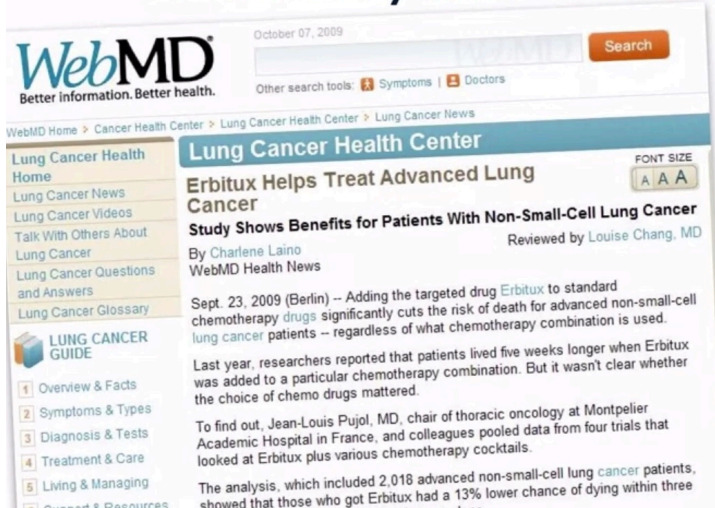


# Information Extraction saved

## Information Extraction

- Goal: Identify and extract fields of interest from free text



Eribitux helps treat lung cancer

Author: Charlene Laino

Reviewer: Louise Chang, MD

Sept. 23, 2009

Berlin ...

## Fields of Interest

- Named entities
  - **[NEWS]** People, Places, Dates, ...
  - **[FINANCE]** Money, Companies, ...
  - **[MEDICINE]** Diseases, Drugs, Procedures, ...

# Named Entity Recognition

- **Named entities:** Noun phrases that are of specific type and refer to specific individuals, places, organizations, ...
- **Named Entity Recognition:** Technique(s) to identify all mentions of pre-defined named entities in text
  - Identify the mention / phrase: *Boundary detection*
  - Identify the type: *Tagging / classification*

## Approaches to identify named entities

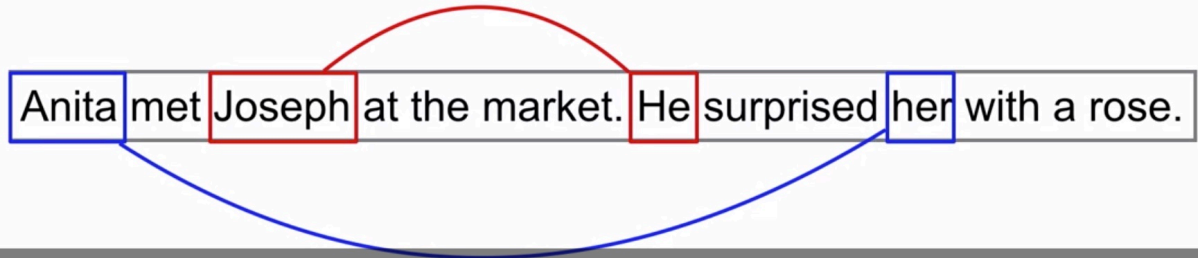
- Depends on kinds of entities that need to be identified
- For well-formatted fields like date, phone numbers:  
Regular expressions (Recall Week 1)
- For other fields: Typically a machine learning approach

## Person, Organization, Location/GPE

- Standard NER task in NLP research community
- Typically a four-class model
  - PER
  - ORG
  - LOC / GPE
  - Other / Outside (any other class)

## Co-reference resolution

- Disambiguate mentions and group mentions together



## Question Answering

- Given a question, find the most appropriate answer from the text
  - What does Erbitux treat?
  - Who gave Anita the rose?
- Builds on named entity recognition, relation extraction, and co-reference resolution

## Take Home Concepts

- Information Extraction is important for natural language understanding and making sense of textual data
- Named Entity Recognition is a key building block to address many advanced NLP tasks
- Named Entity Recognition systems extensively deploy supervised machine learning and text mining techniques discussed in this course

2. If the shortest distance between words A and B in the WordNet hierarchy is 6, the path-based similarity measure  $\text{PathSim}(A,B)$  would be:

☐ 6

☐  $1/6 = 0.167$

☐  $1 - 1/5 = 5/6 = 0.833$

☒  $1/(6+1) = 1/7 = 0.143$