

## Regular Expression

**Club Moderators:** Elton Vasconcelos, Peter Mulhair, Euan McDonnell, Chew Cheng, and Dapeng Wang

# Topics to be addressed - Survey Result



# Regular Expression

~~$\log_2(FC) = 0$~~

→ Not what it is about!

Computer Science: computational language for representing patterns on textual strings  
→ interpreted by most programming languages (Java, Python, R, **PERL** ...)

## Metacharacters

char	meaning
<code>^</code>	beginning of string
<code>\$</code>	end of string
<code>.</code>	any character except newline
<code>*</code>	match 0 or more times
<code>+</code>	match 1 or more times
<code>?</code>	match 0 or 1 times; <i>or</i> : shortest match
<code> </code>	alternative
<code>()</code>	grouping; “storing”
<code>[]</code>	set of characters
<code>{}</code>	repetition modifier
<code>\</code>	quote or special

## Repetition

<code>a*</code>	zero or more <i>a</i> 's
<code>a+</code>	one or more <i>a</i> 's
<code>a?</code>	zero or one <i>a</i> 's (i.e., optional <i>a</i> )
<code>a{m}</code>	exactly <i>m</i> <i>a</i> 's
<code>a{m,}</code>	at least <i>m</i> <i>a</i> 's
<code>a{m,n}</code>	at least <i>m</i> but at most <i>n</i> <i>a</i> 's
<code>repetition?</code>	same as <i>repetition</i> but the <i>shortest</i> match is taken

Read the notation *a*'s as “occurrences of strings, each of which matches the pattern *a*”. Read *repetition* as any of the repetition expressions listed above it. Shortest match means that the shortest string matching the pattern is taken. The default is “greedy matching”, which finds the longest match. The *repetition?* construct was introduced in Perl version 5.

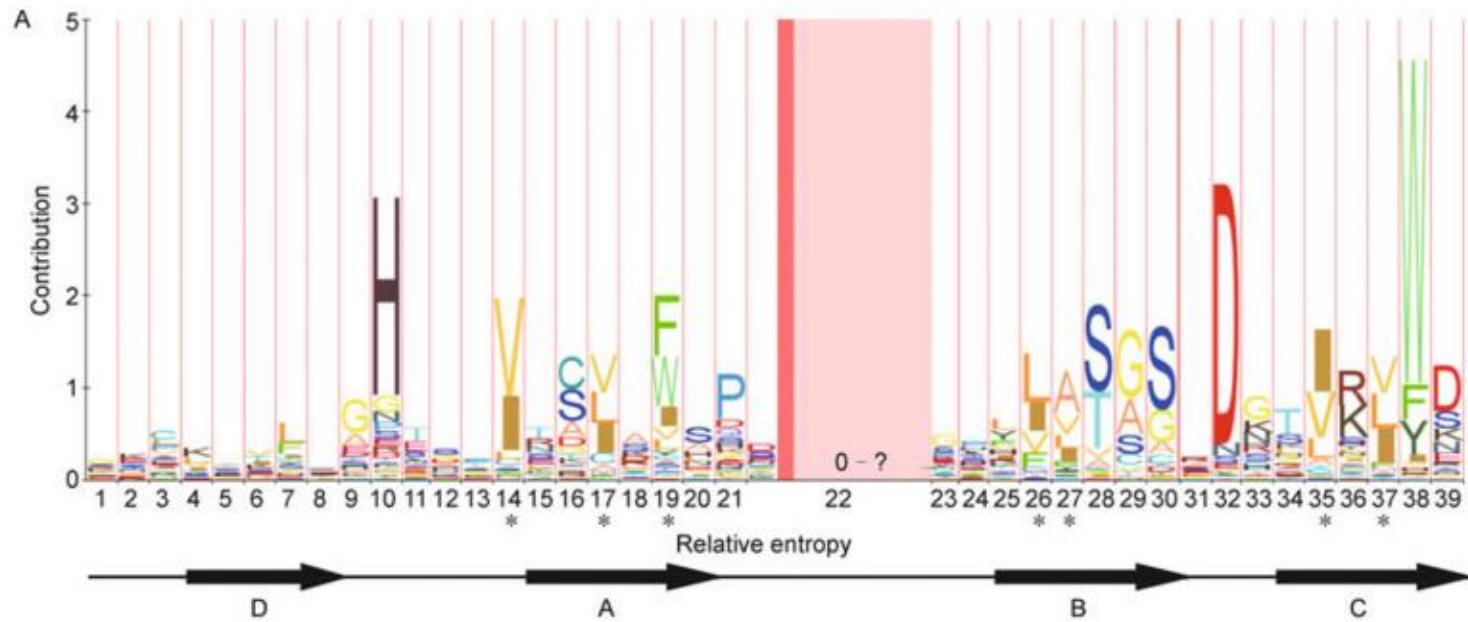
## Special notations with \

Single characters		“Zero-width assertions”	
<code>\t</code>	tab	<code>\b</code>	“word” boundary
<code>\n</code>	newline	<code>\B</code>	not a “word” boundary
<code>\r</code>	return (CR)		
<code>\xhh</code>	character with hex. code <i>hh</i>		

## Matching

<code>\w</code>	matches any <i>single</i> character classified as a “word” character (alphanumeric or “_”)
<code>\W</code>	matches any non-“word” character
<code>\s</code>	matches any whitespace character (space, tab, newline)
<code>\S</code>	matches any non-whitespace character
<code>\d</code>	matches any digit character, equiv. to <code>[0-9]</code>
<code>\D</code>	matches any non-digit character

To present a metacharacter as a data character standing for itself, precede it with `\` (e.g. `\.` matches the full stop character `.` only).



```
$ wget ftp://ftp.ensembl.org/pub/release-98/fasta/homo_sapiens/pep/Homo_sapiens.GRCh38.pep.all.fa.gz
$ gunzip Homo_sapiens.GRCh38.pep.all.fa.gz
$ perl fasta-single-line.pl Homo_sapiens.GRCh38.pep.all.fa >Homo_sapiens.GRCh38.pep.all.fa2
$ grep -B1 -P '[A-Z]{8}GH[A-Z]{21}D[A-Z]{5}WD' Homo_sapiens.GRCh38.pep.all.fa2 >WD40-proteins-Hsa.fa
```

Bring your issues on!