know the features
- ▶ Features should evolve across iterations (future work)