

# Natural Language Processing - Week02

Assist. Prof. Dr. Tuğba YILDIZ

İSTANBUL BİLGİ UNIVERSITY  
Department of Computer Engineering

March 1, 2017

## 1 Natural Language Processing (NLP)

## 2 Levels of NLP

## 3 Why NLP is Hard?

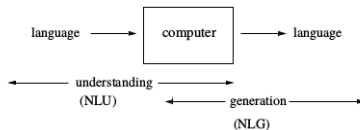
## 4 Models and Algorithms

## 5 NLP - An Interdisciplinary Field

## 6 Regular Expressions and Automata

# Natural Language Processing (NLP)

- NLP is a field for computers to analyze, understand, and/or generating human language.
- computers using natural language as input and/or output
- concerns understanding and generating spoken and written text
- we will mostly concern with **written text** (not speech).



# Natural Language Processing (NLP)

- Natural Language Understanding : Mapping the given input in the natural language into a useful representation.
- Natural Language Generation : Producing output in the natural language from some internal representation.

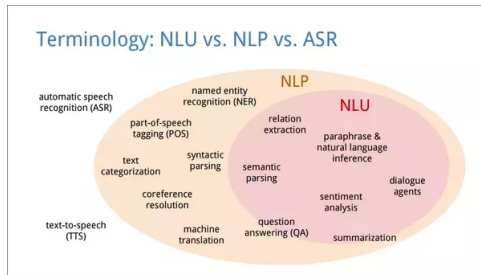


Fig.1 NLU in NLP <sup>1</sup>

<sup>1</sup><http://nlp.stanford.edu/wcmac/papers/20140716-UNLU.pdf>

# Levels of Natural Language Processing (NLP)

- To process written text, we need:
  - Phonetics and Phonology - the study of linguistic sounds and their relations to words
    - Phonology : concerns units of speech
    - Phonetics : deals with how speech is organized
  - Morphology - knowledge of the meaningful components of words
  - Syntax - knowledge of the structural relationships between words
  - Semantics - knowledge of meaning
  - Pragmatics - knowledge of the relationship of meaning to the goals and intentions of the speaker (The meaning of the sentence depends on an understanding of the context and the speaker's intent.)
  - Discourse - knowledge about linguistic units larger than a single utterance

# Levels of Natural Language Processing (NLP)

- **Phonetics and Phonology** - The study of linguistic sounds and their relations to words
- Turkish orthography is highly regular and a word's pronunciation is always completely identified by its spelling.

Turkish	IPA	English approximation	Turkish	IPA	English approximation
A a	/a/	As a in <i>father</i>	M m	/m/	As m in <i>man</i>
B b	/b/	As b in <i>boy</i>	N n	/n/	As n in <i>nice</i>
C c	/dʒ/	As j in <i>Joy</i>	O o	/o/	As o in <i>more</i>
Ç ç	/tʃ/	As ch in <i>chair</i>	Ö ö	/ø/	As e in <i>learn</i> , with lips rounded
D d	/d/	As d in <i>dog</i>	P p	/p/	As p in <i>pin</i>
E e	/e/ <sup>[1]</sup>	As e in <i>red</i>	R r	/r/ <sup>[2]</sup>	As r in <i>ring</i>
F f	/f/	As f in <i>far</i>	S s	/s/	As s in <i>song</i>
G g	/g/, /ɟ/	As g in <i>got</i>	Ş ş	/ʃ/	As sh in <i>show</i>
Ğ ğ	/ɰ/, /ɰ̯/, /ɰ̯̯/	<sup>[3]</sup>	T t	/t/	As t in <i>tick</i>
H h	/h/	As h in <i>hot</i>	U u	/u/	As oo in <i>too</i>
I ı	/ɯ/	As e in <i>open</i>	Ü ü	/y/	As e in <i>new</i>
İ i	/i/	As ee in <i>feet</i>	V v	/v/	As v in <i>vat</i>
J j	/ʒ/	As s in <i>measure</i>	Y y	/j/	As y in <i>yes</i>
K k	/k/, /c/	As k in <i>kit</i>	Z z	/z/	As z in <i>zigzag</i>
L l	/l/, /ɭ/	As l in <i>love</i>			

Fig.2 It presents the Turkish letters, the sounds they correspond to in International Phonetic Alphabet and how these can be approximated more or less by an English speaker. <sup>2</sup>

<sup>2</sup>[https://en.wikipedia.org/wiki/Turkish\\_alphabet](https://en.wikipedia.org/wiki/Turkish_alphabet)

# Levels of Natural Language Processing (NLP)

- **Morphology** - The study of internal structures of words and how they can be modified
- Parsing complex words into their components

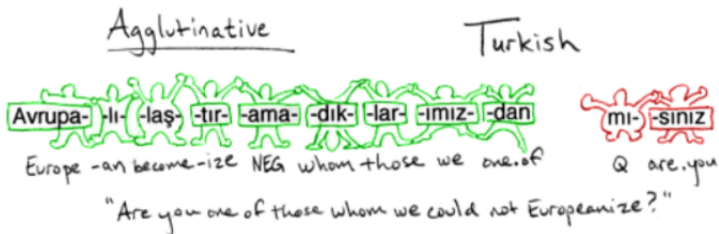


Fig.2 Morphological Parsing <sup>3</sup>

<sup>3</sup>[https://hpi.de/fileadmin/user\\_upload/fachgebiete/plattner/teaching/NaturalLanguageProcessing/Week02/NLP\\_Levels\\_of\\_NLP.pdf](https://hpi.de/fileadmin/user_upload/fachgebiete/plattner/teaching/NaturalLanguageProcessing/Week02/NLP_Levels_of_NLP.pdf)

# Levels of Natural Language Processing (NLP)

- **Morphology** - The study of internal structures of words and how they can be modified
- Relatively simple for English. But for some languages such as Turkish, it is more difficult.
  - uygarlaştıramadıklarımızdanmışsınızcasına
  - uygar-laş-tır-ama-dık-lar-ımız-dan-mış-sınız-casına
  - uygar +BEC +CAUS +NEGABLE +PPART +PL +P1PL +ABL +PAST +2PL +Aslf
  - “(behaving) as if you are among those whom we could not civilize/cause to become civilized”
    - +BEC is “become” in English
    - +CAUS is the causative voice marker on a verb
    - +PPART marks a past participle form
    - +P1PL is 1st person plural possessive marker
    - +2PL is 2nd person plural
    - +ABL is the ablative (from/among) case marker
    - +Aslf is a derivational marker that forms an adverb from a finite verb form
    - +NEGABLE is “not able” in English
- Inflectional and Derivational Morphology.
- Common tools: Finite-state transducers



# Levels of Natural Language Processing (NLP)

- **Morphology** - The study of internal structures of words and how they can be modified
- Part of Speech Tagging
  - Each word has a part-of-speech tag to describe its category.
  - Part-of-speech tag of a word is one of major word groups
  - open classes – noun, verb, adjective, adverb
  - closed classes – prepositions, determiners, conjunctions, pronouns, participles
  - POS Taggers try to find POS tags for the words.
  - duck is a verb or noun? (morphological analyzer cannot make decision).
  - A POS tagger may make that decision by looking the surrounding words.
  - Duck! (verb)
  - Duck is delicious for dinner. (noun)
- Inflectional and Derivational Morphology.
- Common tools: Finite-state transducers

# Levels of Natural Language Processing (NLP)

- **Syntax** - The study of the structural relationships between words in a sentence
- Parsing - converting a flat input sentence into a hierarchical structure that corresponds to the units of meaning in the sentence.
- There are different parsing formalisms and algorithms.
- Most formalisms have two main components:
  - grammar – a declarative representation describing the syntactic structure of sentences in the language.
  - parser – an algorithm that analyzes the input and outputs its structural representation (its parse) consistent with the grammar specification.
- CFGs are in the center of many of the parsing mechanisms. But they are complemented by some additional features that make the formalism more suitable to handle natural languages.

# Levels of Natural Language Processing (NLP)

- **Syntax** - The study of the structural relationships between words in a sentence
- Parsing complex words into their components

```
(ROOT
(S
(NP
(NP (NNP Surgical) (NN resection) (NNS specimens))
(P (IN of)
(NP
(NP (CD 85) (JJ invasive) (JJ ductal) (NNS carcinomas))
(P (IN of)
(NP
(NP (CD 85) (NNS women))
(SBAR
(whNP (WP who))
(S
(VP (VBD had)
(VP (VBN undergone)
(NP (CD 3D) (NN ultrasound))))))))))
(VP (VBD were)
(VP (VBN included)))
(. .)))
```

Fig.2 Morphological Parsing <sup>4</sup>

<sup>4</sup>[https://hpi.de/fileadmin/user\\_upload/fachgebiete/plattner/teaching/NaturalLanguageProcessing/Week02/NLP\\_Levels\\_of\\_NLP.pdf](https://hpi.de/fileadmin/user_upload/fachgebiete/plattner/teaching/NaturalLanguageProcessing/Week02/NLP_Levels_of_NLP.pdf)

# Levels of Natural Language Processing (NLP)

- **Semantic** - The study of the meaning of words, and how these combine to form the meanings of sentences
- Assigning meanings to the structures created by syntactic analysis.
- Mapping words and structures to particular domain objects in way consistent with our knowledge of the world.
- Semantic can play an import role in selecting among competing syntactic analyses and discarding illogical analyses.
- I robbed the bank – bank is a river bank or a financial institution
- We have to decide the formalisms which will be used in the meaning representation.

# Levels of Natural Language Processing (NLP)

- **Semantic** - The study of the meaning of words, and how these combine to form the meanings of sentences
- Synonymy: fall - autumn
- Hypernymy - Hyponymy (is a): animal - dog
- Meronymy - Holonym (part-whole): finger - hand
- Antonymy: big - small

# Levels of Natural Language Processing (NLP)

- **Pragmatics** - Social use of language
- The study of how language is used to accomplish goals, and the influence of context on meaning
- Understanding the aspects of a language which depends on situation and world knowledge
- Example:
  - Give me the salt!
  - Could you please give me the salt?

# Levels of Natural Language Processing (NLP)

- **Discourse** - The study of linguistic units larger than a single statement
- John reads a book. He borrowed it from his friend.

**Berlin** (/bəˈrɪn/, German: [bɛʁˈliːn] (listen)) is the capital of Germany, and one of the 16 states of Germany. With a population of 3.5 million people,<sup>[4]</sup> Berlin is Germany's largest city. It is the second most populous city proper and the seventh most populous urban area in the European Union.<sup>[5]</sup> Located in northeastern Germany on the banks of River Spree, it is the center of the Berlin-Brandenburg Metropolitan Region, which has about 6 million residents from over 180 nations.<sup>[6][7][8][9]</sup> Due to its location in the European Plain, Berlin is influenced by a temperate seasonal climate. Around one third of the city's area is composed of forests, parks, gardens, rivers and lakes.<sup>[10]</sup>

Fig.2 NLP Levels with an Example <sup>5</sup>

<sup>5</sup><http://nlp.stanford.edu/wcmac/papers/20140716-UNLU.pdf>

# Natural Language Processing (NLP)

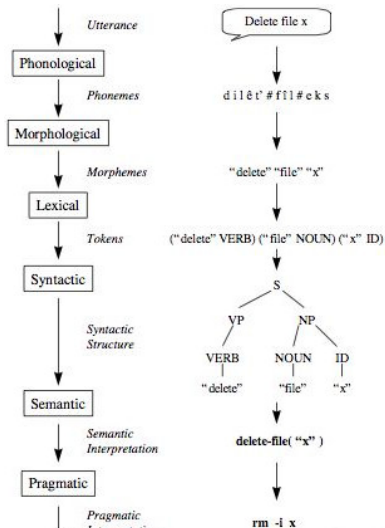


Fig.2 NLP Levels with an Example <sup>6</sup>



# Natural Language Processing (NLP)

- **Natural Language Understanding** : Mapping the given input in the natural language into a useful representation.
- Different level of analysis required:
  - Words
  - Morphological Analysis : Morphologically analyzed words
  - Syntactic Analysis : Syntactic structure
  - Semantic analysis : Context-independent meaning representation
  - Discourse analysis : Final meaning representation

# Natural Language Processing (NLP)

- **Natural Language Generation** : Producing output in the natural language from some internal representation.
- Different level of synthesis required:
  - Utterance
  - Utterance Planning : meaning representations for sentences
  - Sentence Planning and Lexical Choice : Syntactic structures of sentences with lexical choices
  - Sentence Generation : Morphologically analyzed words
  - Morphological Generation
  - Words

# NLP Applications

- Applications:
  - automatic summarization
  - machine translation
  - named entity recognition
  - relationship extraction
  - sentiment analysis
  - speech recognition
  - topic segmentation
  - dialogue systems

# NLP Applications: Dialogue Systems

- 2001: A Space Odyssey (directed by Stanley Kubrick)
  - -Dave: Open the pod bay doors, HAL.
  - -HAL: I'm sorry Dave, I'm afraid I can't do that.
  - -Dave: What's the problem?
  - -HAL: I think you know what the problem is just as well as I do.

# Natural Language Processing (NLP)

- HAL must be able to
  - **Phonetics and Phonology** - recognize words from an audio signal and to generate audio signal from a sequence of words.
  - Morphology - produce and recognize words and other variations of individual words (I can't, doors)
    - ["I'm", "sorry", "Dave", "I'm", "afraid", "I", "can't", "."]
  - Syntax - use structural knowledge to properly string together the words that constitute its response.
    - "I'm I do, sorry that afraid Dave I'm can't."
  - Semantics - know meaning of the words
  - Syntax - use structural knowledge to properly string together the words that constitute its response.
    - "Yes, I'd like to hear it, HAL. Sing it for me!"

# Natural Language Processing (NLP)

- HAL must be able to
  - Phonetics and Phonology - recognize words from an audio signal and to generate audio signal from a sequence of words.
  - **Morphology** - produce and recognize words and other variations of individual words (I can't, doors)
    - ["I'm", "sorry", "Dave", "I'm", "afraid", "I", "can't", "."]
  - Syntax - use structural knowledge to properly string together the words that constitute its response.
    - "I'm I do, sorry that afraid Dave I'm can't."
  - Semantics - know meaning of the words
    - "Yes, I'd like to hear it, HAL. Sing it for me!"

# Natural Language Processing (NLP)

- HAL must be able to
  - Phonetics and Phonology - recognize words from an audio signal and to generate audio signal from a sequence of words.
  - Morphology - produce and recognize words and other variations of individual words (I can't, doors)
    - ["I'm", "sorry", "Dave", "I'm", "afraid", "I", "can't", "."]
  - **Syntax** - use structural knowledge to properly string together the words that constitute its response.
    - "I'm I do, sorry that afraid Dave I'm can't."
  - Semantics - know meaning of the words
    - "Yes, I'd like to hear it, HAL. Sing it for me!"

# Natural Language Processing (NLP)

- HAL must be able to
  - Phonetics and Phonology - recognize words from an audio signal and to generate audio signal from a sequence of words.
  - Morphology - produce and recognize words and other variations of individual words (I can't, doors)
    - ["I'm", "sorry", "Dave", "I'm", "afraid", "I", "can't", "."]
  - Syntax - use structural knowledge to properly string together the words that constitute its response.
    - "I'm I do, sorry that afraid Dave I'm can't."
  - **Semantics** - know meaning of the words
    - "Yes, I'd like to **hear** it, HAL. **Sing** it for me!"



# Natural Language Processing (NLP)

- HAL must be able to
  - Phonetics and Phonology - recognize words from an audio signal and to generate audio signal from a sequence of words.
  - Morphology - produce and recognize words and other variations of individual words (I can't, doors)
    - ["I'm", "sorry", "Dave", "I'm", "afraid", "I", "can't", "."]
  - Syntax - use structural knowledge to properly string together the words that constitute its response.
    - "I'm I do, sorry that afraid Dave I'm can't."
  - **Semantics** - know meaning of the words
    - "Yes, I'd like to hear it, HAL. Sing it for me!"

# Natural Language Processing (NLP)

## ■ HAL must be able to

- **Pragmatics** - have knowledge about the kind of actions that speakers intend by their use of sentence.

-Dave: Open the pod bay doors, HAL.

-HAL: No or No, I won't open the door. (-)

-HAL: I'm sorry Dave, I'm afraid I can't do that. (+)

- **Discourse** - need to examine the earlier questions that were asked.

-Dave: Yes, I'd like to hear it, HAL. Sing it for me.

-HAL: It's called "Daisy." Daisy, Daisy, give me your answer, do. I'm half crazy, all for the love of you. It won't be a stylish marriage.

# Natural Language Processing (NLP)

- HAL must be able to
  - Pragmatics - have knowledge about the kind of actions that speakers intend by their use of sentence.
    - Dave: Open the pod bay doors, HAL.
    - HAL: No or No, I won't open the door. (-)
    - HAL: I'm sorry Dave, I'm afraid I can't do that. (+)
  - **Discourse** - need to examine the earlier questions that were asked.
    - Dave: Yes, I'd like to hear it, HAL. Sing it for me.
    - HAL: It's called "Daisy." Daisy, Daisy, give me your answer, do. I'm half crazy, all for the love of you. It won't be a stylish marriage.

# Why NLP is Hard?

- Ambiguity is at all levels of analysis:
  - Phonetics and phonology : "I scream" vs. "ice cream" / "recognize speech" vs. "wreck a nice beach"
  - Morphology : unionized = union + ized? / un + ionized?
  - Syntax:
    - John saw the man on the mountain with a telescope.
    - Who is on the mountain? John, the man, or both? Who has the telescope? John, the man, or the mountain?
  - Semantics: The astronomer loves the star. (star in the sky / celebrity)
  - Pragmatic:
    - "I just came from New York."
    - Would you like to go to New York today?
    - Would you like to go to Boston today?
    - Boy, you look tired.
  - Discourse: John and Mary took two trips around France. They were both wonderful.

# Ambiguity

- I made her duck.
  - How many different interpretations does this sentence have?
  - What are the reasons for the ambiguity?
  - The categories of knowledge of language can be thought of as ambiguity resolving components.
  - How can each ambiguous piece be resolved?

# Ambiguity (cont)

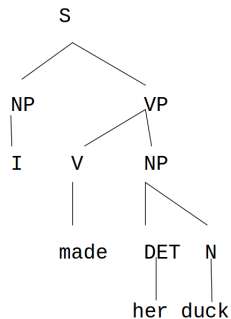
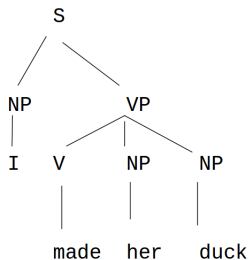
- Some interpretations of : I made her duck.
  - 1 I cooked duck for her.
  - 2 I cooked duck belonging to her.
  - 3 I created a toy duck which she owns.
  - 4 I caused her to quickly lower her head or body.
  - 5 I used magic and turned her into a duck.
- duck – morphologically and syntactically ambiguous: noun or verb.
- her – syntactically ambiguous: dative or possessive.
- make – semantically ambiguous: cook or create.
- make – syntactically ambiguous

## Ambiguity (cont)

- I made her duck.
- part-of-speech tagging - Deciding whether duck is verb or noun.
- word-sense disambiguation - Deciding whether make is create or cook.
- lexical disambiguation - Resolution of part-of-speech and word-sense ambiguities are two important kinds of lexical disambiguation.
- syntactic ambiguity - her duck is an example of syntactic ambiguity, and can be addressed by probabilistic parsing.

# Ambiguity (cont)

## ■ I made her duck.





# Example: Ambiguity in Turkish

## ■ masalı

- 1 masal +Noun+A3sg+Pnon+Acc (= the story)
- 2 masal +Noun+A3sg+P3sg+Nom (= his story)
- 3 masa +Noun+A3sg+Pnon+NomaDB+Adj+With (= with tables)

# Models and Algorithms

- state machines, rule systems, logic, probabilistic models, and vector-space models.
  - state machines (deterministic and non-deterministic finite-state automata and finite-state transducers)
  - regular grammars and regular relations, context-free grammars
  - first order logic, also known as the predicate calculus
  - hidden Markov models or HMMs
  - vector-space models, based on linear algebra
  - machine learning tools like classifiers (decision trees, support vector machines, logistic regression, etc.)

# NLP - An Interdisciplinary Field

- NLP borrows techniques and insights from several disciplines.
  - Linguistics: How do words form phrases and sentences? What constraints the possible meaning for a sentence?
  - Computational Linguistics: How is the structure of sentences are identified? How can knowledge and reasoning be modeled?
  - Computer Science: Algorithms for automata, parsers.
  - Engineering: Stochastic techniques for ambiguity resolution.
  - Psychology: What linguistic constructions are easy or difficult for people to learn to use?
  - Philosophy: What is the meaning, and how do words and sentences acquire it?

# Regular Expressions and Automata

- Regular Expressions (RE): the standard notation for characterizing text sequences, a language for specifying text search strings
- RE is used for specifying text strings in situations like
  - Web-search
  - information retrieval
  - word processing
  - computing frequency from text

# Regular Expressions and Automata

- RE is a formula in a special language that is used for specifying simple classes of **strings**.
- A **string** is a sequence of symbols, any sequence of alphanumeric character
- RE search requires a **pattern** that we want to search in text.

# Regular Expressions and Automata

- How can we search for any of these?
  - woodchucks

RE	Example Patterns Matched
/woodchucks/	“interesting links to <u>woodchucks</u> and lemurs”
/a/	“ <u>M</u> ary Ann stopped by Mona’s”
/Claire_says,/	“Dagmar, my gift please,” <u>Claire</u> says,”
/DOROTHY/	“SURRENDER <u>DOROTHY</u> ”
/!/	“You’ve left the burglar behind again!” said No

# Regular Expressions

- woodchuck
- Woodchuck
- Letters inside square brackets

RE	Match	Example Patterns
/[wW]oodchuck/	Woodchuck or woodchuck	“ <u>Woodchuck</u> ”
/[abc]/	‘a’, ‘b’, or ‘c’	“In uomini, in soldat <u>i</u> ”
/[1234567890]/	any digit	“plenty of <u>7</u> to 5”

# Regular Expressions

## ■ Ranges

RE	Match	Example Patterns Matched
/ [A - Z] /	an uppercase letter	“we should call it ‘ <u>D</u> renched Blossoms”
/ [a - z] /	a lowercase letter	“ <u>m</u> y beans were impatient to be hoed!”
/ [0 - 9] /	a single digit	“Chapter <u>1</u> : Down the Rabbit Hole”



# Regular Expressions

- Negations
- Carat means negation only when first in [].

RE	Match (single characters)	Example Patterns Matched
[^A-Z]	not an uppercase letter	"Oyfn pripetchik"
[^Ss]	neither 'S' nor 's'	"I have no exquisite reason for't"
[^\.]	not a period	" <u>our</u> resident Djinn"
[e^]	either 'e' or '^'	"look up <u>^</u> now"
a^b	the pattern 'a^b'	"look up <u>a^b</u> now"

# Regular Expressions

## ■ Optional previous character

RE	Match	Example Patterns Matched
woodchucks?	woodchuck or woodchucks	<u>“woodchuck”</u>
colou?r	color or colour	<u>“colour”</u>

# Regular Expressions

- Woodchucks is another name for groundhog!
- The pipe | for disjunction

Pattern	Matches
<code>groundhog   woodchuck</code>	
<code>yours   mine</code>	yours mine
<code>a   b   c</code>	= <code>[abc]</code>
<code>[gG]roundhog   [Ww]oodchuck</code>	

# Regular Expressions

## ■ Stephen C Kleene (Kleene \*, Kleene +)

colou?r	optional previous char	color colour
ba!*	0 or more of previous char	ba ba! ba!! ba!!! ba!!!!
ba!+	1 or more of previous char	ba! ba!! ba!!! ba!!!!
beg.n	1 any char	begin begun beg1n beg.n

# Regular Expressions

- $\{n\}$  : The preceding item is matched exactly  $n$  times.
- $\{n,\}$  : The preceding item is matched  $n$  or more times
- $\{,m\}$  : The preceding item is matched at most  $m$  times.
- $\{n,m\}$  : The preceding item is matched at least  $n$  times, but not more than  $m$  times.
- for the special characters use backslash character

# Regular Expressions

## ■ Anchors

Pattern	Matches
<code>^[A-Z]  </code>	<u>P</u> alo Alto
<code>^[\^A-Za-z]</code>	<u>1</u> <u>"Hello"</u>
<code>\. \$</code>	The end <u>.</u>
<code>. \$</code>	The end? <u>.</u> The end! <u>!</u>

# Regular Expressions

- Example:
  - Find the words which has 4 characters
  - Find the words which has 4 characters

# Regular Expressions

- The process we just went through was based on fixing two kinds of errors
  - Matching strings that we should not have matched (there, then, other)
    - False positives (Type I)
  - Not matching things that we should have matched (The)
    - False negatives (Type II)
- Reducing the error rate :
- Increasing accuracy or precision (minimizing false positives)
- Increasing coverage or recall (minimizing false negatives).



# Regular Expressions

- Regular expressions play a surprisingly large role
- Sophisticated sequences of regular expressions are often the first model for any text processing text
- For many hard tasks, we use machine learning classifiers
- But regular expressions are used as features in the classifiers
- Can be very useful in capturing generalizations

# Finite State Automaton

- Regular Expressions (RE): the standard notation for characterizing text sequences
- RE is used for specifying text strings in situations like
  - Web-search
  - information retrieval
  - word processing
  - computing frequency from text

# References I