

Natural Language Processing Handout

CSCE 420

November 19, 2019

1 Regular Expression

- Disjunctions `[]` :

Pattern	Matches
<code>[Ww]oodchuck</code>	woodchuck, Woodchuck
<code>[0123456789]</code>	Any single digit

- Disjunctions `|`:

Pattern	Matches
<code>abc def</code>	Find 'abc' or 'def'.
<code>a b ab</code>	Find 'a' or 'b' or 'ab'. Example: 'abc'

- Ranges:

Pattern	Matches
<code>[A - Z]</code>	An uppercase letter.
<code>[a - z]</code>	A lowercase letter.
<code>[0 - 9]</code>	A single digit.

- Negation `^`. (Note: Carat means negation only when its first in `[]`)

Pattern	Matches
<code>[^A - Z]</code>	Not upper case
<code>[^Ss]</code>	Not 'S' nor 's'
<code>[^e^]</code>	Not 'e' nor '^'
<code>a^b</code>	Search for the pattern 'a^b'

- Other notations.

<code>?</code>	0 or 1 of previous character
<code>*</code>	0 or more of previous character
<code>+</code>	1 or more of previous character
<code>.</code>	Any character
<code>^</code>	Start anchor
<code>\$</code>	End anchor
<code>\</code>	Escape character

2 Text Classification

Training set:

#	Text	Class
1	Carla Betty Carla	a
2	Carla Carla Suzanne	a
3	Carla Matt	a
4	Taylor Jessica Carla	b

Document d_5 : *Carla Carla Carla Taylor Jessica*