## Computational Structures in Data Science

# Lecture #13:
# Regular Expressions

UC Berkeley EECS
Adj. Assistant Prof.
Dr. Gerald Friedland

Facebook... *sigh*

April 13th, 2018

---

## Speaking of Facebook...



https://www.youtube.com/watch?v=bqWuioPHhz0

http://www.teachingprivacy.org

---

## Computational Concepts Toolbox

- **Data type: values, literals, operations,**
- **Expressions, Call expression**
- **Variables**
- **Assignment Statement**
- **Sequences: tuple, list**
- **Dictionaries**
- **Data structures**
- **Tuple assignment**
- **Function Definition Statement**
- **Conditional Statement**
- **Iteration: list comp, for, while**
- **Lambda function expr.**

- **Higher Order Functions**
  - as Values, Args, Results
- **Higher order function patterns**
  - **Map, Filter, Reduce**
  - **Function factories**
- **Recursion**
  - **Linear, Tail, Tree**
- **Abstract Data Types**
- **Mutation**
- **Iterators and Generators**
- **Object Oriented Programming, Classes**
- **Exceptions**
- **Declarative Programming**
- **Regular Expressions**

---

## On Languages...

- **Human (Natural Language)**
  - **Developed by (social) evolution**
  - **Used to communicate between humans**
  - **Change "brain state" of recipient of communication**
- **Mathematics**
  - **Formal language developed out of philosophy**
  - **Syntax (structure) and semantics (meaning) well defined**
  - **Used to communicate scientific results between humans more rigorously**
- **Programming Languages**
  - **Formalized grammar allows automatic translation**
  - **Syntax and semantics unambiguous but limited**
  - **Used to communicate between humans and computer**

# Syntax, Grammar, Semantics

- **Syntax (programming language):**
  - Set of rules that defines the combinations of symbols that are considered to be a correctly structured document or fragment in a programming language.

- **Grammar:**
  - Formalism (language) that defines the syntax of a programming language.

- **Semantics:**
  - Consequence (meaning) attached to a sequence of symbols.

# Grammars (Chomsky)

- **Grammars consist of:**
  - Terminals (literals)
  - Non-terminals
  - Production rules
- **Only define Syntax!**
- **Example:**

**terminals**
{generate, hate, great, green, ideas, linguists}

**nonterminals**
{*SENTENCE, NOUNPHRASE, VERBPHRASE, NOUN, VERB, ADJ*}

**production rules**
$SENTENCE \rightarrow NOUNPHRASE\ VERBPHRASE$
$NOUNPHRASE \rightarrow ADJ\ NOUNPHRASE$
$NOUNPHRASE \rightarrow NOUN$
$VERBPHRASE \rightarrow VERB\ NOUNPHRASE$
$VERBPHRASE \rightarrow VERB$
$NOUN \rightarrow$ ideas
$NOUN \rightarrow$ linguists
$VERB \rightarrow$ generate
$VERB \rightarrow$ hate
$ADJ \rightarrow$ great
$ADJ \rightarrow$ green

# Another Grammar

$$S \rightarrow aSb$$
$$S \rightarrow \varepsilon$$

- **S non-terminal**
- **a,b terminal**
- **Epsilon: empty!**

**Question: What are valid words in this grammar?**

# Recap: Language Structures (Python)

- **Variables** and **literals**
  - with some internal representation, e.g. Integers, Floats, Booleans, Strings, …
  
  In Python: Implicit data types!

- **Operations** on variable and literals of a **type**
  - **e.g.** `+, *, -, /, %, //, **`
  - `==, <, >, <=, >=`

- **Expressions** are valid well-defined sets of operations on variables and literals that produce a value of a type.
  - `x=4*3`

## Recap: `while` statement – iteration control

- **Repeat a block of statements until a predicate expression is satisfied**

```
<initialization statements>
while <predicate expression>:
   <body statements>

<rest of the program>
```
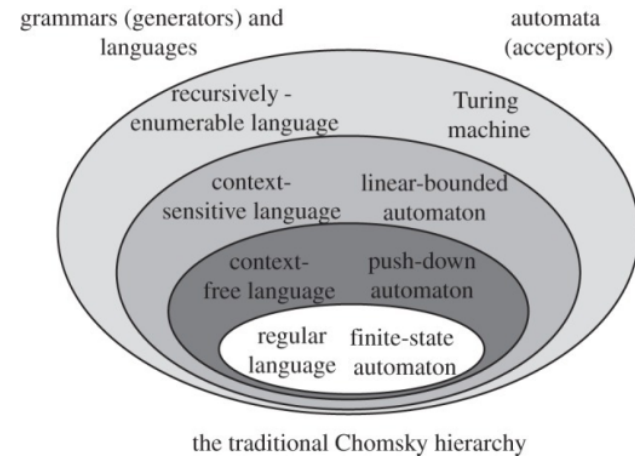
## Grammar Hierarchy



grammars (generators) and languages          automata (acceptors)

recursively - enumerable language          Turing machine

context-sensitive language          linear-bounded automaton

context-free language          push-down automaton

regular language          finite-state automaton

the traditional Chomsky hierarchy

## Grammar Hierarchy

| Language | Grammar | Automaton |
|---|---|---|
| Regular | A → a<br>A → aB | Finite state machine |
| Context-free | A → γ | Non-deterministic pushdown automaton |
| Context-sensitive | αAβ → αγβ | Linear-bounded non-deterministic Turing machine |
| Recursively enumerable | α → β<br>no restrictions | Turing machine |

- **Programming Languages usually context-free**
- **Many filter tools for data usually regular languages (i.e. form regular expressions)**

## Regular Expressions

- **Easy to parse**
- **Semantics typically:**
  - Find <string>
  - Find <string1>/Replace with <string2>
- **Widely available in:**
  - Python
  - Unix command line (bash, flex, sed)
  - Various editors (emacs, vi)
  - Most data science tools
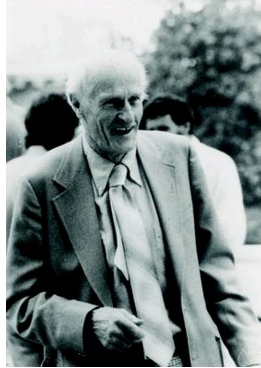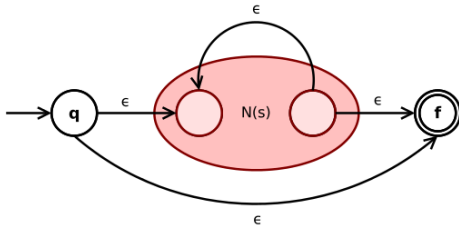
- **Important mechanism for 'cleaning' or 'extracting' data**

## Regular Expressions: Math

- **Invited by S. Kleene**
- **Most important operator: \***
  **(say 'star' or Kleene star)**
- **s\* means: 's' n-times**
  **(for varying n)**

## Regular Expressions: Python

- **re package**
- **Regular expressions are compiled into object**
- **Then various methods are available:**

```
>>> import re
>>> p = re.compile('s*')
>>> p
re.compile('s*')

>>> print(p.match(""))
None

>>> m = p.match('sssss')
>>> m
<_sre.SRE_Match object; span=(0, 5),
match='sssss'>
```

## More on REs in Python

```
>>> re.split('\W+', 'Words, words, words.')
['Words', 'words', 'words', '']
>>> re.split('(\W+)', 'Words, words, words.')
['Words', ', ', 'words', ', ', 'words', '.', '']
>>> re.split('\W+', 'Words, words, words.', 1)
['Words', 'words, words.']
>>> re.split('[a-f]+', '0a3B9', flags=re.IGNORECASE)
['0', '3', '9']

>>> re.split('(\W+)', '...words, words...')
['', '...', 'words', ', ', 'words', '...', '']

>>> re.split('x*', 'foo')
['foo']
>>> re.split("(?m)^$", "foo\n\nbar\n")
['foo\n\nbar\n']
```

## Full RE syntax

## Summary: RegEx

- **Programming Languages are very formally defined using grammars**
- **Types of grammars need different complexity**
- **Regular expressions very simple and effective tool for data filtering**
- **More in the lab!**

- **Next Lecture: Information and Bits**