

LECTURE 1

□ Sentence Segmentation:

Detecting a sentence correctly.

• CHALLENGES

- "." cannot always detect (USA)
- "!", → "Alas! he died"
- "?" → "What!?"

• SOLUTION

- Rules, &
 - Large set of rules are hard to maintain.
- machine Learning
 - 99% perfect.

□ Tokenization: Detect a word?

But on what?

• Challenges

- punctuation, C++, C#
- emoticons. 😊 ☹
- New York vs York!
- Phrasal verb → ~~work~~ work out
- Different Language.

Solution:

- Each language needs it's own solⁿ. So no general solⁿ.

□ Stemming

অবাক লেভ এ আগুন নাগায় দেওয়া

Equivalence → Equival.

Fast, but inaccurate.

Organization → Organ

UN লেভ হলে শুধু!

□ Lemmatization

- Hand built lexicon for all word forms. (Walked → walk).

- Accurate but slower.

- Chicken egg scenario with POS tagging.

// Lemmatization is like figuring out what bird (the root word) an egg (the word in the sentence) will become. But to do that accurately, you need to know the context in which it is used, which is provided by POS tagging. //



context $\rightarrow 0$ $\rightarrow 0$ $\rightarrow 0$

□ Embedding

— Word as a vector.

$\rightarrow \text{Dog} = [0.5; 0.4; 0.1]$

$\text{Dogs} = [0.5; 0.4; 0.2]$

□ Parts of speech.

◦ Closed class \rightarrow Suppose pronouns.

◦ Open class

\rightarrow Names! \rightarrow "Eshita the gach guru"

Challenge

* One word can have diff. pos tag based on it's use.



Named Entity Recognition

Identify phrases that are named people, locations, organizations.

Person: Turing is the . . .

Org: ACM is a . . .

location: The Mt. Everest is

Geo pol. entity: U.S.A is the . . .

Challenge

* Ambiguity

- Washington was born into slavery
- Washington went up 2 games
- I arrived Washington today
- Washington passed a pri . . . law.

What/who
is actually
WASHINGTON?!



Solⁿ

BIO → Begin name
Inside "
Outside "

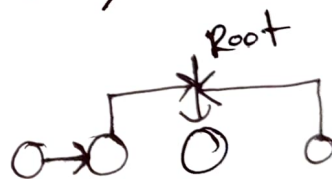
BILOU → Begin "
Inside "
Last "
Outside "
Unit-length.

Parsing & Syntactic Representation.

✦ Constituency tree



* Dependency tree



Challenge:

Attachment Ambiguity.

One morning, I saw Messi in my Pajamas.

Who was in my PAJAMAS?

— me?

— messi?

Coordination Ambiguity.

Old men and women.

— old (men and women) ?

— old (men) and women ?

Solⁿ

* Probabilistic grammar parsing

* Transition based parsing.