



Regular Expressions

Jeffrey Leek
Johns Hopkins Bloomberg School of Public Health

Regular expressions

- Regular expressions can be thought of as a combination of literals and *metacharacters*
- To draw an analogy with natural language, think of literal text forming the words of this language, and the metacharacters defining its grammar
- Regular expressions have a rich set of metacharacters

Literals

Simplest pattern consists only of literals. The literal “nuclear” would match to the following lines:

```
Ooh. I just learned that to keep myself alive after a  
nuclear blast! All I have to do is milk some rats  
then drink the milk. Aweosme. :}
```

```
Laozi says nuclear weapons are mas macho
```

```
Chaos in a country that has nuclear weapons -- not good.
```

```
my nephew is trying to teach me nuclear physics, or  
possibly just trying to show me how smart he is  
so I'll be proud of him [which I am].
```

```
lol if you ever say "nuclear" people immediately think  
DEATH by radiation LOL
```

Literals

The literal “Obama” would match to the following lines

```
Politics r dum. Not 2 long ago Clinton was sayin Obama  
was crap n now she sez vote 4 him n unite? WTF?  
Screw em both + McCain. Go Ron Paul!
```

```
Clinton conceeds to Obama but will her followers listen??
```

```
Are we sure Chelsea didn't vote for Obama?
```

```
thinking ... Michelle Obama is terrific!
```

```
jetlag..no sleep...early mornig to starbux..Ms. Obama  
was moving
```

Regular Expressions

- Simplest pattern consists only of literals; a match occurs if the sequence of literals occurs anywhere in the text being tested
- What if we only want the word “Obama”? or sentences that end in the word “Clinton”, or “clinton” or “clinto”?

Regular Expressions

We need a way to express

- whitespace word boundaries
- sets of literals
- the beginning and end of a line
- alternatives (“war” or “peace”) Metacharacters to the rescue!

Metacharacters

Some metacharacters represent the start of a line

```
^i think
```

will match the lines

```
i think we all rule for participating  
i think i have been outed  
i think this will be quite fun actually  
i think i need to go to work  
i think i first saw zombo in 1999.
```

Metacharacters

\$ represents the end of a line

```
morning$
```

will match the lines

```
well they had something this morning  
then had to catch a tram home in the morning  
dog obedience school in the morning  
and yes happy birthday i forgot to say it earlier this morning  
I walked in the rain this morning  
good morning
```


Character Classes with []

We can list a set of characters we will accept at a given point in the match

```
[Bb][Uu][Ss][Hh]
```

will match the lines

```
The democrats are playing, "Name the worst thing about Bush!"  
I smelled the desert creosote bush, brownies, BBQ chicken  
BBQ and bushwalking at Molonglo Gorge  
Bush TOLD you that North Korea is part of the Axis of Evil  
I'm listening to Bush - Hurricane (Album Version)
```

Character Classes with []

```
^[Ii] am
```

will match

```
i am so angry at my boyfriend i can't even bear to  
look at him
```

```
i am boycotting the apple store
```

```
I am twittering from iPhone
```

```
I am a very vengeful person when you ruin my sweetheart.
```

```
I am so over this. I need food. Mmmm bacon...
```

Character Classes with []

Similarly, you can specify a range of letters [a-z] or [a-zA-Z]; notice that the order doesn't matter

```
^[0-9][a-zA-Z]
```

will match the lines

```
7th inning stretch  
2nd half soon to begin. OSU did just win something  
3am - cant sleep - too hot still.. :(  
5ft 7 sent from heaven  
1st sign of starvagtion
```

Character Classes with []

When used at the beginning of a character class, the “^” is also a metacharacter and indicates matching characters NOT in the indicated class

```
[^?.]$\n
```

will match the lines

```
i like basketballs\n6 and 9\ndont worry... we all die anyway!\nNot in Baghdad\nhelicopter under water? hmmm\n
```