

Decentralized Clustering for Botnet Detection

Alessio Guerrieri

18 Novembre 2010

OUTLINE

- 1 INTRODUCTION
 - Botnets
 - Clustering
- 2 BOTMINER: CENTRALIZED BOTNET DETECTION FRAMEWORK
 - Overview
 - Framework
 - Results
- 3 DECENTRALIZED BOTNET DETECTION FRAMEWORK
 - Intro
 - Framework
 - Results
 - Future work

BOTNETS

DEFINITION

A group of nodes forms a botnet if they

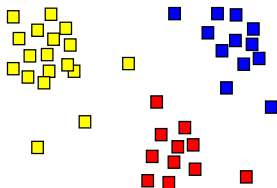
- communicate using a similar pattern
 - commit similar malicious activities
-
- No specialization
 - We can try to clusters nodes with similar behaviour

CLUSTERING ALGORITHMS

CLUSTERING PROBLEM

Given a set of items, divide them in subsets such that items in the same subsets are “similar”

- Items are usually represented as points in R^n
- Similarity measured with distance function



K-MEANS CLUSTERING

- Needs K (number of clusters) as input
- Partitions the points to minimize the within-cluster sum of squares
- Problem is NP-HARD

STANDARD ALGORITHM (LLOYD 1982)

- Choose K centroids at random
- While changes:
 - Assign each point to closest centroids
 - Move each centroid to the average of the points assigned to it

NOTES ON K-MEANS CLUSTERING

- Results depends on starting centroids
- Can terminate at local optima
- Time complexity is $2^{\Omega(n)}$
- Smoothed analysis shows polynomial time

XMEANS (PELLEG 2000)

- Does not need the exact number of clusters
- Uses Bayesian Information Criterion to compare models
- Receives MAX= maximum number of clusters

ALGORITHM

- Run 2-means on the data to get 2 clusters
- While number of cluster less than MAX:
 - For each cluster:
 - Run 2-means locally
 - If improvement, replace cluster with result of 2-means
 - If no replacement, replace worst clusters
- Return best clustering found at any iteration
- Can be speeded up with kd-trees and caching

BOTMINER

- From paper “BotMiner: Clustering Analysis of Network Traffic for Protocol- and Structure-Independent Botnet Detection” (Gu 2008)

FRAMEWORK OVERVIEW

- Each node is mapped to two points in R^n , one created from communication pattern (C-Plane), one from malicious activities (A-Plane)
- The C-Plane and A-Plane are clustered independently
- We look for subset of nodes that are in the same cluster in both planes.

MAPPING

C-PLANE

- number of flows per hour
- the number of packets per flow
- the average size of the packets
- the number of bytes sent per second.

A-PLANE

- Scan activity
- Spam activity
- Binary downloading
- Exploit activity

For each activity, keep additional features (ex: ports)

CLUSTERING

C-PLANE

Two-step clustering using X-Means:

- Global clustering on 8 features
- Local clustering using all features.

A-PLANE

First cluster according to activities, then using the features.

CROSS-CLUSTERING CORRELATION

“BotNet Score” is computed using two components:

- amount of malicious activities of the node
- how many nodes share the same characteristics (Cross-clustering correlation)

CROSS-CLUSTERING CORRELATION

If node n is in clusters A_i and C_j :

$$ccc(n) = \frac{A_i \cap C_j}{A_i \cup C_j}$$

- Nodes with score bigger than threshold are reported

RESULTS

- Used ten botnet traces embedded in anormal traffic trace

Botnet	Bots	Detected?	Clustered	Det. Rate	False Pos.	FP Rate
IRC-rbot	4	YES	4	100%	1/2	0.003
IRC-sdbot	4	YES	4	100%	1/2	0.003
IRC-spybot	4	YES	3	75%	1/2	0.003
IRC-N	259	YES	258	99.6%	0	0
HTTP-1	4	YES	4	100%	1/2	0.003
HTTP-2	4	YES	4	100%	1/2	0.003
P2P-Storm	13	YES	13	100%	0	0
P2P-Nugache	82	YES	82	100%	0	0

MODEL

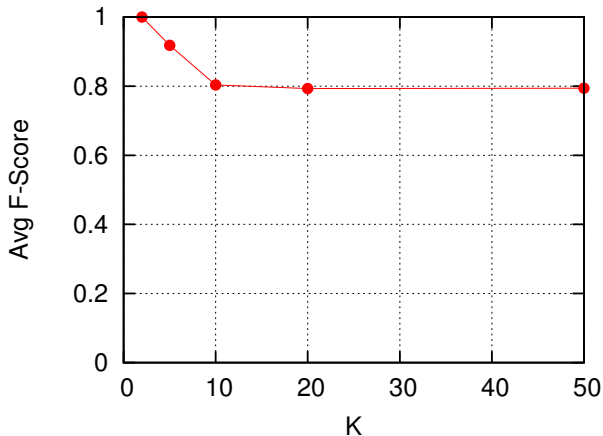
- Network divided in many subnetworks
- Each subnetwork contains one computing node and many monitor nodes
- Monitor nodes observe the subnetwork and send data to computing node.
- Subnetworks do not want to share sensitive data

DISTRIBUTED K-MEANS

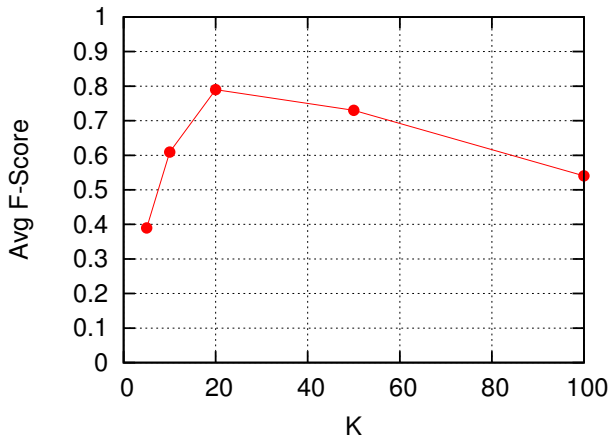
From “Clustering distributed data streams in peer-to-peer environments” (Bandyopadhyay 2006)

- Each node use same starting centroid
- Until nodes or neighbours do not change
 - Map local points to centroid
 - Compute new position of centroids
 - Ask neighbours for their centroids
 - Answer requests if any
 - Replace local centroids with average of centroids received

DISTRIBUTED K-MEANS WITH CORRECT K

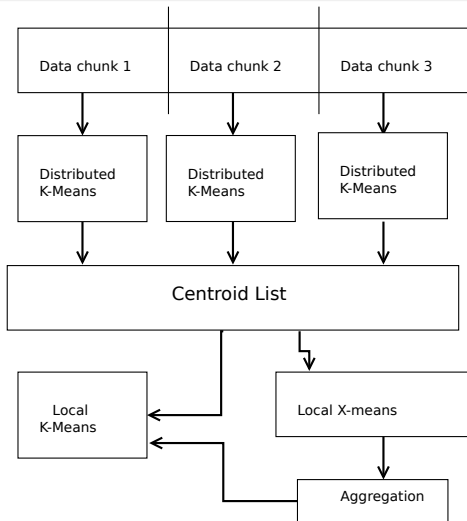


DISTRIBUTED K-MEANS WITH REAL $K=20$

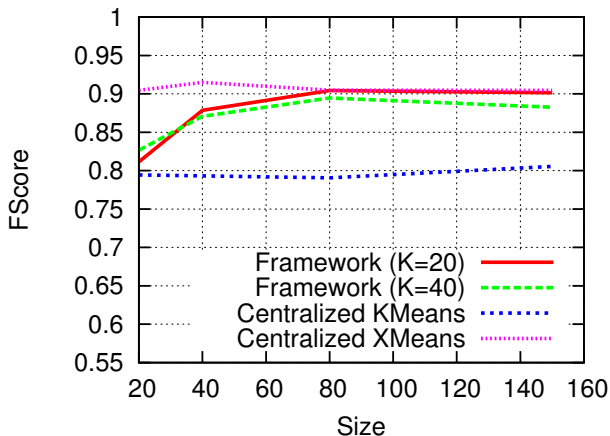


FRAMEWORK IN BRIEF

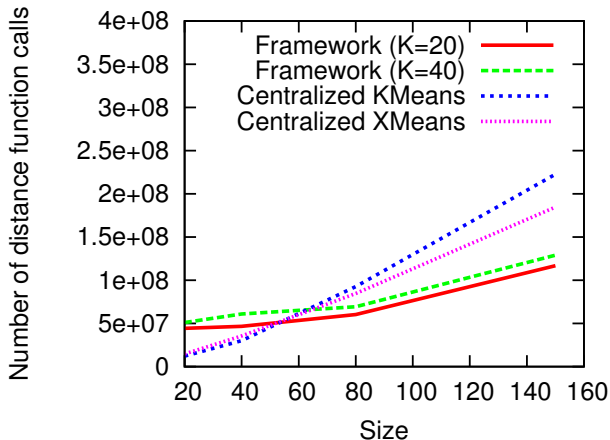
- Each node computes points in A-Plane and C-Plane
- Distributed clustering according to scheme
- Cross clustering Correlation TBD



PRECISION OF FRAMEWORK



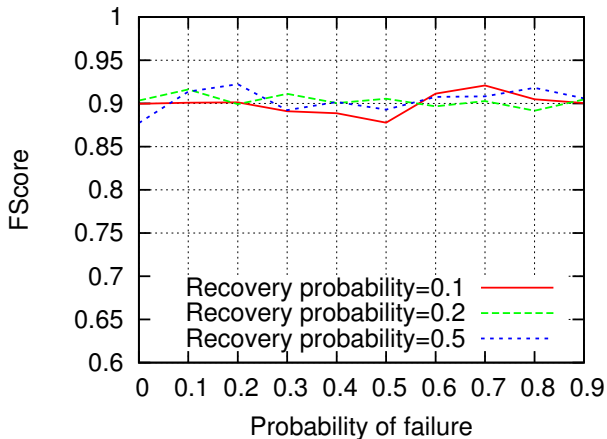
NUMBER OF DISTANCE FUNCTION CALLS



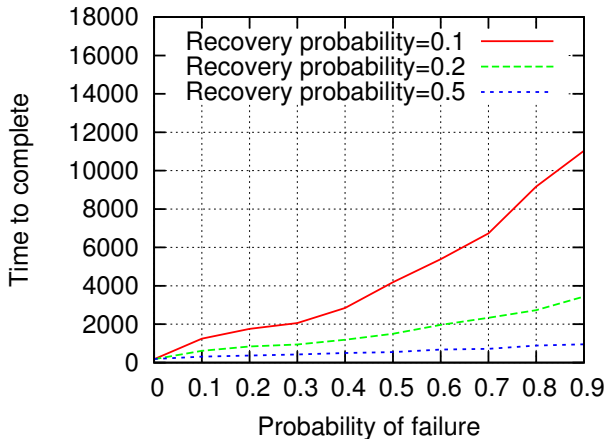
ROBUSTNESS SETTING

- We use two parameters, P_u , P_d
- Every 10 units of time
 - Each up node goes down with prob. P_d
 - Each down node goes up with prob. P_u
- Down nodes do not analyze data or answer messages

ROBUSTNESS



ROBUSTNESS



NEXT STEPS

- Development of Distributed Cross Clustering Correlation
- Testing of framework with real data
- Work on online features