

Problem Def

ML Paper (Sept 26, 2016)

- detection is based on:

- URL
- flow duration
- number of bytes transferred from client to server & other way

- user agent
- referer
- MIME-type
- HTTP status

- The n-dimensional feature vector represents each proxy log & used to differentiate b/w legit & malicious traffic

- Paper model only analysis single proxy log, & skips temporal features

- Attacking domains change frequently but behaviour doesn't,

How

- The proxy logs originating at a particular user machine are grouped into bags based on the domain in the URL.

- The bags are labeled according to the domain.

↳ if domain is in any blacklist, the bag has a +ve label

↳ if not, bag has a -ve label

MIL

- flow is described by a vector of features

$$x \in X \subseteq \mathbb{R}^d \text{ \& a label } y \in Y = \{+1, -1\}$$

↑
malicious

↓
not

- Network traffic monitored in a given period is fully described by the completed annotated data

$$D_{\text{comp}} = \{(x_1, y_1), \dots, (x_m, y_m)\}$$

$$\in (X \times Y)^m$$

independant,
identically distributed

assumed to be generated from i.i.d.
random vars with unknown dist

$$p(x, y)$$

- Annotating everything is expensive, thus we use bags of flows

- The weakly annotated data

$$D_{\text{bag}} = \{ \underbrace{x_1, \dots, x_m}_{\text{features}}, \underbrace{(b_1, z_1), \dots, (b_n, z_n)}_{\substack{\text{assignment} \\ \text{to labeled} \\ \text{bags}}} \}$$

$$\{(b_1, z_1), \dots, (b_n, z_n)\} \in (P \times Y)^m$$

P = set of all partitions of
indices $\{1, \dots, m\}$.

- The i th bag is a set of flow features $\{x_j \mid j \in B_i\}$ label by $z_i \in Y$.

- D_{bag} carries partial info about D_{emp} .

Assumptions:

1) Flow features $\{x_1, \dots, x_m\}$ are the same in both.

2) Negative bags contains a single instance, & the label is correct.

$\Rightarrow z_i = -1$ implies $|B_i| = 1$ & $y_i = -1$

3) +ve bags have a variable size & at least 1 instance is positive

$\Rightarrow z_i = +1$ implies $\exists j \in B_i$ s.t. $y_j = +1$