



IT Systems Engineering | Universität Potsdam

# Natural Language Processing

*Relation Extraction*

Potsdam, 14 June 2012

**Saeedeh Momtazi**  
Information Systems Group

based on the slides of the course book

# Outline

2

- 1 Introduction
- 2 Task
- 3 Pattern Extraction
- 4 Supervised Learning
- 5 Semi-supervised Learning

# Outline

3

- 1 Introduction
- 2 Task
- 3 Pattern Extraction
- 4 Supervised Learning
- 5 Semi-supervised Learning

# Information Extraction

- Named entity recognition
  
- Relation Extraction

# Named Entity Recognition

5

HPI is affiliated to the Potsdam University and located in Potsdam near Berlin. It was founded in 1998 by Hasso Plattner, one of the co-founders of the European software company, SAP AG.

ORG	HPI
ORG	Potsdam University
LOC	Potsdam
LOC	Berlin
DATE	1998
PER	Hasso Plattner
ORG	SAP AG

# Relation Extraction

6

HPI is affiliated to the Potsdam University and located in Potsdam near Berlin. It was founded in 1998 by Hasso Plattner, one of the co-founders of the European software company, SAP AG.

located (ORG-LOC)	HPI - Potsdam
near (ORG-LOC)	HPI - Berlin
near (LOC-LOC)	Potsdam - Berlin
founded (ORG-DATE)	HPI - 1998
founder (ORG-PER)	HPI - Hasso Plattner
co-founder (ORG-PER)	SAP AG - Hasso Plattner

# Information Extraction

7



- Home page
- Contents
- Random content
- Current events
- Random article
- Random in language
- Interaction
- Help
- About Wikipedia
- Community portal
- Recent changes
- Contact Wikipedia
- Toolbox
- Print/export
- Languages
- Wikisource
- Wikibooks
- Wikiquote
- Wikispecies
- Wikiversity
- Wikivoyage
- Wiktionary
- Wikimedia Commons
- Wikimedia Library
- Wikimedia Manual
- Wikimedia News
- Wikimedia Reference desk
- Wikimedia Reports
- Wikimedia Wikiquote
- Wikimedia Wikisource
- Wikimedia Wikispecies
- Wikimedia Wikiversity
- Wikimedia Wikivoyage
- Wikimedia Wiktionary
- Wikimedia Commons
- Wikimedia Library
- Wikimedia Manual
- Wikimedia News
- Wikimedia Reference desk
- Wikimedia Reports

**Potsdam**  
 From Wikipedia, the free encyclopedia

For other uses, see **Potsdam** (disambiguation).  
**Potsdam** (German pronunciation: [ˈpɔʦtəm]) is the capital city of the German federal state of Brandenburg. The name "Potsdam" originally seems to have been "Plothogin" from a West Slavic Polesian but several claims to national and international notability. In Germany, it is of Sarmatian, the largest stone heritage site in Germany. The Potsdam Conference and Babelsberg, in the south-eastern part of Potsdam, was a major film production studio. Potsdam developed into a centre of science in Germany from the 19th century. The

- Contents** [hide]
- 1 Geography
  - 2 History
  - 3 Governance of Potsdam
  - 4 Potsdam in the film industry
  - 5 Parks
  - 6 Administration
  - 7 International relations
  - 8 Twin towns – sister cities
  - 9 Education and research
  - 10 Water rights
  - 11 Transport
  - 12 Sport
  - 13 Famous people
  - 14 References
  - 15 Sources
  - 16 External links

**Geography**



The area was formed from a series of *Tauern*. See *Angermünde*, *Talsperren*. Potsdam is divided into seven historical districts by larger areas of former buildings.



**History**

The area around Potsdam shows clear signs of first mentioned in a document in 99 from around 1246. In 1273, it was all a continuous Hohenzollern possession: Potsdam baronies.



After the Elect of Prussia in 1665, Potsdam. Later, the city became a full residence in 1744, leased for its formal gardens and in 1815, at the formation of the Prussia then upgrade Potsdam and Potsdam to



**Governments of Potsdam**

Between 1815 and 1945 the city of Potsdam situated between *Müritzer* and the merged to urban district into the govern ment by the GDR. The southwestern has, capital.<sup>[?]</sup>





**Coordinates** ⊕ **52°24′0″N 13°40′E**

Administration	
<b>Country</b>	<b>Germany</b>
<b>State</b>	<b>Brandenburg</b>
<b>District</b>	<b>Urban district</b>
<b>Lord Mayor</b>	<b>Jann Jakobs (SPD)</b>

Basic statistics	
<b>Area</b>	<b>187.28 km<sup>2</sup> (72.31 sq mi)</b>
<b>Elevation</b>	<b>35–114 m</b>
<b>Population</b>	<b>156,906 (31 December 2010)<sup>[1]</sup></b>
<b>- Density</b>	<b>838 /km<sup>2</sup> (2,170 /sq mi)</b>

**Other information**

<b>Time zone</b>	<b>CET/CEST (UTC +1/+2)</b>
<b>Licence plate</b>	<b>P</b>
<b>Postal codes</b>	<b>14401–14482</b>
<b>Area code</b>	<b>0331</b>
<b>Website</b>	<b><a href="http://www.potsdam.de">www.potsdam.de</a></b> <span style="font-size: 1.2em;">ⓘ</span>

Log in / create account

Home [Edit](#) [View history](#)

Coordinates: ⊕ **52°24′0″N 13°40′E**

major cultural landmarks, in particular the parks and palaces

go back to this article in the world.



islands, such as the Havel, the Grabelsberg, Tegelberg, See,

many of historical buildings, the south of the city is dominated

entirely as a settlement of the Havelbar situated on a castle i site. By 1317, it was mentioned as a small town, it gained its

a later became the Kingdom of Prussia, it also housed

population growth and economic recovery.

is "without center", by Georg Bräunlein von Kriebitzschdorf,

0. The province comprised two governorates named after

of Preußisch, and the greater part of the *Middle March*. It was the governorate of Potsdam between 1815 and 1822, then it was transferred by the *North Elbe* and the *Havel*, and on the north and southwest, divided into *Prussian* royal districts, named after their





# Motivation

- Creating new structured data sources (knowledge bases)
  - DBPedia
  - Freebase
  - Yago
  
- Answering complex questions using multiple sources

Which soccer player married a Spice Girls star?

("?x" is-a "soccer player")  
("?x" married "?y")  
("?y" member "Spice Girls")

# Outline

10

- 1 Introduction
- 2 Task**
- 3 Pattern Extraction
- 4 Supervised Learning
- 5 Semi-supervised Learning

# Relation Representation

- Representing data as triples

(Argument1    RelationType    Argument2)  
(Subject       Predicate       Object)

- Resource Description Framework (RDF)

# Relation Types

- Having various relation types based on the type of arguments
  - PER-PER
  - ORG-LOC
  - PER-ORG
  - PER-LOC
  - PER-DATE
  
  - Drugs-Disease

# Relation Types



Spouse  
Parent  
Child  
Friendship  
Colleague

..



# Relation Types

14



Place of birth  
Place of death  
Buried in

..



# Relation Types

15



Founder  
Co-founder  
Owner  
Employee  
Student/Alum  
Professor

..



# Relation Types

16



Located  
Near  
Founded-location  
Headquarter

..



# Approaches

- Manually created patterns
- Supervised machine learning
- Semi-supervised learning

# Outline

18

- 1 Introduction
- 2 Task
- 3 Pattern Extraction**
- 4 Supervised Learning
- 5 Semi-supervised Learning

# Pattern Extraction

- 19
- What are the potential words to express a relation type?  
(PER Member ORG)  
("?x" Member "?y")

x is a member of y.  
x is an employee of y.  
x works at y.  
x is a staff of y.

x is (a|an) (member|employee|staff) of y.  
x (works) at y.

# Pattern Extraction

20

- Advantages

- Having high precision results

- Disadvantages

- Having low recall
- Finding all possible patterns is labor intensive
- Covering all relations is very difficult

# Outline

21

- 1 Introduction
- 2 Task
- 3 Pattern Extraction
- 4 Supervised Learning**
- 5 Semi-supervised Learning

# Supervised Classification

22

- Training data
  - Defining a fix set of relation types
  - Choosing the corresponding named entities
  - Selecting a set of text as training data
  - Recognizing the named entities in the text
  - Labeling the relations between named entities manually

# Task

23

## ■ Input

- A pair of entities
- A context that this pair appears in
- Possible relation types

## ■ Output

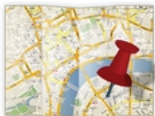
- Type of relation between two entities, if there exist any



Thomas Edison

Thomas Edison died on October 18, 1931, in New Jersey due to complications of diabetes.

Place of birth  
 Place of death ✓  
 Buried in  
 ..



New Jersey

# Feature Selection

24

Thomas Edison died on October 18, 1931, in New Jersey due to complications of diabetes.

- The target entities
  - $T1$ : Thomas Edison
  - $T2$ : New Jersey
  
- Surrounding words of target entities
  - $T1_{+1}$ : died
  - $T2_{-1}$ : in
  - $T2_{+1}$ : due
  
- All words between the target entities (bag-of-word)
  - died
  - on
  - October
  - 18
  - ,
  - 1931
  - ,
  - in

# Feature Selection

25

Thomas Edison died on October 18, 1931, in New Jersey due to complications of diabetes.

## ■ Other Features

- The named entity label of the target words
  - $NE(T1)$ : PER
  - $NE(T2)$ : LOC

## ■ The syntactic structure of the words between target words

# Classification Algorithm

26

- Applying any kinds of classifiers
  - $K$  Nearest Neighbor
  - Support Vector Machines
  - Naïve Bayes
  - Maximum Entropy
  - Logistic Regression
  - ...

# Supervised Classification

27

## ■ Advantages

- Very good performance if
  - Having enough training data
  - Having test data similar to training data

## ■ Disadvantages

- Manual labeling of training data is labor expensive
- Difficult to get good results for other genres

# Outline

28

- ① Introduction
- ② Task
- ③ Pattern Extraction
- ④ Supervised Learning
- ⑤ Semi-supervised Learning

# Semi-supervised Learning

29

- Having no large training data
- Producing a small training data (seed data)
  - A set of triples
- Using the seed data to find further entity pairs with the same relation

Bootstrapping

# Bootstrapping

30

- Using the collected seed data
- Finding sentences which contain at least one entity pairs
- Extracting the common contexts of the pair
- Creating patterns from the extracted context
- Using the pattern to grep more pairs and add them to seed data

# Bootstrapping

31

- Using the collected seed data

( Thomas Edison Spouse Mina Mille)

# Bootstrapping

32

- Using the collected seed data
- Finding sentences which contain at least one entity pairs

Thomas Edison married Mina Mille.

Edison married a young woman named Mina Mille.

In 1871, Thomas Edison married Mina Mille.

Thomas Edison marries Mina Mille on December 25.

# Bootstrapping

33

- Using the collected seed data
- Finding sentences which contain at least one entity pairs
- Extracting the common contexts of the pair
- Creating patterns from the extracted context

Thomas Edison married Mina Mille.

Edison married a young woman named Mina Mille.

In 1871, Thomas Edison married Mina Mille.

Thomas Edison marries Mina Mille on December 25.

# Bootstrapping

34

- Using the collected seed data
- Finding sentences which contain at least one entity pairs
- Extracting the common contexts of the pair
- Creating patterns from the extracted context
- Using the pattern to grep more pairs and add them to seed data

( Albert Eistein Spouse “?” )

Einstein marries his cousin Elsa Löwenthal on June 2.

Einstein married Elsa Löwenthal in Berlin.

Einstein married Elsa Löwenthal on 2 June 1919.

After their divorce in 1919, Einstein married Elsa Löwenthal in the same year.

Albert Einstein was married to Elsa Löwenthal for 17 years.

Einstein marries Elsa Löwenthal.

In the same year Albert Einstein married Elsa Löwenthal.

⇒ ( Albert Eistein Spouse Elsa Löwenthal )

# Bootstrapping

35

- Using the collected seed data

( Thomas Edison Spouse Mina Mille)

( Albert Eistein Spouse Elsa Löwenthal )

# Bootstrapping

36

- Using the collected seed data
- Finding sentences which contain at least one entity pairs
- Extracting the common contexts of the pair
- Creating patterns from the extracted context

Albert Einstein's wife, Elsa Löwenthal, was his first cousin.  
Elsa Löwenthal was the wife of Albert Einstein.  
Einstein's wife was named Elsa Löwenthal.

# Further Reading

37

- Speech and Language Processing
  - Chapters 22.2