

# From NLP to MLP

Peter Teufel/Udo Payer/Günther Lackner  
Graz University of Technology, Austria



# What to expect...

- Knowledge mining in various areas:
  - Event correlation, RDF data analysis, WIFI privacy, malware analysis
  - ... and e-Participation
- Common model based on ML and AI

Intro

Idea

MLP vs.  
NLP

Lexical  
Parser

POS(C)  
Tagging

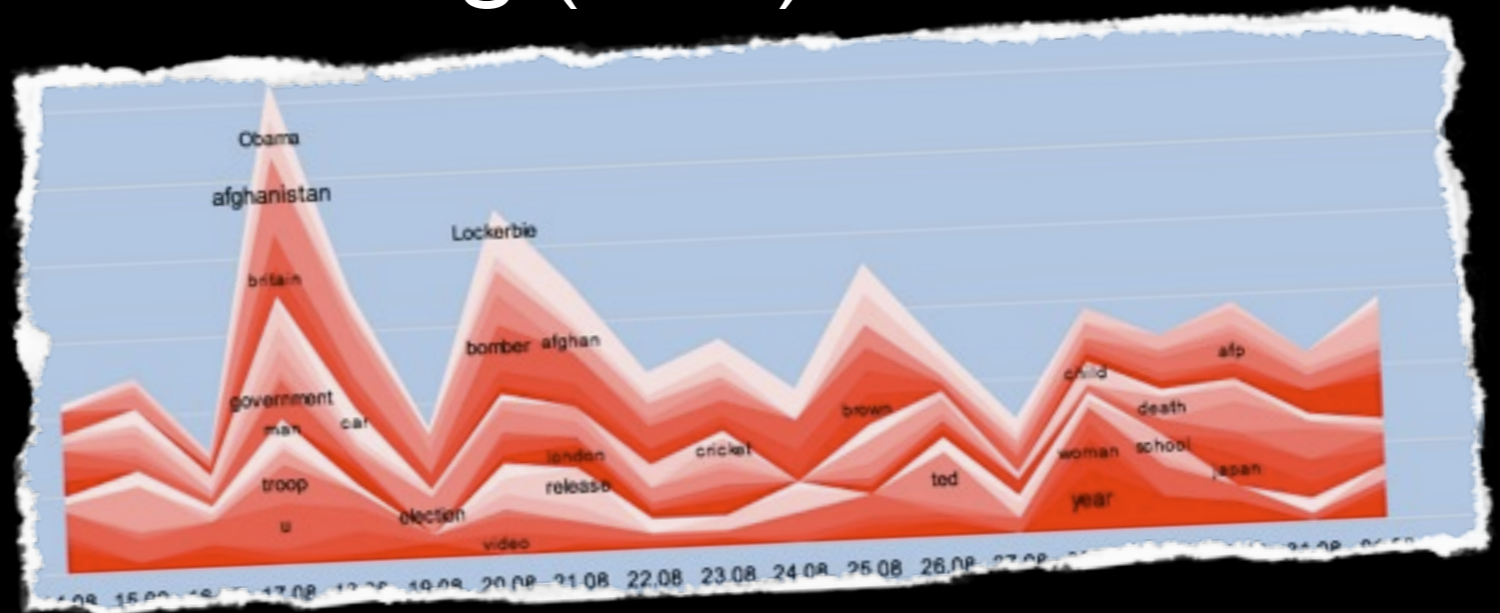
Lemmatization

Activation  
Patterns

Analysis

# e-Participation

- Natural Language Processing (NLP)
- Analysis of text
  - Clustering
  - Relations
  - Semantic aware search



Intro

Idea

MLP vs.  
NLP

Lexical  
Parser

POS(C)  
Tagging

Lemmatization

Activation  
Patterns

Analysis

# Idea...

- In previous work, analysis of polymorphic shellcodes
- Neural networks for the detection of shellcode decryption loops
- Code vs. natural language?

Intro

Idea

MLP vs.  
NLP

Lexical  
Parser

POS(C)  
Tagging

Lemmatization

Activation  
Patterns

Analysis

# NLP

- Natural Language Processing
- Various techniques are available
- Sentence analysis
- POS tagging
- Word sense disambiguation

Intro

Idea

MLP vs.  
NLP

Lexical  
Parser

POS(C)  
Tagging

Lemmatization

Activation  
Patterns

Analysis

# MLP

- Machine Language Processing
- In this talk: focus on assembler, but applicable to any other language
- Assuming that machine language has similar properties as real language
- Semantics, grammar, (word/instruction) disambiguation

Intro

Idea

MLP vs.  
NLP

Lexical  
Parser

POS(C)  
Tagging

Lemmatization

Activation  
Patterns

Analysis

# MLP

- Malware: Signatures? Understanding?
- Fuzzy analysis
  - Group code with similar behavior
  - Semantic search for similar code
  - Semantic relations within code

Intro

Idea

MLP vs.  
NLP

Lexical  
Parser

POS(C)  
Tagging

Lemmatization

Activation  
Patterns

Analysis

# Goal

MLP vs. NLP

- Can we map NLP to MLP???
- NLP research highly developed, mature frameworks, techniques
- Problems in language analysis easier to spot, since we know our language very well

Intro

Idea

MLP vs.  
NLP

Lexical  
Parser

POS(C)  
Tagging

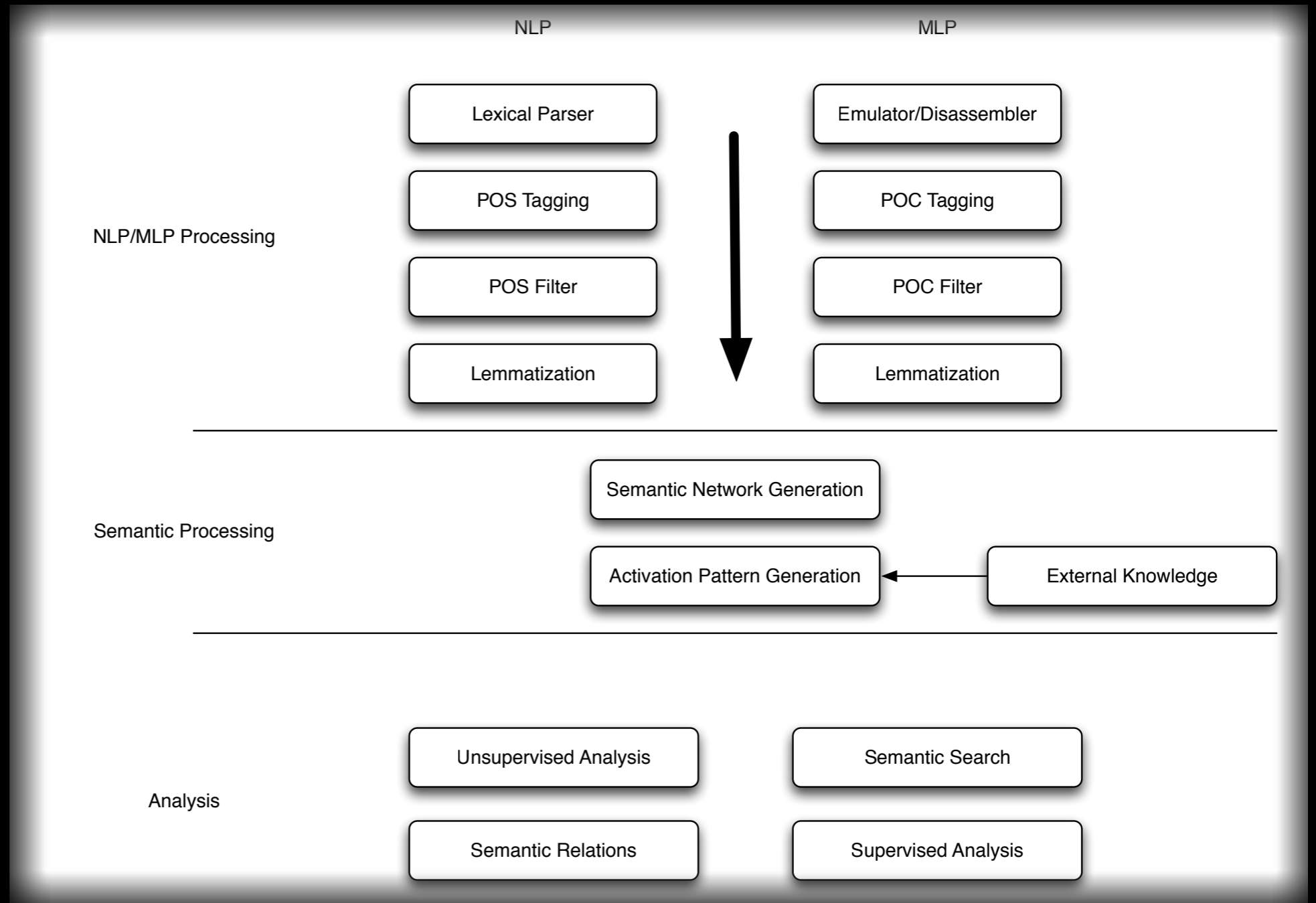
Lemmatization

Activation  
Patterns

Analysis



# Layers



Intro

Idea

MLP vs. NLP

Lexical Parser

POS(C) Tagging

Lemmatization

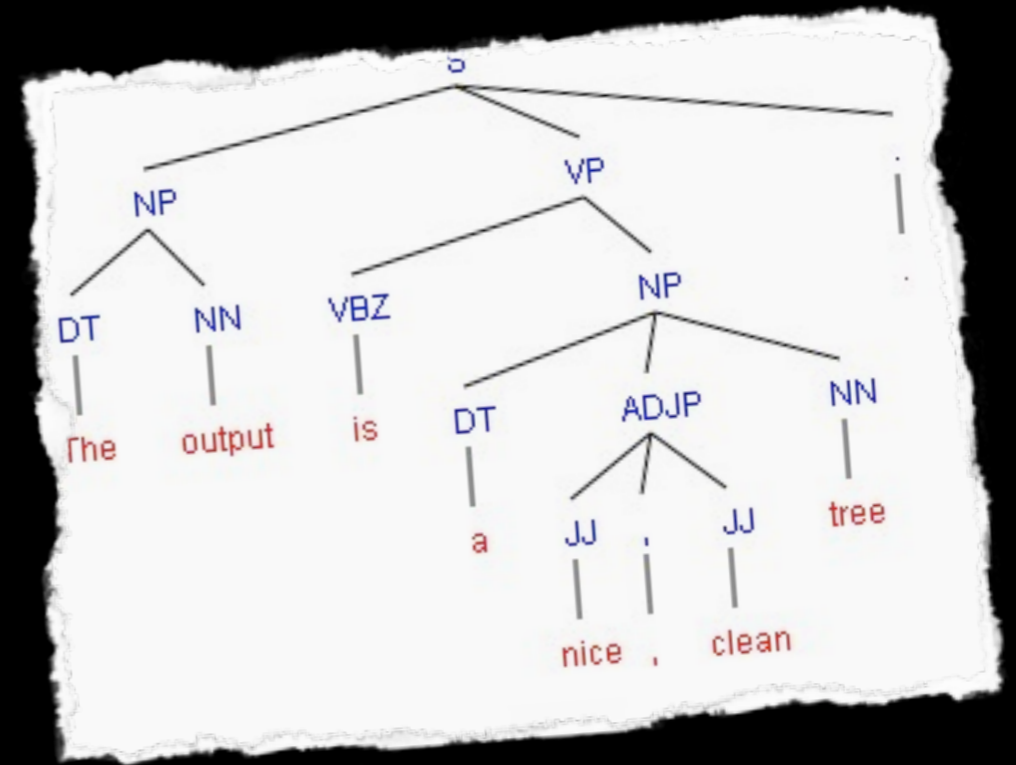
Activation Patterns

Analysis

Unsupervised Analysis

Supervised Analysis

# Lexical parser



- NLP:
  - Extract sentences, analyze terms, find relations (e.g. Stanford parser)
  - This beautiful\_A city\_N is\_V called\_V St. Petersburg\_SN.

Intro

Idea

MLP vs.  
NLP

Lexical  
Parser

POS(C)  
Tagging

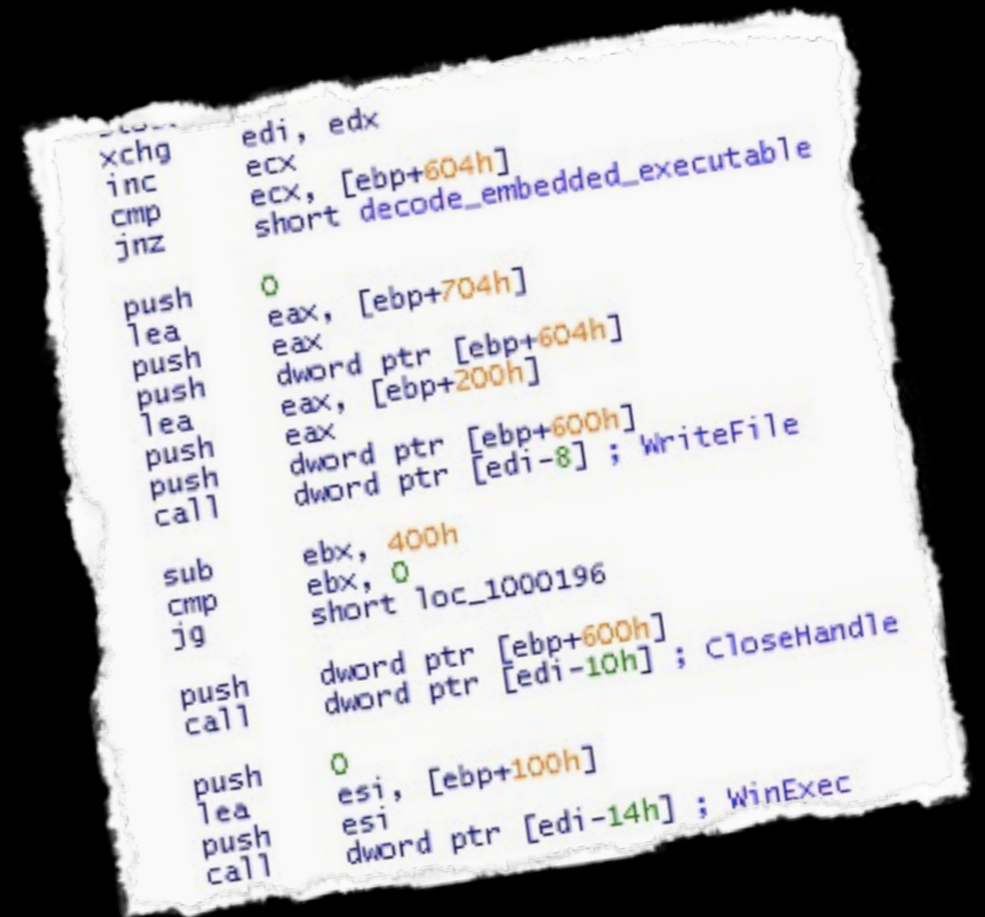
Lemmatization

Activation  
Patterns

Analysis

# Lexical parser

- MLP:
  - instruction sequences (mov, sub, xor...)
  - relations between typical instructions (modifying the same register, variable)
  - e.g. interrupt preparation, loops e.g.



# POS tagging, filter

- NLP:
  - This beautiful\_A city\_N is\_V called\_V St. Petersburg\_SN.
  - Tagging terms
  - Filtering (stop words etc.): beautiful, city, St. Petersburg

Intro

Idea

MLP vs.  
NLP

Lexical  
Parser

POS(C)  
Tagging

Lemmatization

Activation  
Patterns

Analysis

# POC tagging, filter

- MLP:
  - jmp, call, jz... (branch type)
  - add, sub... (arithmetic)
  - xor... (logic)

Intro

Idea

MLP vs.  
NLP

Lexical  
Parser

POS(C)  
Tagging

Lemmatization

Activation  
Patterns

Analysis

# Lemmatization

- NLP:
  - bought => buy
  - cars => car
- MLP:
  - drop parameters (mov ax,4)
  - group instructions (arithmetic, logic)

Intro

Idea

MLP vs.  
NLP

Lexical  
Parser

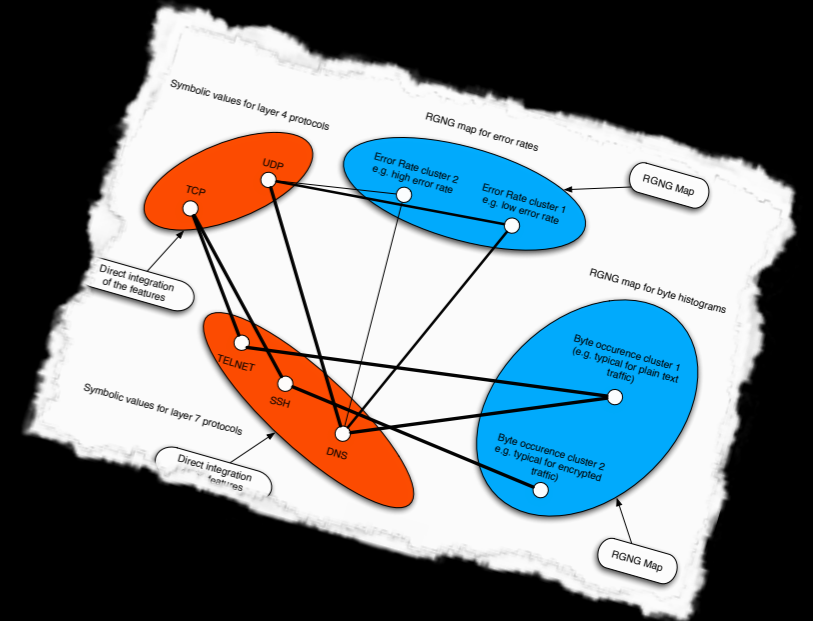
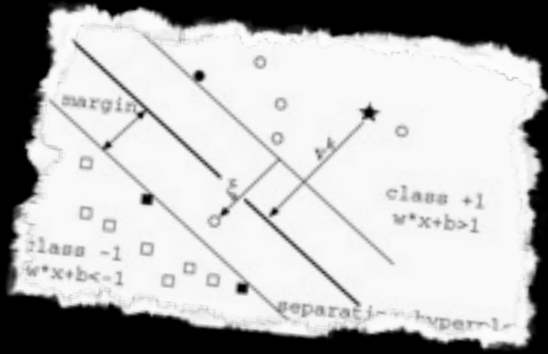
POS(C)  
Tagging

Lemmatization

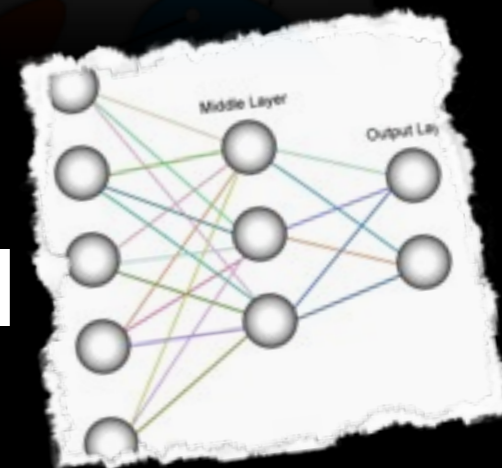
Activation  
Patterns

Analysis

# Activation Patterns

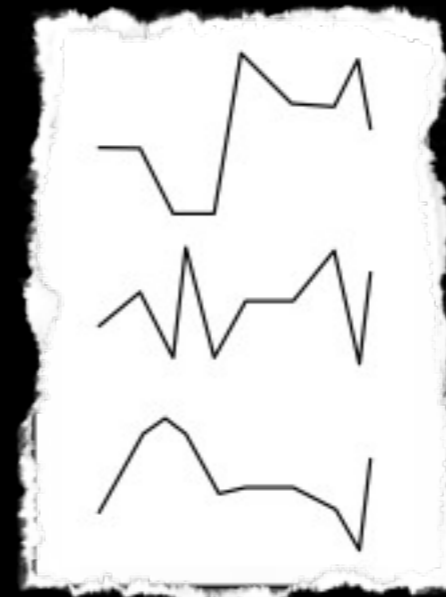


- Based on semantic networks, spreading activation, machine learning
- Allows us to analyze arbitrary combination of features (symbolic, real values)
- Patterns are the basis for a wide range of analysis methods



# Analysis

- Unsupervised learning (clustering)
- Supervised learning (classifiers)
- Semantic search
- Semantic relations
- Anomaly detection
- Feature relevance



Intro

Idea

MLP vs.  
NLP

Lexical  
Parser

POS(C)  
Tagging

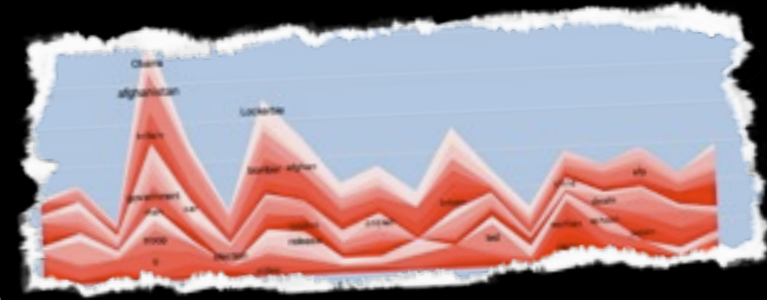
Lemmatization

Activation  
Patterns

Analysis



# Applications



RDF Analysis

Event Correlation

Twitter Mining

Malware  
Analysis

e-Participation  
Text analysis

User Tracking in  
WIFI Networks

Intro

Idea

MLP vs.  
NLP

Lexical  
Parser

POS(C)  
Tagging

Lemmatization

Activation  
Patterns

Analysis

# MLP Example

- Metasploit shellcodes
- Decoder loops of various decryption engines
- Clustering, semantic search and relations

Intro

Idea

MLP vs.  
NLP

Lexical  
Parser

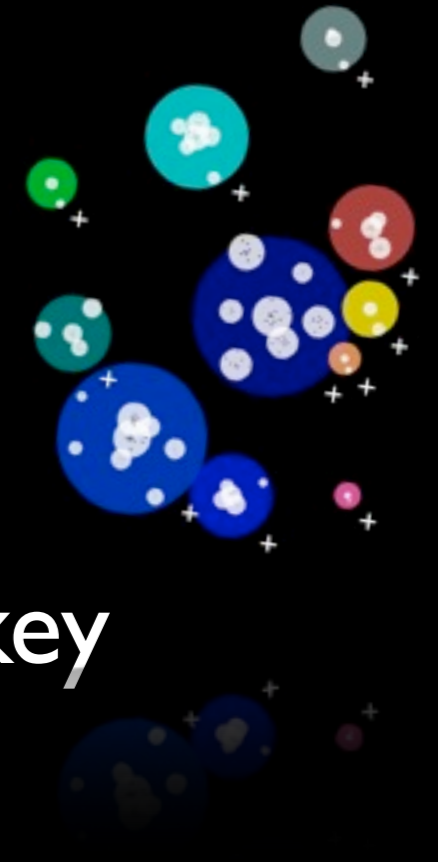
POS(C)  
Tagging

Lemmatization

Activation  
Patterns

Analysis

# Clustering



- NLP:
  - Group similar documents, extract key concepts from clusters, gain quick overview
- MLP:
  - Group similar code (instruction sequences, functions, decryption loops)

# Semantic Search

- NLP:
  - St. Petersburg is beautiful. The city was founded in 1703.

| Result | Decoder        | Instruction chain                        | Description   |
|--------|----------------|--|---------------|
| 1      | shikata-ga-nai | xor add add loop                         | Decoder       |
| 2      | shikata-ga-nai | xor mov fnstenv pop mov xor add add loop | Decoder setup |
| 3      | nonalpha       | pop mov add mov cmp jge                  | Decoder setup |
| 4      | fnstenv-mov    | xor sub loop                             | Decoder       |
| 5      | countdown      | xor loop                                 | Decoder       |

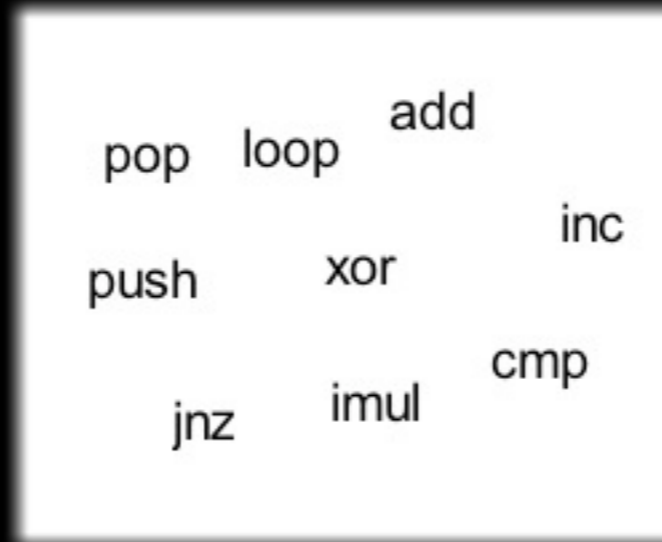
- MLP:
  - Search for similar concepts (decryption loops)

# Semantic Relations

- NLP



- MLP



# Thoughts...

- Same basic principles for NLP and MLP
- Use the existing knowledge, bring it to another domain
- Apply it to arbitrary language
- Understand malware/programs???

Thank you for your  
attention!