

CSCE 590-1: From Data to Decisions with Open Data: A Practical Introduction to AI

Lecture 19/20: Text Data and Analysis

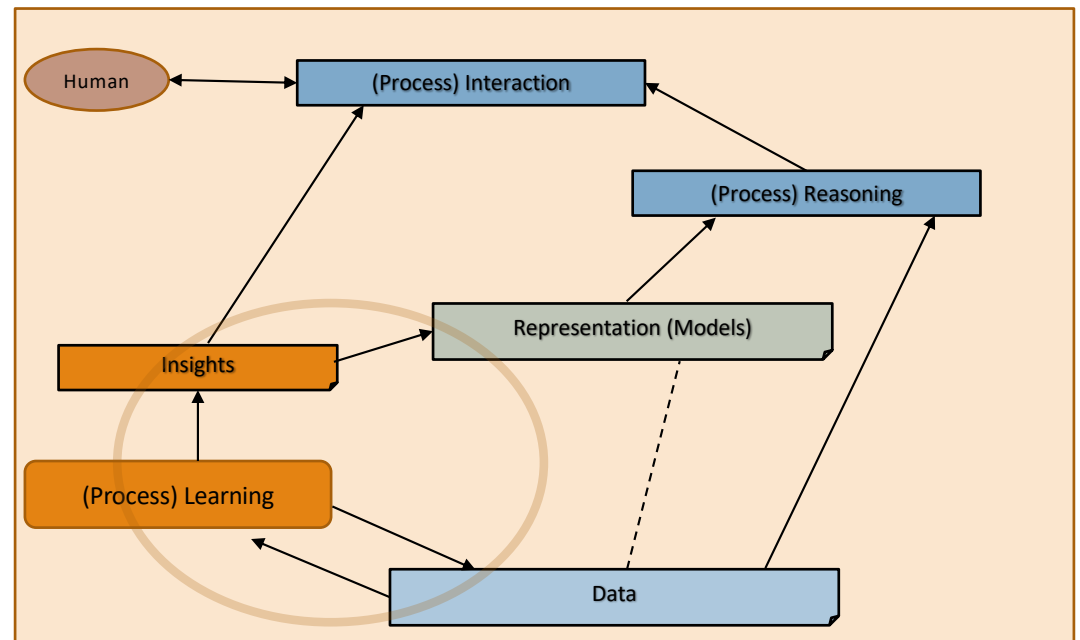
PROF. BIPLAV SRIVASTAVA, AI INSTITUTE

18TH, 23RD MAR, 2021

Carolinian Creed: “I will practice personal and academic integrity.”

Organization of Lecture 20

- Introduction Segment
 - Recap/ Discussion of Lecture 18
 - Reading material – Lecture 19
- Main Segment
 - What is text ?
 - Words and forms
 - Multi-lingual
 - Nature of analysis possible
 - How it complements numerical analysis
- Concluding Segment
 - About Next Lecture – Lecture 21
 - Quiz 3
 - Ask me anything



Introduction Segment

Recap of Lecture 18

- We looked at a learning agent
- Reinforcement learning method
 - Various variations
- Bayesian Optimization

Example Situation – Course Selection

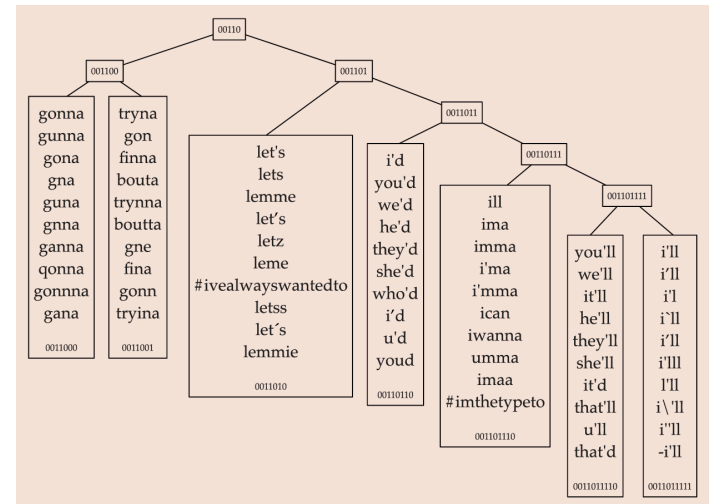
- A person wants to pass an academic program in two majors: A and B
- There are three subjects: A, B and C, each with three levels (*1, *2, *3). There are thus 9 courses: A1, A2, A3, B1, B2, B3, C1, C2, C3
- To graduate, at least one course at beginner (*1) level is needed in major(s) of choice(s), and two courses at intermediate levels (*2) are needed
- **Role of Text:** *Course descriptions, pre-requisites, test grading, ...*
- **Answer questions**
 - Q1: What are the topics to be covered in a course ?
 - Q2: What topics has a student already studied in previous program that can be reused?
 - ...

Lecture 19 – Reading

- Due to adverse weather, we had to change plan

Contextual Word Representations: Putting Words into Computers, CACM 2020,

<https://cacm.acm.org/magazines/2020/6/245162-contextual-word-representations/fulltext>



Review - Contextual Word Embeddings

- Words as discrete
- Words with distributional assumptions:
 - Context: given a word, its nearby words or sequences of words
 - Words used in similar ways are likely to have related meanings; i.e., words used in the same (similar) context have related meanings
 - No claim about meaning except relative similarity v/s dis-similarity of words
- Two main strategies
 - Compare with words in a manually-created taxonomy, e.g., Wordnet
 - Learn context and representation from data

Credit:

Contextual Word Representations: Putting Words into Computers", by Noah Smith, CACM
June 2020

Main Segment

Understanding Concepts - Words

What is a Word ?

- Unix command - `man wc`

“A word is defined as a string of characters delimited by white space characters.”

- Example
 - Content = “CSCE 771: Computer Processing of Natural Language
Lecture 3: Words, Morphology, Lexicons
Prof. Biplav Srivastava, AI Institute
31st Aug 2020 ”
 - Command = “`wc -w content.txt`”
Result = “`20 content.txt`” (stored in file - result.txt)
 - “CSCE 771: Computer Processing of Natural Language (7)
Lecture 3: Words, Morphology, Lexicons (12)
Prof. Biplav Srivastava, AI Institute (17)
31st Aug 2020 ” (20)

Types of Words in English

Content words (open-class – i.e., continuously changing):

- **Nouns:** student, university, knowledge,...
- **Verbs:** write, learn, teach,...
- **Adjectives:** difficult, boring, hard,
- **Adverbs:** easily, repeatedly,...

Function words (closed-class – fixed):

- **Articles:** a, an, the
- **Prepositions:** in, with, under,...
- **Conjunctions:** and, or,...
- **Determiners:** a, the, every,...

Another Language - Turkish

A Turkish word

Chinese: 我开始写小说 = 我 开始 写 小说
I start(ed) writing novel(s)

uygarlaştıramadıklarımızdanmışsınızcasına
uygar_laş_tır_ama_dık_lar_ımız_dan_mış_sınız_casına

"as if you are among those whom we were not able to civilize (=cause to become civilized)"

uygar: civilized

_laş: become

_tır: cause somebody to do something

_ama: not able

_dık: past participle

_lar: plural

_ımız: 1st person plural possessive (our)

_dan: among (ablative case)

_mış: past

_sınız: 2nd person plural (you)

_casına: as if (forms an adverb from a verb) K. Oflazer pc to J&M

A strict reliance on spaces will make us miss useful parts of text

Common Definitions

- **Corpus** (plural corpora): a computer-readable corpora collection of text or speech.
- **Lemma**: A lemma is a set of lexical forms having the same stem, the same major part-of-speech, and the same word sense. [Example: Cat and cats have same lemma.](#)
- **Word form**: The word form form is the full inflected or derived form of the word. [Example: Cat and cats have different word forms.](#)
- **Word type**: Types are the number of distinct words in a corpus. if the set of words is V , the number of types is the word token vocabulary size $|V|$.
- **Word tokens**: The total number N of running words in the sentence / document of interest.
- **Code switching**: use multiple languages in a code switching single communicative act – [Example: Hindlish \(Hindi English\), Spanish \(Spanish English\)](#)

“They picnicked by [the](#) pool, then lay back on [the](#) grass and looked at [the](#) stars.”

- 16 tokens, 14 word types

Source: Jurafsky & Martin

Lexical Meaning – Common Terms

- **Synonym:** same/ similar meaning
 - start-begin, finish-end, far-distant
- **Antonym:** opposite meaning
 - Far – near, clever - stupid, high - low, big – small
- **Homonym:** identical in spelling and pronunciation
 - bear, bank, ...
- **Homophones:** sounds identical but are written differently
 - site-sight, piece-peace.
- **Homograph:** written identically but sound differently
 - Potato, tomato, lead, wind, minute
- **Polysemy:** a word or phrase which has two(or more) rated meanings
 - Duck, sharp

Source: Mausam

Knowing About Words

_____ Of **course** he wants to **take** the advanced **course** **too**. _____
He already **took** **two** beginners' **courses**.

- Words – set of characters separated by spaces
- Word forms –
 - Spelling differences - specialize v/ specialise
 - Meaning similarity/differences - Take/ took, course/ courses, two/ too
- Word types – distinct words

Pop Quiz: Are word tokens and word types same in the example above?

Word Variety

- **Inflection:** creates different forms of the same word
 - Verbs: to be, being, I am, you are, he is, I was,
 - Nouns: one book, two books
- **Derivation:** creates different words from the same lemma
 - grace ⇒ disgrace ⇒ disgraceful ⇒ disgracefully
- **Compounding:** new words from combinations

“ice cream”, “website”, “web site”, “New York-based”

- **Clitics** - *a clitic is a morpheme that has syntactic characteristics of a word, but depends phonologically on another word or phrase. In this sense, it is syntactically independent but phonologically dependent ...*

English: “doesn’t” , “I’m” ,

Italian: “dirglielo” = dir + gli(e) + lo // tell + him + it

New words over time:

Google ⇒ Googler, to google, to ungoogle,
to misgoogle, googlification, ungooglification,
googlified, Google Maps,
Google Maps service, ...

Review: Regular Expression

Metacharacter	Explanation
<code>^</code>	Matches the starting position within the string
<code>.</code>	Matches any single character
<code>[]</code>	Matches a single character that is contained within the brackets
<code>[^]</code>	Matches a single character that is not contained within the brackets.
<code>\$</code>	Matches the ending position of the string
<code>*</code>	Matches the preceding element zero or more times
<code>+</code>	Matches the preceding element one or more times
<code> </code>	Separates choices

Regex	Matches any string that
<code>hello</code>	contains {hello}
<code>gray grey</code>	contains {gray, grey}
<code>gr(a e)y</code>	contains {gray, grey}
<code>gr[ae]y</code>	contains {gray, grey}
<code>b[aeiou]bble</code>	contains {babble, bebble, bibble, bobble, bubble}
<code>[b-chm-pP]at ot</code>	contains {bat, cat, hat, mat, nat, oat, pat, Pat, ot}
<code>colou?r</code>	contains {color, colour}
<code>rege(x(es)? xps?)</code>	contains {regex, regexes, regexp, regexps}
<code>go*gle</code>	contains {ggle, gogle, google, gooogle, goooogle, ...}
<code>go+gle</code>	contains {gogle, google, gooogle, goooogle, ...}
<code>g(oog)+le</code>	contains {google, googoogle, googoogoogle, googoogoogle, ...}
<code>z{3}</code>	contains {zzz}
<code>z{3,6}</code>	contains {zzz, zzzz, zzzzz, zzzzzz}
<code>z{3,}</code>	contains {zzz, zzzz, zzzzz, ...}

Example Source: <https://cs.lmu.edu/~ray/notes/regex/>

Implementation: Finding Words in Python

- Python has extended Regex specifications for convenience
- Useful for
 - Matching patterns
 - Information extraction
 - Content manipulation (e.g., substitution)
 - Error (e.g., spelling) correction

```
data = "The CSCE 771 course is taught at  
University this Fall!"  
pattern = "[tT]+\w"  
m = re.findall(pattern, data)  
print(m)
```

```
['Th', 'ta', 'ty', 'th']
```

Details: <https://docs.python.org/3/library/re.html>

Code Examples

- Regular expressions
 - <https://github.com/biplav-s/course-d2d-ai/blob/main/sample-code/l20-text-overview/WordLesson-Examples.ipynb>
- Supporting multiple languages – encoding
 - <https://github.com/biplav-s/course-d2d-ai/blob/main/sample-code/l20-text-overview/Multiple%20Languages.ipynb>

Morphology

Morphemes: The small meaningful units that make up words

Stems: The core meaning-bearing units

Affixes: Bits and pieces that adhere to stems

Morphemes: stems, affixes

dis-grace-ful-ly
prefix-stem-suffix-suffix

Many word forms consist of a **stem** plus a number of **affixes** (*prefixes* or *suffixes*)

Infixes are inserted inside the stem.

Circumfixes (German gesehen) surround the stem

Morphemes: the smallest (meaningful/grammatical) parts of words.

Stems (grace) are often **free morphemes**.

Free morphemes can occur by themselves as words.

Affixes (dis-, -ful, -ly) are usually **bound morphemes**.

Bound morphemes have to combine with others to form words.

- Plural nouns add -s to singular:
 - book-books,
- but:
 - box-boxes, fly-flies, child-children
- Past tense verbs add -ed to infinitive:
 - walk-walked,
- but:
 - like-liked, leap-leapt

Source: Julia Hirschberg

Morphological Generation

- Generate legal variations.
 - For **grace** (**stem**): grace**ful**, grace**fully**, **dis**grace, **dis**grace**ful**, **dis**grace**fully**, ungraceful, ungracefully, undisgraceful, **un****dis**grace**fully**,...
- But avoid ungrammatical variations
 - *grace**ly****ful**, *gracefully, *disungracefully,...

Source: Julia Hirschberg

Advanced Topic –Language Formalism

An **alphabet** Σ is a **set of symbols**:

e.g. $\Sigma = \{a, b, c\}$

A **string** ω is a **sequence of symbols**, e.g. $\omega = abcb$.

The **empty string** ϵ consists of zero symbols.

The Kleene closure Σ^* ('**sigma star**') is the (**infinite**) **set of all strings** that can be formed from Σ :

$\Sigma^* = \{\epsilon, a, b, c, aa, ab, ba, aaa, \dots\}$

A **language** $L \subseteq \Sigma^*$ over Σ is also a set of strings.

Typically we only care about **proper subsets** of Σ^* ($L \subset \Sigma^*$).

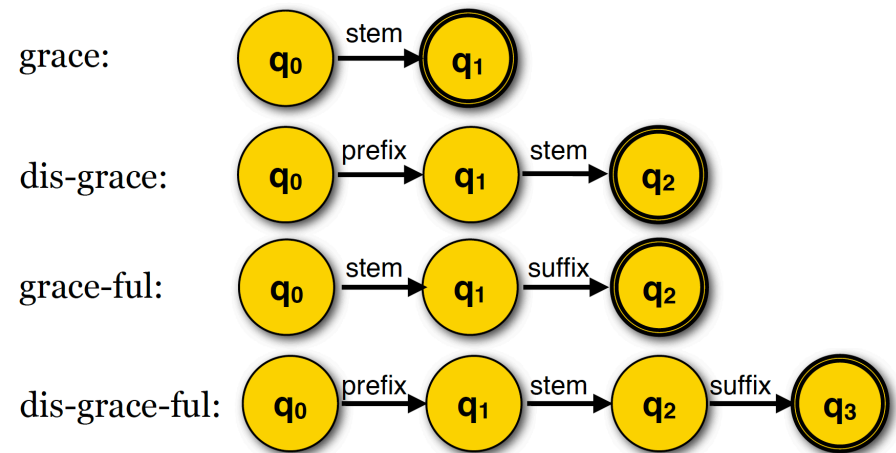
- Automata
- Finite State Automata
- Deterministic Finite State Automata (DFSA)
- Non-Deterministic Finite State Automata (NDFS A)

Source: Julia Hirschberg

Advanced Topics – Recognizing as Automata

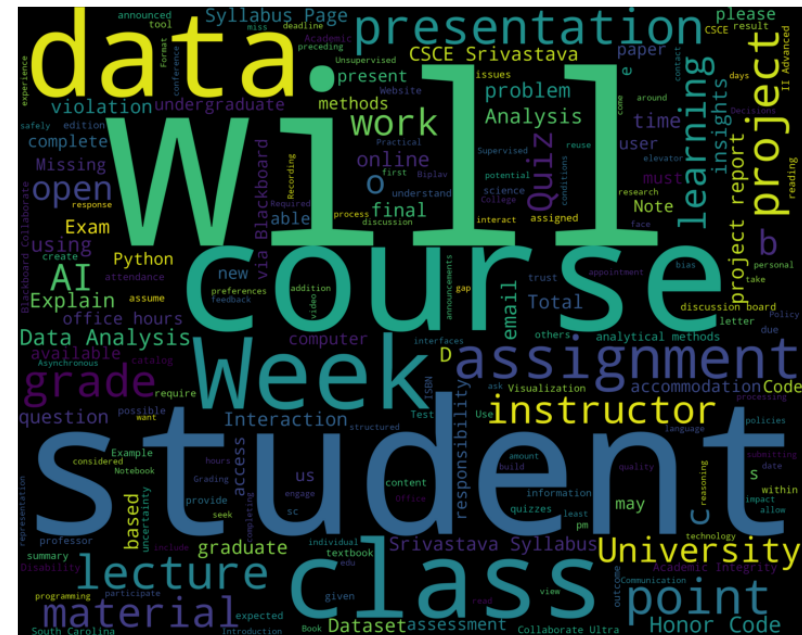
- Automata –
 - an abstract model of a computer which reads an input string, and changes its internal state depending on the current input symbol. It can either accept or reject the input string.
 - Hence, an automata defines a language
- Finite State Automata – regular expressions

Source: Julia Hirschberg



Analysis: Insights About a Course

Course Description:
CSCE 590-1



How NLP Complements Quantitative Analysis

- Quantitative data:
 - Captures precise information about well-defined attributes
 - Allows all the tools of mathematics for analysis
 - Example: average student scores
- Textual data
 - Captures multi-dimensional information
 - Needs careful consideration about the context of information
 - Language can be imprecise
 - Example: topics covered in a course
- Both have strengths

Common NLP Tasks

- Sentiment
- Summarization
- Machine translation
- Natural Language Interface to Databases
- Natural Language Generation

Code Examples

- Word Tag Cloud, translation
 - https://github.com/biplav-s/course-d2d-ai/blob/main/sample-code/l20-text-overview/FirstLook_ClassSyllabusData.ipynb

Review - Contextual Word Embeddings Paper

- Words as discrete
- Words with distributional assumptions:
 - Context: given a word, its nearby words or sequences of words
 - Words used in similar ways are likely to have related meanings; i.e., words used in the same (similar) context have related meanings
 - No claim about meaning except relative similarity v/s dis-similarity of words
- Two main strategies
 - Compare with words in a manually-created taxonomy, e.g., Wordnet
 - Learn context and representation from data

Credit:

Contextual Word Representations: Putting Words into Computers", by Noah Smith, CACM
June 2020

Lecture 20: Concluding Comments

- Word representation paper – keep for overall context of text processing
- We looked at text word, forms, types
- Multiple languages
- Applications: Word tag cloud

Concluding Segment

Upcoming Classes

15	Mar 4 (Th)	Reasoning and Search	Semester - Midpoint
16	Mar 9 (Tu)	Agent – Optimization	
17	Mar 11 (Th)	Agent – Handling Uncertain World	
18	Mar 16 (Tu)	Agent – Learning	
19	Mar 18 (Th)	Reading – Word representation	Change due to weather
20	Mar 23 (Tu)	Text: Data	Quiz 3
21	Mar 25 (Th)	Review: project presentations, Discussion	
22	Mar 30 (Tu)	Text: Summary, Sentiment	
23	Apr 1 (Th)	Text: Visualization, Explanation	
24	Apr 6 (Tu)	Paper presentations – Graduate students	Final assignment for Graduate students
25	Apr 8 (Th)	Case Study 1: Water (Structured+Text)	Quiz 4
26	Apr 13 (Tu)	Case Study 2: Finance (Structured+Text)	



Grad students – select papers

About Next Lecture – Lecture 21

Lecture 20: Review of Projects

- Each student will prepare a 1 slide presentation covering
 - Problem, User of results, data, approach, status
 - Template on next slide
 - Put slide at:
<https://drive.google.com/drive/folders/1tofaeCDEJToxoMbMpFEcWVtjveiLWh6o?usp=sharing>
- Present for 5 mins + 2 mins question

Project Name:
Student Name:

Problem

User of Results

Data

Approach

Status

Quiz 3
