

LING 520: Computational Analysis of English

Semester: FALL '16

Instructor: Sowmya Vajjala

Iowa State University, USA

15 November 2016

Class Outline

- ▶ NLP Applications: Information Extraction and Machine Translation
- ▶ Information Extraction
- ▶ Machine Translation
- ▶ Reminder: Assignment 5 submission due - tonight.

Regarding Assignment 5 - Question 1

- ▶ If you started working on A5 in the past 1-2 days, you may have noticed something is wrong with Stanford Parser online demo.
- ▶ Here are three options to answer Question 1 in this case:
 1. Use any other parser demo online and do the question 1
 2. use Stanford parser GUI that comes with the download (refer to README to know how)
 3. Write python code to print Stanford parser outputs and analyse that. (which you perhaps did in question 3)

NLP applications - overview

- ▶ *Text classification* - we discussed a few weeks back.
- ▶ *Information extraction*
- ▶ *Machine Translation*
- ▶ Information retrieval (Search)
- ▶ Dialog systems/conversational agents
- ▶ Question Answering/Summarization

Information Extraction

Information Extraction - overview

- ▶ Task: Extract different types of information (names, dates, relationships etc) from text.
- ▶ Let us take an example text:

Supermoon is an event that happens when a full moon is closest to Earth. It orbits our planet in an oval shape so sometimes it comes closer to us than at other times. To us Earth-lings, the moon appears 30 per cent brighter and 14 per cent bigger. By the way, supermoon is not an astrological term. It's scientific name is perigee-syzygy, but supermoon is more catchy. Astrologer Richard Nolle first came up with the term and he defined it as "... a new or full moon which occurs with the moon at or near (within 90% of) its closest approach to Earth in a given orbit", according to earthsky.org.

When is the next supermoon?

Monday, November 14. This supermoon will be the biggest and brightest in 70 years, so it will definitely be worth a look. The "undeniably beautiful" astronomical event will not come again until November 25, 2034, NASA said.

Information Extraction - tasks

- ▶ identify and classify "named entities" in the text - Named Entity Recognition (Supermoon, perigee-syzygy, Richard Nolle, NASA)
- ▶ Not sufficient to just say something is a proper noun. What sort of proper noun is it? person, event? organization? place?

Information Extraction - tasks

- ▶ identify and classify "named entities" in the text - Named Entity Recognition (Supermoon, perigee-syzygy, Richard Nolle, NASA)
- ▶ Not sufficient to just say something is a proper noun. What sort of proper noun is it? person, event? organization? place?
- ▶ Event detection. Key events in our example: supermoon, its next occurrence, etc.

Information Extraction - tasks

- ▶ identify and classify "named entities" in the text - Named Entity Recognition (Supermoon, perigee-syzygy, Richard Nolle, NASA)
- ▶ Not sufficient to just say something is a proper noun. What sort of proper noun is it? person, event? organization? place?
- ▶ Event detection. Key events in our example: supermoon, its next occurrence, etc.
- ▶ Temporal Expressions: Monday, November 14, November 25, 2034 etc.

Information Extraction - tasks

- ▶ identify and classify "named entities" in the text - Named Entity Recognition (Supermoon, perigee-syzygy, Richard Nolle, NASA)
- ▶ Not sufficient to just say something is a proper noun. What sort of proper noun is it? person, event? organization? place?
- ▶ Event detection. Key events in our example: supermoon, its next occurrence, etc.
- ▶ Temporal Expressions: Monday, November 14, November 25, 2034 etc.
- ▶ All this information is useful for: Relation extraction (identifying that Richard Nolle is the scientist who coined the term Supermoon)

Information Extraction - Methods

- ▶ Regular expressions (If you know the patterns of named entities, temporal expressions etc)
- ▶ Machine learning (If we do not know the patterns)
- ▶ Example: NER can be seen as a sequence labeling problem like in POS tagging, coupled with gazetteers containing names of persons, organizations etc.

Information Extraction and NLTK

- ▶ NER example (NER.py)
- ▶ Relation extraction example (RelExtraction.py)

Some current IE projects: code and data

- ▶ openIE - <http://knowitall.github.io/openie/>
- ▶ reVerb project - <http://reverb.cs.washington.edu/>
- ▶ Google relation extraction corpus
<https://research.googleblog.com/2013/04/50000-lessons-on-how-to-read-relation.html>

Machine Translation

Machine Translation

- ▶ Task: translate text from one language to another (text can be a word, a sentence, a paragraph, a full document)
- ▶ Primary issues in solving the task:
 1. Different scripts, different word orders, different orthographic rules (no caps, no punctuation etc), different morphological structure

Machine Translation

- ▶ Task: translate text from one language to another (text can be a word, a sentence, a paragraph, a full document)
- ▶ Primary issues in solving the task:
 1. Different scripts, different word orders, different orthographic rules (no caps, no punctuation etc), different morphological structure
 2. Understanding genre specific nuances (translating news versus translating a novel)
 3. advanced issues: stylistic and cultural differences, idioms, metaphors etc.

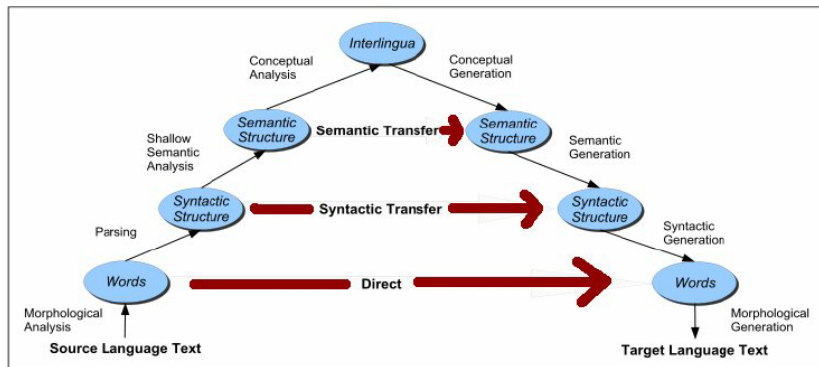
Machine Translation

- ▶ Task: translate text from one language to another (text can be a word, a sentence, a paragraph, a full document)
- ▶ Primary issues in solving the task:
 1. Different scripts, different word orders, different orthographic rules (no caps, no punctuation etc), different morphological structure
 2. Understanding genre specific nuances (translating news versus translating a novel)
 3. advanced issues: stylistic and cultural differences, idioms, metaphors etc.
 4. degrees of automatic translation: rough translation, translation + post editing, high-quality translation for very targeted domain specific language.

Machine Translation - Usage scenarios

- ▶ Rough translation: translating webpages for a general purpose reader
- ▶ Translation with post-editing: translation of software manuals for localization (Computer aided human translation).
Computer does some part of translation followed by human translator edits.
- ▶ domain specific: weather forecasts (where vocabulary is limited, and language patterns are also limited).

Machine Translation - Classical view



source: <https://goo.gl/slXqkS>

Rule based MT

- ▶ Direct translation: do direct word by word translation from source to target. Not used now, but the same intuition of incremental translation works in currently used systems too.
- ▶ Transfer rules based translation: have rules for translating from X to Y, to account for word-order differences, and other such issues in addition to lexical transfer rules.
- ▶ Direct + transfer rules based translation

Machine Translation - Statistical view

Bayes is Back!

- ▶ Idea: If you have a large collection of parallel sentences from source and target languages, you can approximate human translation with a statistical model.
- ▶ If F is a foreign language and we are translating from F to English, probability of best translation:
$$= \operatorname{argmax}_{E \in \text{English}} P(F|E) * P(E)$$
- ▶ Where $P(E)$ is the language model for target language (English), which helps us choose the translation that is most likely in English language
- ▶ $P(F|E)$ is the translation model. A commonly used translation model is "Phrase Based Statistical Machine Translation".
- ▶ Challenge: Creating alignments between the source and target sentences in the training data.

Machine translation - Statistical view

An example phrase table. source: <https://goo.gl/cORL3E>

```

wife | manzela | 0.0107991 0.0203291 0.000693866 0.0012239 | 0-0 | 526 14412 10 |
wife | manzela | 0.00008 0.0002178 0.00062448 0.0002584 | 0-0 | 11250 14412 9 |
wife | ji | 0.0001000 0.0001322 0.000693866 0.000674471 | 0-0 | 5707 14412 10 |
wife | . manzela | 0.482755 0.584211 0.00134283 0.00513742 | 0-1 | 58 14412 28 |
wife | manzela | 0.130435 0.584211 0.00137348 0.0347526 | 0-0 | 507 14412 27 |
wife | cseuul | 0.0379147 0.4128652 0.00562448 0.0006411 | 0-0 | 237 14412 9 |
wife | manzelenkendo | 0.0205479 0.016385 0.00200816 0.001863 | 0-0 | 1460 14412 30 |
wife | starta | 0.001008 0.0023666 0.0014896 0.0012239 | 0-0 | 4160 14412 18 |
wife | chat | 0.366492 0.370192 0.00485706 0.0044877 | 0-0 | 191 14412 70 |
wife | zene , ze | 0.017448 0.339074 0.0041632 0.00028566 | 0-0 | 189 14412 6 |
wife | manzelaive | 0.00021023 0.0001357 0.000485706 0.000408 | 0-0 | 33015 14412 7 |
wife | manzel | 0.00087634 0.001182 0.000495706 0.0005828 | 0-0 | 7986 14412 7 |
wife | zennakka | 0.007471 0.020284 0.0004141 0.0006411 | 0-0 | 81 14412 6 |
wife | zennaty | 0.00713558 0.0151643 0.000445706 0.0010491 | 0-0 | 981 14412 7 |
wife | zene ma | 0.186047 0.218609 0.00555593 0.00051047 | 0-0 | 43 14412 8 |
wife | wife | 0.40875 0.00041632 0.0006411 | 0-0 | 14412 6 |
wife | zennakko | 0.28 0.354839 0.000485706 0.0006411 | 0-0 | 25 14412 7 |
wife | ne zenuo | 0.0101988 0.017739 0.000485706 0.00081939 | 0-1 | 493 14412 18 |
wife | manzela | 0.00063012 0.0006776 0.000346933 0.0002914 | 0-0 | 7935 14412 5 |
wife | zene - | 0.0297919 0.218609 0.000346933 0.00279138 | 0-0 | 168 14412 5 |
wife | wikileaks | 0.000271 | 0.508273 | 0.529163 | 2 2 |
wikileaks founder julian | zakladatel wikileaks julian | 0.143037 0-5 0.180444 | 1-0 0-1 2-2 | 2 4 2 |
wikileaks founder zakladatel wikileaks | 0.000000 | 0.000000 | 0.000000 | 0-0 1-0 | 1 |
wikileaks has | wikileaks | 0.0555556 0.026732 0-5 0.277782 | 0-0 1-0 | 18 2 1 |
wikileaks | surveru wikileaks | 0.5 0.299373 0.047619 0.330862 | 0-0 0-1 | 2 21 1 |
wikileaks | wikileaks | 0.666667 0.298613 0.0852381 0.0385802 | 0-0 0-1 | 3 21 2 |
wikileaks | wikileaks | 0.333333 0.597015 0.047619 0.175554 | 0-0 | 3 21 1 |
wikileaks | wikileaks ne | 0.5 0.597015 0.047619 0.0234101 | 0-0 | 2 21 1 |
wikileaks | wikileaks | 0.722222 0.597015 0.619048 0.555356 | 0-0 | 18 21 3 |
wikileaks | prirode a | 0.0484945 0.0861874 0.0135135 0.00991068 | 0-0 1-1 | 88 296 4 |
wikileaks | prirody a | 0.015907 0.0270334 0.0135135 0.00786363 | 0-0 1-1 | 36 296 6 |
wikileaks | nasa volne sijici a | 0.666667 0.128285 0.0135135 5.15859e-07 | 0-1 0-2 1-3 | 6 296 4 |
wikileaks | volne sijici | 0.618562 0.28895 | 0.770707 0.010011 | 0-1 27 0-1 | 27 |
wikileaks | volne sijici a | 0.25 0.0259009 0.027027 0.000170571 | 0-0 0-1 1-2 | 32 296 8 |
wikileaks | divokyt a | 0.5 0.249007 0.0168919 0.0188663 | 0-0 1-1 | 10 296 5 |
wikileaks | divokych a | 0.3123 0.216521 0.033738 0.0278444 | 0-0 1-1 | 32 296 21 |
wikileaks | divoky a | 0.370371 0.309782 0.0675676 0.0346088 | 0-0 1-1 | 56 296 20 |
wikileaks | divoká | 0.344263 0.300483 0.0704859 0.0504974 | 0-0 1-1 | 61 296 12 |
wikileaks | divoká a | 0.342857 0.313755 0.0405405 0.0292602 | 0-0 1-1 | 35 296 12 |
wikileaks | divoká , | 0.0714286 0.0049893 0.0135135 0.0010169 | 0-0 1-1 | 56 296 4 |
wikileaks | divoci | 0.361311 0.333967 0.049189 0.0119093 | 0-0 1-1 | 36 296 13 |
wikileaks | a | 9.243e-07 2.2291e-05 0.0168919 0.854051 | 0-0 1 | 5.4095e+06 296 5 |
wikileaks | a | 0.00230814 0.00004607 0.0135135 0.0003829 | 0-0 1-0 | 1733 296 6 |
wikileaks | divoce a | 0.0641026 0.0520312 0.0168919 0.0143154 | 0-0 1-1 | 78 296 5 |
wikileaks | prirode a | 0.117647 0.0861874 0.00675676 0.0031373 | 0-0 1-2 | 17 296 2 |
wikileaks | sece divoci a | 0.5 0.333967 0.00675676 0.11991e-05 | 0-0 1-2 | 5 296 2 |
wikileaks | robusnefne a | 0.785714 0.00852105 0.00675676 0.000157316 | 0-0 1-1 | 7 296 2 |
wikileaks | zblida | 0.1333 0.0151 | 0.00675676 0.000235976 | 0-0 1-1 | 296 |
wikileaks | divokos a | 0.153846 0.275795 0.00675676 0.0142368 | 0-0 1-1 | 33 296 2 |
wikileaks | zidovely a | 0.75 0.189787 0.0101351 0.0010121 | 0-0 1-1 | 4 296 3 |
wikileaks | meci planety zaslani a | 0.328234 0.00675676 0.18527e-13 | 0-1 3-3 | 2 296 2 |
wikileaks | volne sijici a | 1 0.220253 0.00675676 0.0013059 | 0-0 0-1 1-2 | 2 296 2 |
wikileaks | tak | 1.60929e-11 1.6247e-06 0.00137038 0.0034894 | 0-0 1-1 | 277063 296 3 |
wikileaks | divokoho a | 0.25 0.305389 0.0101351 0.0174637 | 0-0 1-1 | 12 296 3 |
wikileaks | volne sijicni lososa a | 0.5 0.020268 0.0101351 2.57161e-09 | 0-0 0-1 1-3 | 6 296 3 |
wikileaks | planety zaslani a | 0.038234 0.00675676 0.2618e-06 | 0-0 1-2 | 296 2 |
wikileaks | volne sijicich a | 0.666667 0.274843 0.00675676 0.0168932 | 0-0 0-1 1-2 | 3 296 2 |
wikileaks | masefne a | 0.66667 0.000677 0.0101351 0.00014632 | 0-0 1-1 | 9 296 3 |
wikileaks | divoky kraj , i | 0.5 0.0625458 0.00675676 0.079791e-10 | 0-0 1-3 | 4 296 2 |
wikileaks | neposustane a | 0.666667 0.076808 0.00675676 0.00078658 | 0-0 1-1 | 3 296 2 |
wikileaks | volne sijici lososa a | 0.128285 0.0101351 0.01303e-06 | 0-0 0-1 1-3 | 296 3 |
wikileaks | volne sijici exemplare a | 0.128285 0.00675676 2.404e-09 | 0-0 0-1 1-3 | 2 296 2 |
wikileaks | sijici a | 0.42857 0.0214 0.0137038 0.0922638 | 0-0 1 | 7 296 |
wikileaks | volne sijicim | 0.864647 0.154382 0.00455474 0.00055449 | 0-0 0-1 | 30 6367 20 |
wikileaks | divoka | 0.407051 0.403013 0.0181882 0.0591269 | 0-0 | 1280 6367 551 |
wikileaks | volne sijici | 0.876288 0.285407 0.013501 0.00129006 | 0-0 | 3367 6367 85 |
wikileaks | volne sijicijo | 0.664384 0.271795 0.0152348 0.000215077 | 0-0 0-1 | 146 6367 31 |
wikileaks | divokaj | 0.336389 0.00562448 0.0003548 | 0-0 | 197 6367 58 |
wikileaks | divokoho | 0.368421 0.409594 0.0263841 0.0204453 | 0-0 | 456 6367 165 |
wikileaks | divokam | 0.262486 0.330357 0.00785299 0.0066163 | 0-0 | 177 6367 50 |
wikileaks | volne sijicich | 0.439787 0.368683 0.0906225 0.0197099 | 0-0 0-1 | 1312 6367 557 |
wikileaks | prirode | 0.180825 0.115596 0.0164913 0.0116043 | 0-0 | 972 6367 105 |
wikileaks | divokaj | 0.277778 0.319773 0.0164913 0.0160251 | 0-0 | 198 6367 105 |
wikileaks | divoka | 0.413666 0.420815 0.0417779 0.0342605 | 0-0 | 643 6367 266 |
wikileaks | sijicich | 0.322581 0.499149 0.0424062 0.108031 | 0-0 | 837 6367 270 |
wikileaks | divokos | 0.380743 0.370143 0.0213601 0.0166597 | 0-0 | 77 6367 136 |
wikileaks | divokej | 0.704545 0.729167 0.00486886 0.0032234 | 0-0 | 44 6367 31 |
wikileaks | sijici | 0.040774 0.060515 0.016172 0.0642844 | 0-0 | 2855 6367 103 |
wikileaks | wid | 0.800971 0.804867 0.0259149 0.0173144 | 0-0 | 206 6367 165 |
wikileaks | prirody | 0.071732 0.0362845 0.0142924 0.0082098 | 0-0 | 2448 6367 91 |
wikileaks | divoci | 0.457536 0.474748 0.018619 0.0129658 | 0-0 | 258 6367 115 |
wikileaks | volne sijici | 0.704172 0.172058 0.0874823 0.0112784 | 0-0 0-1 | 791 6367 557 |
wikileaks | divokych | 0.2812 0.290402 0.040487 0.0326057 | 0-0 | 3049 6367 295 |
wikileaks | divoce | 0.0525101 0.0697853 0.0142924 0.0167618 | 0-0 | 1733 6367 91 |
wikileaks | divoky | 0.392814 0.415486 0.051154 0.0405231 | 0-0 | 835 6367 328 |
wikileaks | volne | 0.050316 0.238101 0.039263 0.182446 | 0-0 | 63446 6367 230 |
wikileaks | volne prirode | 0.0594461 0.0602298 0.0014354 2.24439e-05 | 0-0 0-1 | 154 6367 9 |
```

Evaluating MT systems

- ▶ Human raters (in terms of: correctness, clarity, naturalness, grammaticality etc)
- ▶ automatic evaluation - BLEU score (ngram similarity based measure between a translated sentence and a gold standard human translation)
- ▶ Other automatic measures: TER, METEOR, NIST etc.

Useful Resources

- ▶ Chapter 25 in Jurafsky and Martin (very comprehensive overview)
- ▶ <http://www.statmt.org/> - comprehensive website on readings, software, corpora related to developing and testing statistical machine translation systems.
- ▶ <https://www.apertium.org> - open source machine translation platform that lets you create rule-based MT models.
- ▶ Google Translate, Bing Translate etc - MT applications in real life.

Exercise in Analyzing machine translation

- ▶ Two popular machine translation software online: Google Translate and Bing translate
- ▶ Task: Try out how both these tools work for translating from English to your native language and viceversa (Native English speakers: choose another language you know. If you know only English, pair up with some other person).
- ▶ Spend some time with these and we can discuss your observations towards the end of the class.

Next class

- ▶ NLP for CALL - Introduction
- ▶ Assignment 5 discussion
- ▶ Final projects - status and discussion
- ▶ Recommended Readings: Burstein (2009), Chapelle and Chung (2010). Meurers (2013). Read atleast two of them. Large amount of thursday's class invovles discussion about these papers.
- ▶ Optional readings: Chapters on Writing Aids and Tutoring Systems in Language and Computers.