# Semantic NLP
# Clinical NLP Systems

BMI701 Introduction of Biomedical Informatics
Lab Session 6
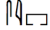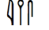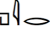
---

Wei-Hung Weng

October 26, 2016

HMS DBMI — MGH LCS

HARVARD
MEDICAL SCHOOL

MASSACHUSETTS
GENERAL HOSPITAL

# Deciphering Hieroglyphs

| | | 🏺🐦🔺 | 𝔑𝔩☐ | 𝔮𝔦𝔫 | ☐𝔩⌣ | 𝔮𝔮🔺 | ⌣𝔩🐟 |
|---|---|---|---|---|---|---|---|
| (knife) | ⌇𝔄𝔩⌇ | 51 | 20 | 84 | 0 | 3 | 0 |
| (cat) | ⌣🐦🔺 | 52 | 58 | 4 | 4 | 6 | 26 |
| **???** | ⌣𝔰𝔩🏺 | **115** | **83** | **10** | **42** | **33** | **17** |
| (boat) | 𝔩🐦𝔦🔺 | 59 | 39 | 23 | 4 | 0 | 0 |
| (cup) | ⌣𝔄𝔩☐ | 98 | 14 | 6 | 2 | 1 | 0 |
| (pig) | ☐𝔩🏺𝔩☐ | 12 | 17 | 3 | 2 | 9 | 27 |
| (banana) | ⌇𝔄⌇𝔄 | 11 | 2 | 2 | 0 | 18 | 0 |

Evert 2010

## Deciphering Hieroglyphs

|  |  | 🔲⬮⚬ | 🅼⬜ | ⬙⚑ | ⬜⬙⚬ | ⬙⬙⚬ | ⬮⬙🐾 |
|---|---|---|---|---|---|---|---|
| (knife) | 〰🦅🔪 | 51 | 20 | 84 | 0 | 3 | 0 |
| (cat) | ⬮⬮⚬ | 52 | 58 | 4 | 4 | 6 | 26 |
| **???** | ⬮𓆑🔲 | 115 | 83 | 10 | 42 | 33 | 17 |
| (boat) | 𓂝⬮⚑⚬ | 59 | 39 | 23 | 4 | 0 | 0 |
| (cup) | ⬮🦅▫ | 98 | 14 | 6 | 2 | 1 | 0 |
| (pig) | ▫🦅🦅⬜ | 12 | 17 | 3 | 2 | 9 | 27 |
| (banana) | 〰🦅〰🦅 | 11 | 2 | 2 | 0 | 18 | 0 |

$$\mathsf{sim}(\text{⬮𓆑🔲}, \text{〰🦅🔪}) = 0.770$$

# Deciphering Hieroglyphs

|  |  | 🔲🐦◠ | 𝌆⬜ | 𝌆𝌆 | 🔲◠ | 𝌆◠ | ◠𝌆 |
|---|---|---|---|---|---|---|---|
| (knife) | 🐦 | 51 | 20 | 84 | 0 | 3 | 0 |
| (cat) | ◠◠◠ | 52 | 58 | 4 | 4 | 6 | 26 |
| **???** | ◠𝌆 | **115** | **83** | **10** | **42** | **33** | **17** |
| (boat) | 𝌆◠ | 59 | 39 | 23 | 4 | 0 | 0 |
| (cup) | ◠🐦 | 98 | 14 | 6 | 2 | 1 | 0 |
| (pig) | 🔲𝌆𝌆 | 12 | 17 | 3 | 2 | 9 | 27 |
| (banana) | 🐦🐦 | 11 | 2 | 2 | 0 | 18 | 0 |

$$\mathsf{sim}(\text{◠𝌆}, \text{🔲𝌆𝌆}) = 0.939$$

Evert 2010

# Deciphering Hieroglyphs

| | | 🐦🐦 | 🐦🐦 | 🐦🐦 | 🐦🐦 | 🐦🐦 | 🐦🐦 |
|---|---|---|---|---|---|---|---|
| (knife) | 🐦 | 51 | 20 | 84 | 0 | 3 | 0 |
| (cat) | 🐦 | 52 | 58 | 4 | 4 | 6 | 26 |
| **???** | 🐦 | 115 | 83 | 10 | 42 | 33 | 17 |
| (boat) | 🐦 | 59 | 39 | 23 | 4 | 0 | 0 |
| (cup) | 🐦 | 98 | 14 | 6 | 2 | 1 | 0 |
| (pig) | 🐦 | 12 | 17 | 3 | 2 | 9 | 27 |
| (banana) | 🐦 | 11 | 2 | 2 | 0 | 18 | 0 |

$$\text{sim}(🐦, 🐦) = 0.961$$

# Deciphering Hieroglyphs

| | get | see | use | hear | eat | kill |
|---|---|---|---|---|---|---|
| knife | 51 | 20 | 84 | 0 | 3 | 0 |
| cat | 52 | 58 | 4 | 4 | 6 | 26 |
| **dog** | 115 | 83 | 10 | 42 | 33 | 17 |
| boat | 59 | 39 | 23 | 4 | 0 | 0 |
| cup | 98 | 14 | 6 | 2 | 1 | 0 |
| pig | 12 | 17 | 3 | 2 | 9 | 27 |
| banana | 11 | 2 | 2 | 0 | 18 | 0 |

verb-object counts from British National Corpus

## Tf-idf Weighting

- Importance of the term in the corpus
- For term $i$ in document $j$

$$w_{i,j} = tf_{i,j} \times \log(\frac{N}{df_i})$$

- $tf_{i,j}$: frequency of $i$ in $j$
- $df_i$: number of documents have $i$
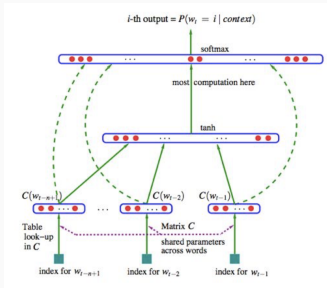- $N$: number of all documents

## Tf-idf Weighting Example

- A: "Dog is so cute."
- B: "I like dog."
- $tfidf_{('dog',A)} = \frac{1}{4} \times \log(\frac{2}{2}) = 0$
- $tfidf_{('dog',B)} = \frac{1}{3} \times \log(\frac{2}{2}) = 0$
- $tfidf_{('cute',A)} = \frac{1}{4} \times \log(\frac{2}{1}) = \frac{\log 2}{4}$
- $tfidf_{('cute',B)} = \frac{0}{3} \times \log(\frac{2}{0}) = 0$
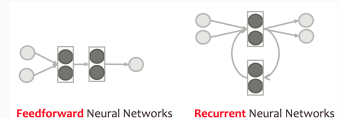
## Semantic Approaches

- Matrix decomposition
    - LSI (Deerwester 1990), NMF (Lee 1999), NTF (Cruys 2010)
    - Using SVD: $U \Sigma V$
    - Fast, unless using NTF
- Language model
    - PLSI (Hofmann 1999), LDA (Blei 2003)
    - Topic modeling, using probability
    - Heavy computation

# Semantic Approaches

- Neural network model
  - NNLM (Bengio 2003), RNN (Mikolov 2010), skip-gram / CBOW (Mikolov 2013)
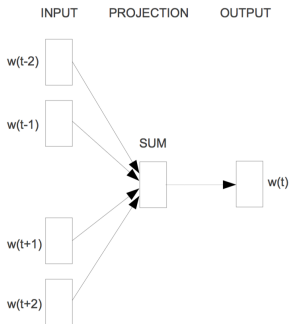  - Heavy computation, hard to implementation
  - Interpretation...?



Bengio 2003



Chang 2015

# Semantic Approaches



Mikolov 2013

Country and Capital Vectors Projected by PCA

Mikolov 2013

| Linguistic preprocessing | → | Morphological analysis | → | Syntatic analysis | → | Semantic analysis |

| Tokenizer Sentence splitter Stopword tagger | Stemmer | POS tagger Shallow parser | UMLS Dictionary lookup |

MetaMap / cTAKES workflow

- Developed by NLM (Aronson 1994)
- Web application of MetaMap
- Java API
- Locally execution
- Download

## cTAKES

- Developed by Mayo NLP (Savova 2010)
- Modularized
- CLI
- Download

## Demonstration

- Topic modeling (`topicmodels`)
- MetaMap / cTAKES
    - Need to download in advance
    - Use CLI or `system()` in R
    - Further processing

## Some Advanced NLP Courses

- NLTK book (very useful!)
- Coursera NLP (Jurafusky)
- Coursera NLP (Radev)
- Coursera NLP provided by Michael Collins is also good, but it's gone now
- Coursera NLP (Collins)
- CS287: Natural Language Processing
- 6.864: Advanced Natural Language Processing

- More text mining techniques
  - Topic modeling, vector space model
- MetaMap / cTAKES
- Contact
  - Github repository
  - ckbjimmy@gmail.com
  - Linkedin: Wei-Hung Weng