# Sentiment Analysis using Word-Graphs

## SUPER

### Social sensors for secUrity assessments and

### Proactive EmeRgencies management

John Violos, Konstantinos Tserpes, Evangelos Psomakelis, Konstantinos Psychas, Theodora Varvarigou

National Technical University of Athens

# SUPER Project

real time social network mapping

strategic planning and emergency management

behavioral analysis

**SUPER is a joint effort of social media and emergency management experts towards introducing a holistic, integrated and privacy-friendly approach to the use of social media in emergencies and security incidents.**

2

## Hurricane Sandy 2012 Oct 27 – Nov 1:  **20M + Tweets**



**Picture from: www.theworld4realz.com**

3

# Social Sensors

- **Event - SubEvent Detection**

- **Topic Community Tracking**

- **Sentiment Analysis**

- **Behaviour Abalysis**

- **Rumour Spreading Identification – Credibility**

- **Intelligent Fusion and Reasoning**

**Improve-Modify**

**Probabilistic Method Component**

**Deep Learning Component**

**Word Graph Representation Component**

**4**

## What is our principal challenge?

To detect Sentiment Polarities of published textual posts in SN!

Sentiments:

+ Positive

-  Negative

= Neutral

5

## How?

SA Problem ⟶ Text Classification Problem

## What is our Contribution?

- Different Reprsentation Model (No BoW)
- Similarity Measures

6

## Neighborhood of the Words
## Sequence of words
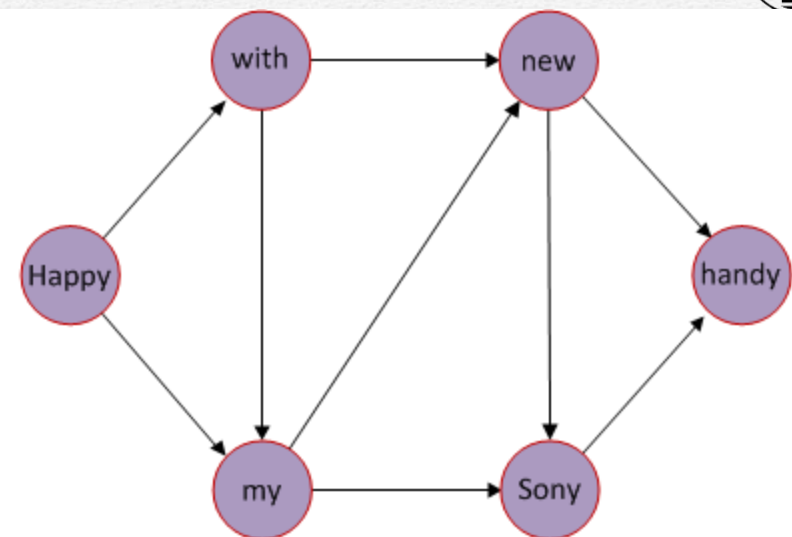
## Example:

\+

The Movie is not boring
I do like it

-

I do not like the Movie
It is boring

## Same Set of Words Different Sentiments
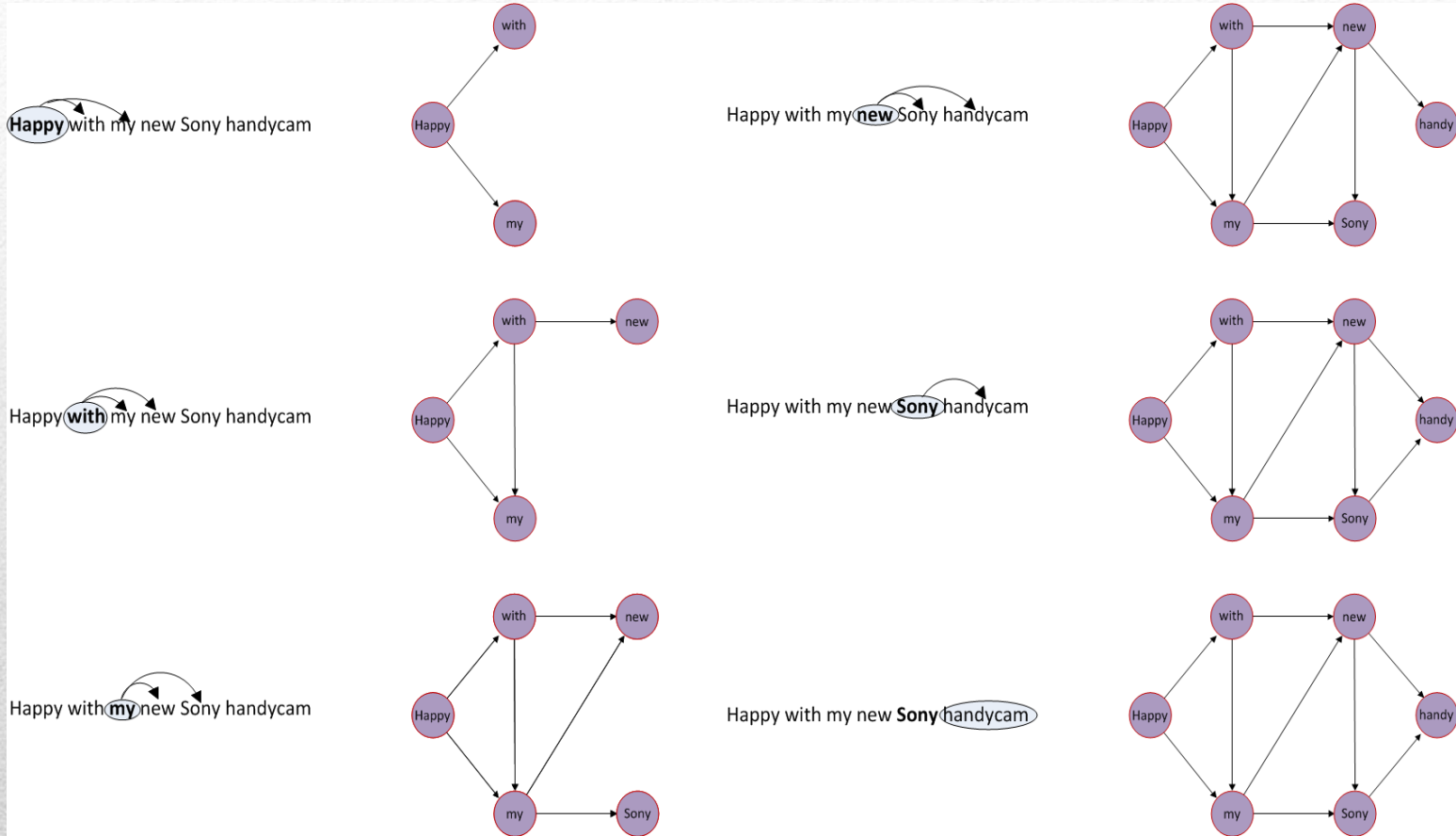## A good solution is the Word Graphs!

7

# Directed & Unweighted Graph

- Nodes ⟶ Words

- Edges ⟶ join neighbor Words

- Vicinity ⟶ Frame N

- Direction ⟶ Sequence of words in the original text

8

Happy with my new Sony handycam



**N=2**

**N=2**

(0.64, 0.08, 0.11)

⬠ The new upcoming tweet

◯ Tweets that have positive polarity

⬛ Tweets that have neutral and negative polarity
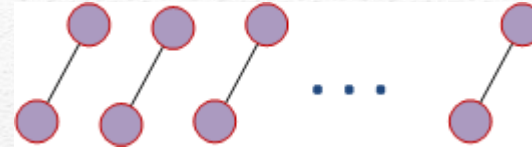
12

- ## Containment Similarity
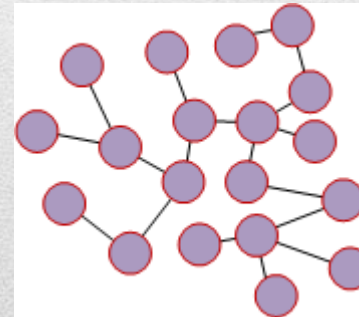  # Common Edges   # Max Edges Normalization

- ## Maximum Common SubGraph
  # Nodes
  # Undirected Edges
  # Directed Edges
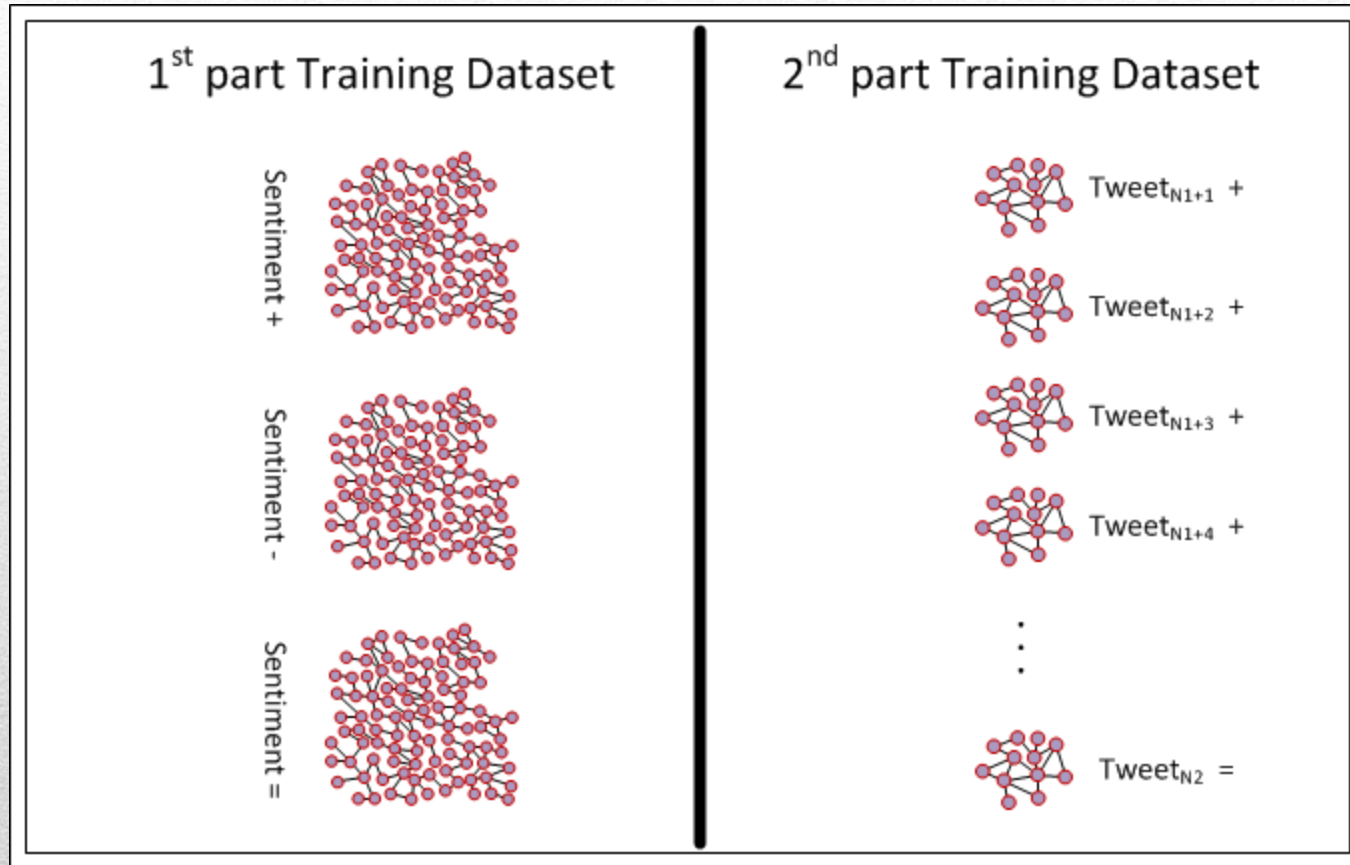
13

Feature Selection techniques:
- Nodes
- Edges

Benefits:
- Improve the Accuracy
- Decrease the Dimensionality & Computation

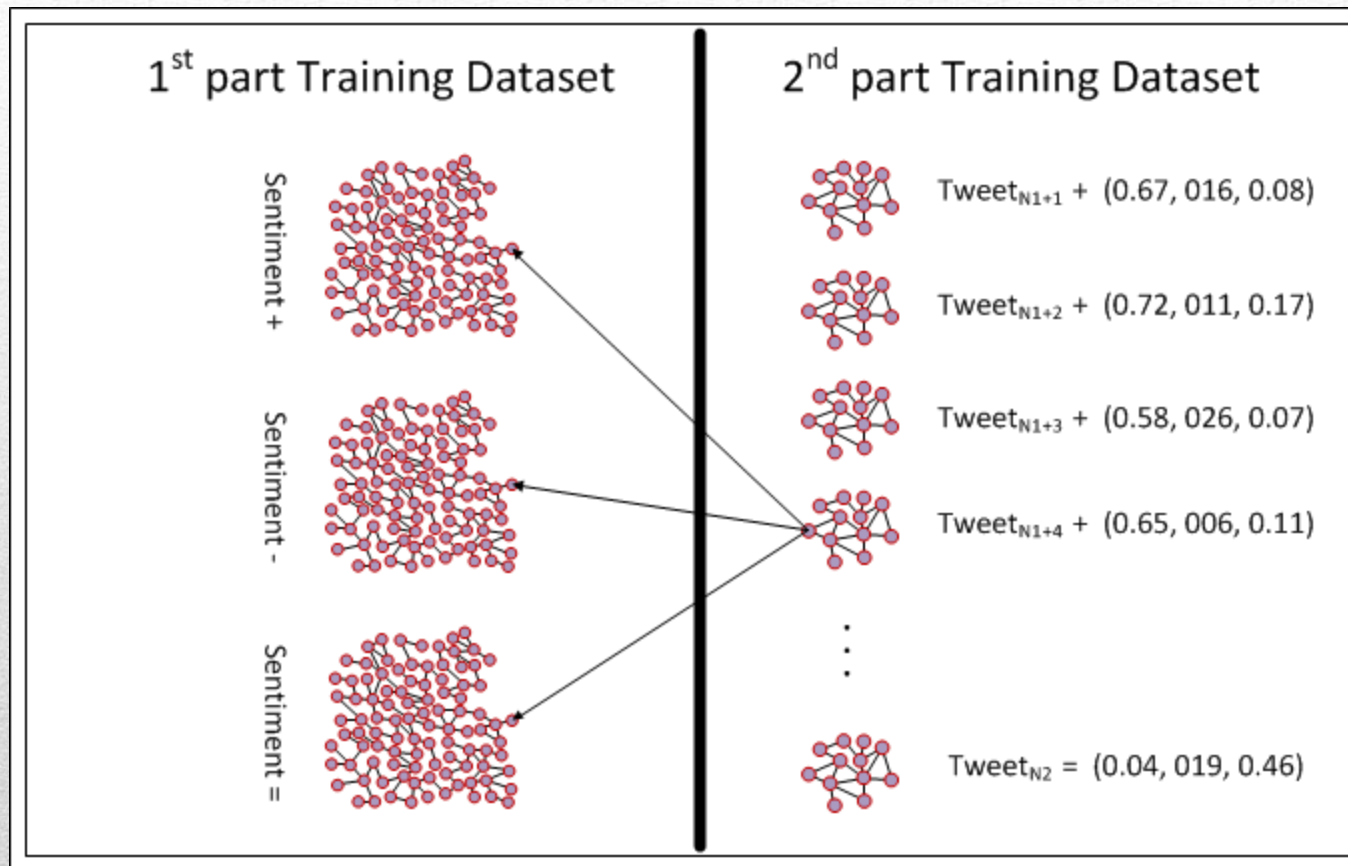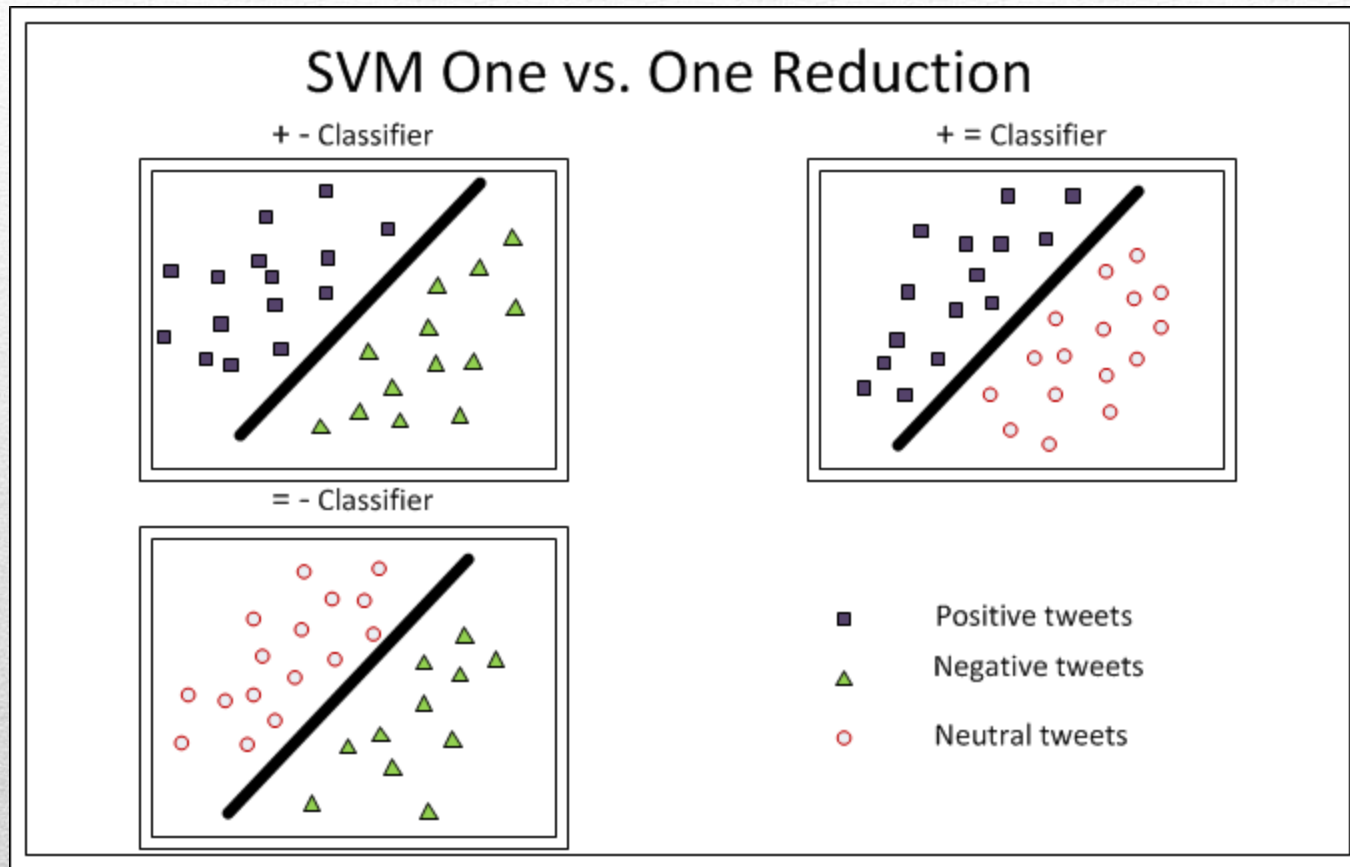Mutual Information:
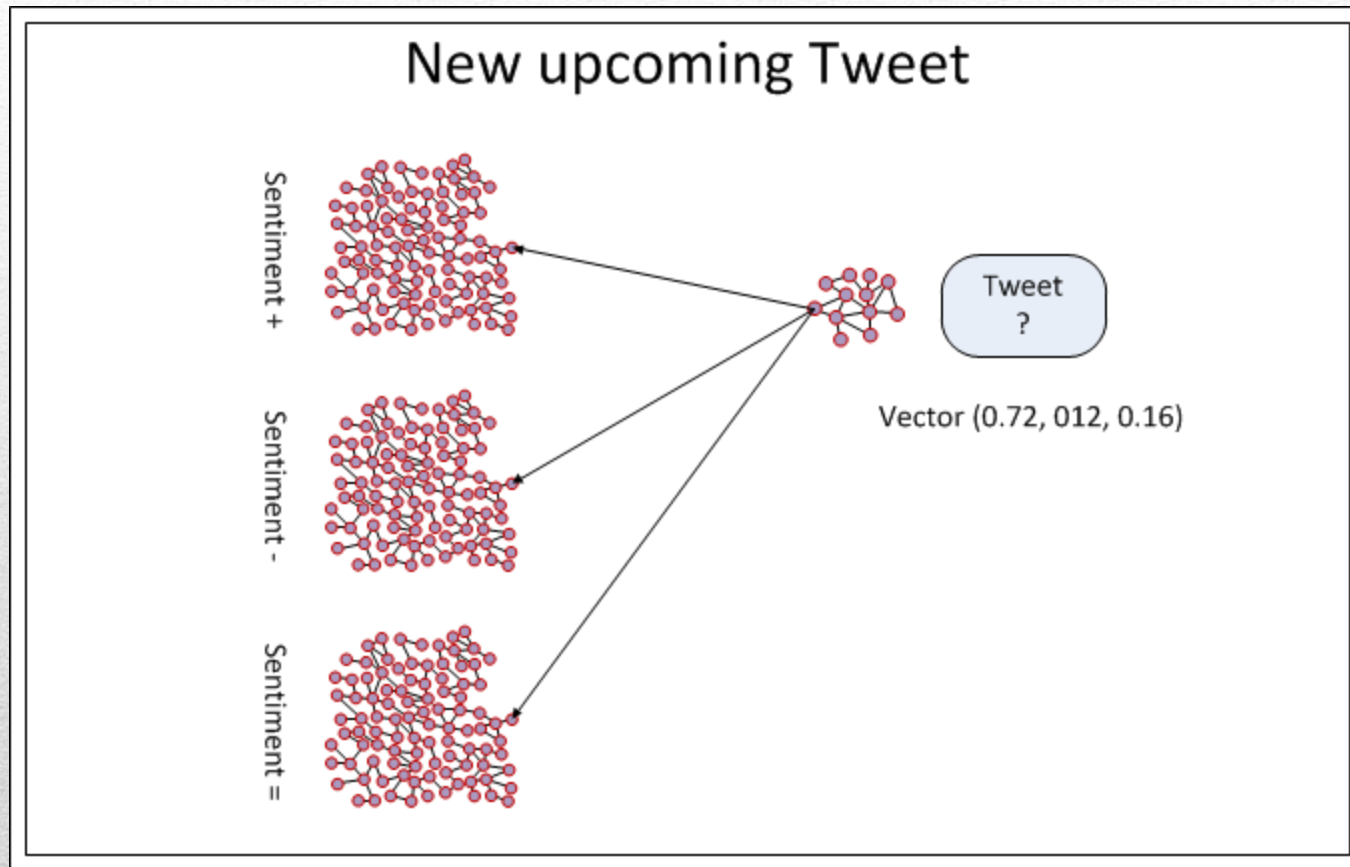The Appropriateness and the Contribution of each Edge for the Accurate Sentiment Prediction + - =

**14**

1st part Training Dataset — 2nd part Training Dataset

1st part Training Dataset: $Tweet_1$ +, $Tweet_6$ -, $Tweet_{11}$ =, $Tweet_2$ +, $Tweet_7$ -, $Tweet_{12}$ =, $Tweet_3$ +, $Tweet_8$ -, $Tweet_{13}$ =, $Tweet_4$ +, $Tweet_9$ -, $Tweet_5$ +, $Tweet_{10}$ -, $Tweet_{N1}$ =

2nd part Training Dataset: $Tweet_{N1+1}$ +, $Tweet_{N1+6}$ -, $Tweet_{N1+11}$ =, $Tweet_{N1+2}$ +, $Tweet_{N1+7}$ -, $Tweet_{N1+12}$ =, $Tweet_{N1+3}$ +, $Tweet_{N1+8}$ -, $Tweet_{N1+13}$ =, $Tweet_{N1+4}$ +, $Tweet_{N1+9}$ -, $Tweet_{N1+5}$ +, $Tweet_{N1+10}$ -, $Tweet_{N2}$ =

16

1st part Training Dataset     2nd part Training Dataset

Sentiment +    Sentiment -    Sentiment =

$Tweet_{N1+1}$ +
$Tweet_{N1+2}$ +
$Tweet_{N1+3}$ +
$Tweet_{N1+4}$ +
$Tweet_{N2}$ =

17

SVM One vs. One Reduction

New upcoming Tweet

Sentiment +

Sentiment -

Sentiment =

Tweet ?

Vector (0.72, 012, 0.16)

Classification of New Upcoming Tweet

Public available Dataset from the research Language-Independent Twitter Sentiment Analysis of Sascha, Hufenhaus and Albayark

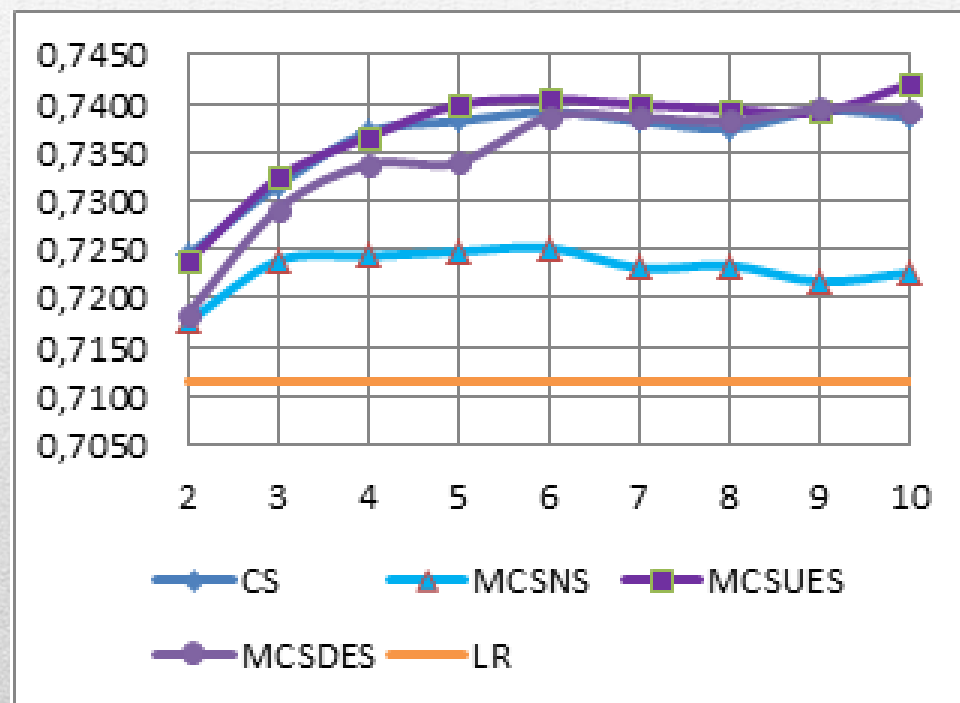The dataset is manually annotated by workers on Amazon's Mehanical Turk

10594 tweets

+ 2334

- 1486

= 6774

22

# SA Accuracy with SOA Methods

| Method | 4Gram | 4Gram Graphs |
|---|---|---|
| Bayesian Network | 0.6788 | 0.6791 |
| C4.5 | 0.6828 | 0.6896 |
| SVM | 0.6777 | 0.6847 |
| **Logistic Regression** | 0.6822 | **0.7115** |
| Simple Logistic Regression | 0.6816 | 0.7109 |
| Multi-Layer Perceptron | 0.6788 | 0.7069 |
| Best-First Tree | 0.6790 | 0.6840 |
| Functional Tree | 0.6822 | 0.7079 |

## 10-Fold Cross Validation
## Various size of N (2 - 10) Frame Window
## Various Graph Similarity Metrics

## No Graph Filering
## SVM Classifier


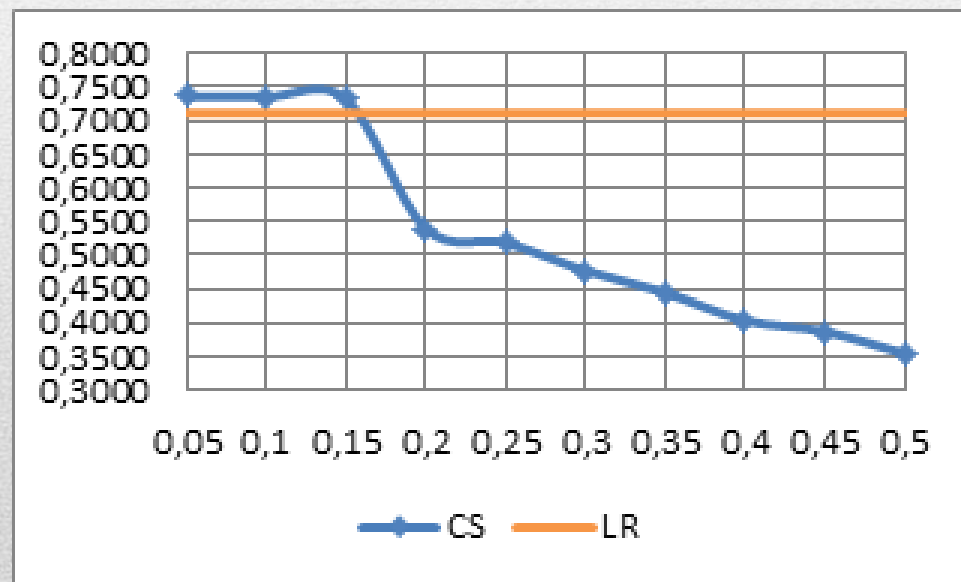
**24**

## 10-Fold Cross Validation
## Four-Words Frame size
## Mutual Information on Edges as Feature Selection technique
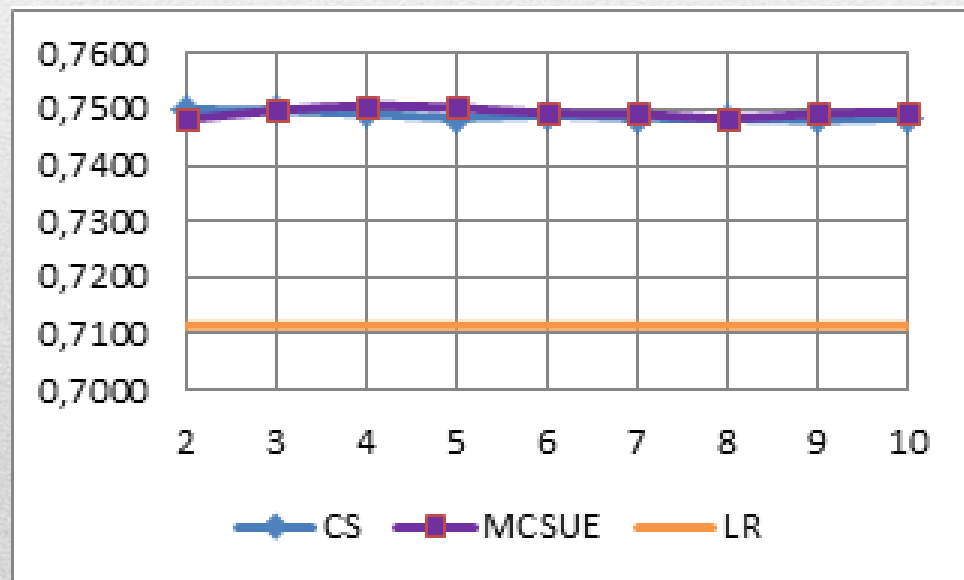## MCS Graph similarity techniques in combination with MI deteriorates the Accuracy of the method

### Contanimen Similarity
### SVM Classifier

## 10-Fold Cross Validation
## Various size of N (2 – 10) Frame Window
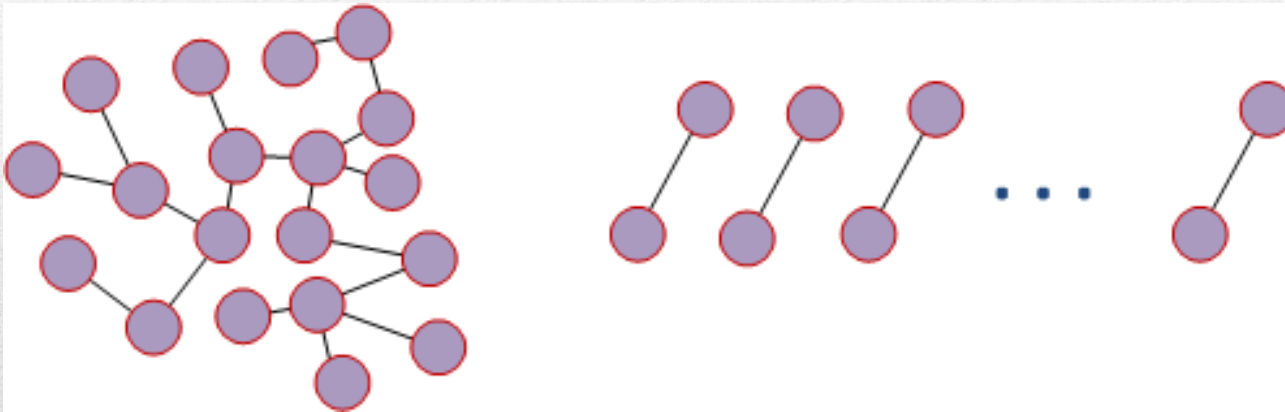## Containment Similarity
## MCSUE Similarity

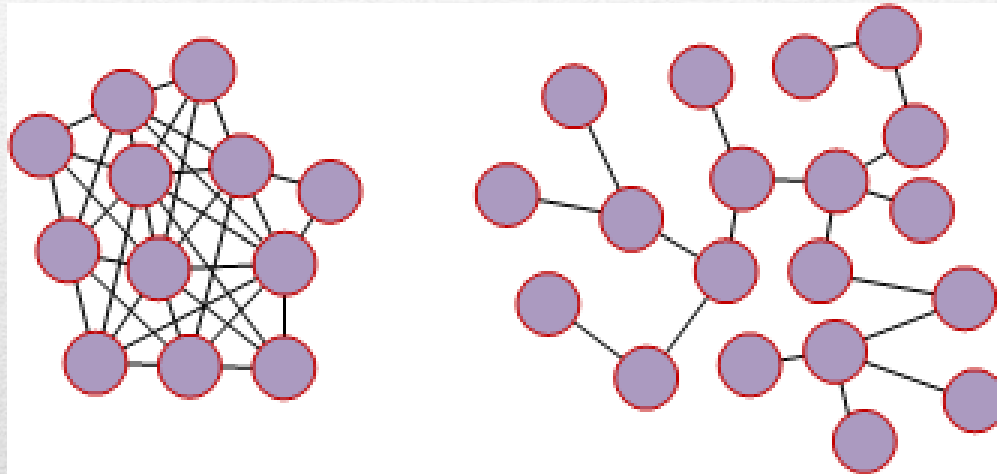## Gaussian Bayes classifier
## No Graph filtering

**WGSA:**

- **Vicinity & sequence of the words-terms.**
- **Writing characteristics & Partial Matching.**
- **Polysemy**
- **Evolution of the Annotated Dataset**
- **The discard of the edges with low MI has as a result to separate the MCS that captures the sentiment.**

27

# MCS>CS



**A correlation among a few common words is more important than the existance of many common word pairs.**

28

# MCSUES>MCSNS



**The Dense Correlation among few words is more important than the sparse correlation between many words.**

29

- **Weighted Graphs**

- **Graph Similarity metrics**

- **Targeted Sentiment Analysis**

- **More Sentiments**

- **Different Training Dataset to Testing Dataset**

30

real time social network mapping

strategic planning and emergency management

behavioral analysis

# Thank you for your attention!