

# Text Analysis in R

## Lecture 1: Introduction

---

Helge Liebert

## Basics

---

# Which programming language to choose?

- *It depends.* Choice is use-case- and taste-specific.
- *Anything* can be done in *any* language. Convenience varies.
- Concepts and toolkits transfer easily most of the time.
- Prior knowledge vs. task suitability.
- Ease of exploratory analysis vs. ease of deployment in production.

# Possible options

- Specialized languages: R, Julia, MATLAB/Octave, Stata, Gauss, ...
- General-purpose languages: Python, Perl, Ruby, C, ...
- Choose a high-level, dynamic, interpreted language unless you are sure you require the extra speed of a compiled language.
- Ideally free and open source. Popular is typically better.
- Research ex ante which libraries are mature and best for solving your specific problem.

- R is the major statistical programming language.
- It is free, used in many sciences and in industry. Good documentation.
- New models are frequently published and implemented first in R.
- Having data processing and analysis in the same language is nice.
- Good library support for common tools (e.g. databases, regular expressions).
- Specific tasks for which high-level wrapper functions are not available may be very cumbersome.
- In recent years, R development has been very active and libraries exist for almost anything.

- General-purpose programming language, supports object-oriented programming.
- Reads like english. Explicit and clear. Whitespace matters, no braces. (*“There should be one obvious way to do it”.*)
- Used extensively in industry and sciences. Good documentation.
- Libraries for almost anything.
- Many science-related libraries exist for other languages, but rarely are they as mature.
- Good and growing support for statistical modeling.
- A bit less suited for interactive data work (but more so for deployment in production).

# This lecture

- The lab sessions utilize R.
- Any task covered by this lecture can be accomplished using R or Python (augmented by shell programs).
- R, Python, SQL and knowing your way around a terminal are highly valued skills on the job market.
- Rule-of-thumb recommendation:
  - Simple data analysis/small text corpora:  
Stick with R. Augment with other tools where required.
  - More involved data processing/larger text corpora:  
Go with Python. You can still analyze data in R.
  - ... and whatever program your colleagues are using.

## Why not Stata, Matlab, Gauss or similar?

- Advantage: Many domain-specific models supported.
- Less support for almost anything else (including text processing).
- Much less flexible for anything not to do with data analysis or numerics.
- Difficult to deploy on a server. Often tied to a GUI.
- Less popular, smaller userbase. Proprietary and expensive.



# A note on text editing

- A script is a set of *plain text* instructions, fed to an interpreter.
- Editing is independent from running code.
- R scripts usually have the suffix `.r`, Python `.py`, Shell `.sh`.
- Proficiency in a text editor makes working with text easier and faster.
- Text editors allow you to integrate your work and edit text efficiently.
- VS Code, Atom, Sublime Text, Vim, Emacs, ...

# A possible setup

The screenshot shows an Emacs editor window with a LaTeX document titled "slides.tex". The document is divided into two main sections: "A note on operating systems" and "Examples".

**A note on operating systems**

- Item MacOS or Linux offer built-in access to a Unix shell (Bash).
- Item Further software is managed via a package management system and distributed via software repositories.
- Item On Linux, use your package manager to install anything you require.
- Item On MacOS, familiarize yourself with Homebrew. Install `iterm2` if you want a fancier terminal.
- Item For Windows, many tools are not available or cumbersome to use. Dependency resolution can be a nightmare.
- Item Windows does not provide proper access to a Unix shell.
- Item Even reliably installing Python was a chore until recently (now use Anaconda).

**Examples**

- Some examples:
  - `vim myscript.r` # edit your R script with vim
  - `R -f myscript.r` # execute your R script
  - `python myscript.py` # execute your python script
  - `git add myscript.py` # stage file for version control
  - `git commit myscript.py` # stage file for version control
  - `man ssh` # display manual pages for the ssh program
  - `ssh myusername@13.438.14.673` # secure shell login to your remote server
- Sounds tedious? It is. But it can also be extremely powerful.
- Convert all your pdf files in a folder to text and search them.
  - `for file in *.pdf; do pdftotext "$file"; done`
  - `grep -icr "keyword" *.txt`

**A note on text editors**

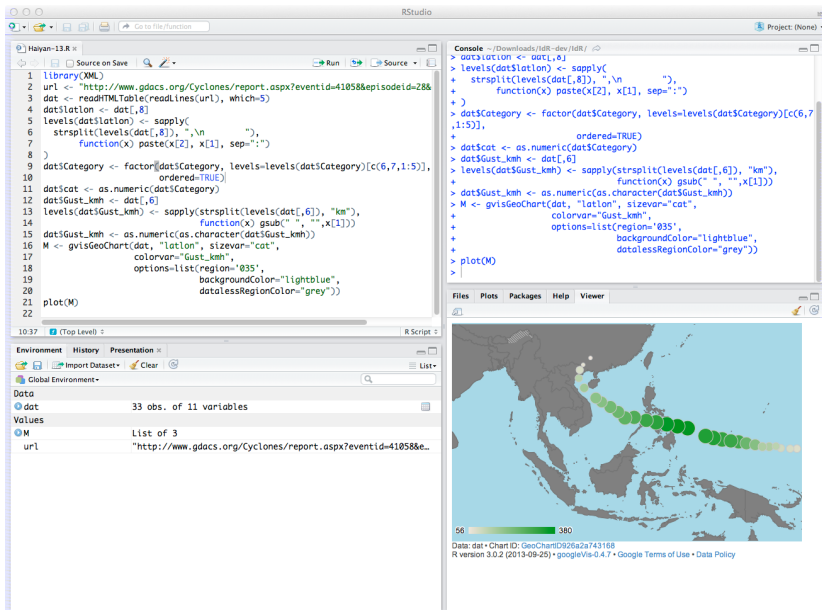
- A script is a set of *plain* text instructions, fed to an interpreter.
- R scripts usually have the suffix `.r`, Python `.py`, Shell `.sh`.
- Much of our work involves working with text files.
- Some text editor is required. A good text editor makes working with text much easier and faster.
- Too many options to list. All are better than Notepad.
- A few suggestions: VS Code, Sublime Text, Atom, Notepad++.
- Learning Vim or Emacs requires you to invest some time.
- Text editors allow you to integrate your work.
- Sometimes IDEs with GUI may be more convenient.
- Features: Efficient text editing, syntax checking, completion, ...

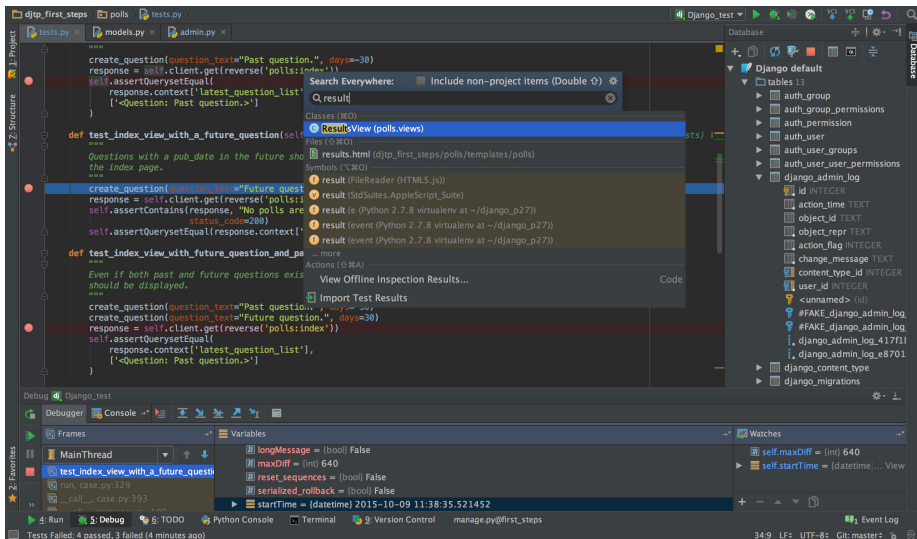
# ... that is universal

```
emacs@helge-x250 ~
1 ## Get maxpages and other information for iteration
2 response <- fromJSON(url, flatten = TRUE)
3 response$pages
4 maxpages <- response$pages$pages
5 records <- response$pages$total
6 columns <- ncol(response$loans)
7
8 ## Open csv, write header
9 header <- names(response$loans)
10 write.table(t(header), file = "data/kliva.csv", sep = ";",
11             col.names = FALSE, row.names = FALSE)
12
13 ## Or collect in data frame (don't do this for large jobs)
14 ## data <- data.frame(matrix(nrow = 0, ncol = columns))
15 ## names(data) <- header
16
17 ## Simple helper function to flatten columns
18 unnest <- function(col) paste(unlist(col), collapse = ", ")
19
20 ## Iterate over pages, limit to first three
21 for (p in seq(1, maxpages, by = 1)[1:3]) {
22   ## Info
23   print(paste0(p, "/", maxpages))
24
25   ## Append page to url
26   pquery <- paste0(url, "&page=", p)
27
28   ## Get data, assert completeness
29   loans <- fromJSON(pquery, flatten = TRUE)$loans
30   stopifnot(nrow(loans) == pagelength)
31   stopifnot(ncol(loans) == columns)
32
33   ## Fix nested list columns ... or just use data.table::fwrite()
34   ## str(loans)
35   loans$tags <- sapply(loans$tags, unnest)
36   ## loans$themes <- sapply(loans$themes, unnest) # missing for older records
37   loans$description_languages <- sapply(loans$description_languages, unnest)
38   ## str(loans)
39
40   ## Collect loans in data frame
41   ## data <- rbind(data, loans)
42
43   ## Append to file
44   write.table(loans, "data/kliva.csv", sep = ";", append = TRUE,
45               col.names = FALSE, row.names = FALSE)
46 }
47
48 ## head(data)
49 ## dim(data)
50
51 ## They work for you, can Example
52 apikey <- "C3WqTbTKAbdQVqrdbBKAjnb"
53 use <- "https://www.theworkforyou.com/api/"
54 format <- "%s"
55
56 ## 5.0 rest-kliva ESS [R] plyr/flymakeR/t[0 0]e[W]mc[K]L
57 loading line: apikey <- "C3WqTbTKAbdQVqrdbBKAjnb"
58
59 6 Vulnerable Groups
60 location.country_code location.country location.geo.level
61 1 KH Cambodia Kampong Cham town
62 2 VN Vietnam Thanh Hoà town
63 4 VN Vietnam Thanh Hoà town
64 4 VN Vietnam Thanh Hoà town
65 5 VN Vietnam Thanh Hoà town
66 6 VN Vietnam Thanh Hoà town
67
68 location.geo.pairs location.geo.type
69 1 12 105.5 point
70 2 19.806692 105.785182 point
71 3 19.806692 105.785182 point
72 4 19.436971 105.374762 point
73 5 19.806692 105.785182 point
74 6 19.806692 105.785182 point
75
76 > [1] 20 25
77
78 >>>>>>> query <- paste0("country_code=", country, "&",
79 + "sectors=", sector, "&",
80 + "borrower_type=", type, "&",
81 + "status=", status, "&",
82 + "sort_by=", sortby)
83 > url <- paste0(baseUrl, method, query)
84 > response <- fromJSON(url, flatten = TRUE)
85 response$pages
86 maxpages <- response$pages$pages
87 records <- response$pages$total
88 columns <- ncol(response$loans)
89 > $page
90 [1] 1
91
92 $total
93 [1] 3361
94
95 $page_size
96 [1] 20
97
98 $pages
99 [1] 169
100
101 >>>> header <- names(response$loans)
102 > write.table(t(header), file = "data/kliva.csv", sep = ";",
103             col.names = FALSE, row.names = FALSE)
104 > unnest <- function(col) paste(unlist(col), collapse = ", ")
105 > for (p in seq(1, maxpages, by = 1)[1:3]) {
106   > print(paste0(p, "/", maxpages))
107   > pquery <- paste0(url, "&page=", p)
108   > loans <- fromJSON(pquery, flatten = TRUE)$loans
109   > stopifnot(nrow(loans) == pagelength)
110   > stopifnot(ncol(loans) == columns)
111   > loans$tags <- sapply(loans$tags, unnest)
112   > loans$description_languages <- sapply(loans$description_languages, unnest)
113   > write.table(loans, "data/kliva.csv", sep = ";", append = TRUE,
114               col.names = FALSE, row.names = FALSE)
115 }
116
117 > apikey <- "C3WqTbTKAbdQVqrdbBKAjnb"
118 [1] "1/169"
119 [1] "2/169"
120 [1] "3/169"
121 >>>
```

# ... that is universal

<p>each school-track, classroom formation is conditionally ignorable with respect to SN status, as students from different primary school districts are mixed and their SN status is not observed by secondary school administrators.<sup>\footnote{Using PISA data from secondary schools in Switzerland, \cite{Vardardottir2015} shows that track-by-school fixed effects render peer group composition conditionally uncorrelated with a large set of students' characteristics, while track fixed effects and school fixed effects do not.} Neither primary schools nor the SPS share information with secondary schools for equity reasons and to avoid stigma when transitioning between schools.</sup></p> <p>174 % NEU: Bei bitte lesen 175 % changed some small things. OK for me. 176 Beyond this anecdotal evidence, we formally test the validity of the identification strategy with four balancing tests, which are presented and discussed extensively in Appendix B. First, we examine whether the proportion of SN peers predicts individual baseline characteristics (gender, native speaker, and age). The aim of this test is to detect potential selection into classrooms. We also conduct this test separately for SN and non-SN students. None of the baseline characteristics are statistically significant at conventional levels, either considered individually or jointly. Second, we regress the indicator for SN status on class fixed effects, which should be jointly insignificant if assignment to <b>classrooms</b> is ignorable with respect to SN status \cite{chettyEtal2011}. We also conduct this test to check for ignorable assignment of SN students to teachers. We find no evidence for systematic assignment of SN students to either classes or teachers. Third, we conduct a simulation exercise in the spirit of \cite{carrell2010}. We re-sample classes and thereby assign SN students randomly to <b>classrooms</b>, and test whether the observed distribution of SN students differs from the simulated one. In addition, we compute the interquartile range of the proportion SN students across classes for each simulation, and compare these simulated interquartile ranges with the one we observe in the data. Neither simulation procedure uncovers any worrisome pattern in the assignment of SN students to classes. Fourth, we decompose the variation in the fraction of SN peers across and within schools. To do so, we examine the residual variation in the proportion of SN peers after partialling out the school-track-year fixed effects. We find that the residual distribution in the proportion of SN peers is consistent with variation from a random process. Overall, the balancing tests we performed indicate that the key identification assumption of (conditionally) ignorable assignment of SN students to classes is plausible.</p> <p>178 Our identification relies on variation between classes within school-track-years. Although families can potentially choose their district of residence and thereby influence schooling options for their children, possible selection into schools does not confound our results.<sup>\footnote{Endogenous class formation could still occur if parents request to transfer their children to a class with a lower SN fraction. To investigate this potential threat, we acquired the official education statistics from the Swiss Federal Statistical Office (SOL, \textit{Statistik der Lernenden} in German) for the years 2012-2015. Importantly, the SOL has a classroom ID which allows us to reconstruct the classes within each school-year-track. For the state of St. Gallen we find that no</sup></p> <p>A: ● 104k xSource/manuscript_R1_v3.tex 77:0 37% LF UTF-8 LaTeX/FPS ○ 1</p> <p>✎ EditDiff Control Panel* diff. 4 of 15 Quick Help Auto-refining is ON</p>	<p>different education tracks, and classes are strictly separated between tracks. Within each school-track, classroom formation is conditionally ignorable with respect to SN status, as students from different primary school districts are mixed and their SN status is not observed by secondary school administrators.<sup>\footnote{Using PISA data from secondary schools in Switzerland, \cite{Vardardottir2015} shows that track-by-school fixed effects render peer group composition conditionally uncorrelated with a large set of students' characteristics, while track fixed effects and school fixed effects do not.} Neither primary schools nor the SPS share information with secondary schools for equity reasons and to avoid stigma when transitioning between schools.</sup></p> <p>172 % NEU: Bei bitte lesen 173 Beyond this anecdotal evidence, we formally test the validity of the identification strategy with four balancing tests, which are presented and discussed extensively in Appendix B. First, we examine whether the proportion of SN peers predicts individual baseline characteristics (gender, native speaker, and age). The aim of this test is to detect potential selection into classrooms. We also conduct this test separately for SN and non-SN students. None of the baseline characteristics are statistically significant at conventional levels, either considered individually or jointly. Second, we regress the indicator for SN status on class fixed effects, which should be jointly insignificant if assignment to <b>classroom</b> is ignorable with respect to SN status \cite{chettyEtal2011}. We also conduct this test to check for ignorable assignment of SN students to <b>specific</b> teachers. We find no evidence for systematic assignment of SN students to either classes or teachers. Third, we conduct a simulation exercise in the spirit of \cite{carrell2010}. We re-sample classes, randomly assign SN students to <b>classes</b>, and test whether the observed distribution of SN students differs from the simulated one. In addition, we simulate random classroom assignment within schools, compute the interquartile range of the proportion SN across classes, and compare the simulated interquartile range with the one we observe in the data. Neither simulation procedure uncovers any worrisome pattern in the assignment of SN students to classes. Fourth, we decompose the variation in the fraction of SN peers across and within schools. To do so, we examine the residual variation in the proportion of SN peer after partialling out the school-track-year fixed effects. We find that the residual distribution in the proportion of SN peers is consistent with variation from a random process. Overall, the balancing tests we performed indicate that the key identification assumption of (conditionally) ignorable assignment of SN students to classes is plausible.</p> <p>175 Our identification relies on variation between classes within school-track-years. Although families can potentially choose their district of residence and thereby influence schooling options for their children, possible selection into schools does not confound our results.<sup>\footnote{Endogenous class formation could still occur if parents request to transfer their children to a class with a lower SN fraction. To investigate this potential threat, we acquired the official education statistics from the Swiss Federal Statistical Office (SOL, \textit{Statistik der Lernenden} in German) for the years 2012-2015. Importantly, the SOL has a classroom ID which allows us to reconstruct the classes within each school-year-track. For the state of St. Gallen we find that no</sup></p> <p>B: ● 104k x010/manuscript_R1_v2.tex 77:0 37% LF UTF-8 LaTeX/FPS ○ 11</p> <p>type ? for help</p>
----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------





## A note on text editing

- Text editors allow you to integrate your work and edit text efficiently.
- Sometimes IDEs with GUI may be more convenient.
- Features: Regex search and replace, diff, syntax checking, formatting, completion, persistent undo, documentation lookup, version control support ...
- Consider using version control. Git is the predominant version control software used today. [ProGit](#) is a good and free resource.

# Regular expressions

---



# Introduction

- Text data is now commonly used as a modeling input in both research and industry.
- Often involves substantial amounts of data preparation.
- Text processing is a common task in IT - many helpful tools available.
- Working with text requires understanding **regular expressions**.
- Language-independent concept.

# Regular expressions

- A regular expression is a character sequence that describes a set of strings.
- Regular expressions are constructed analogously to arithmetic expressions, by using various operators to combine smaller expressions.
- Usually used for find/replace operations on strings, or for validation.
- Pervasive in Unix text processing programs (**grep** was originally written by Ken Thompson).
- **grep**: globally search a regular expression and print.
- Not limited to command line search tools.

# Regular expressions

- *Pattern matching*: Find one of a specified set of strings in text.
- Examples:
  - Diagnoses in medical records.
  - Addresses or zip codes in concatenated admin records.
  - Sequences within a genome, e.g. a virus signature.
  - Validate data-entry fields (URL, date, email, credit card #).
  - [Example using a regex tester](#).

# Obligatory xkcd

WHENEVER I LEARN A NEW SKILL I CONCOCT ELABORATE FANTASY SCENARIOS WHERE IT LETS ME SAVE THE DAY.

OH NO! THE KILLER MUST HAVE FOLLOWED HER ON VACATION!



BUT TO FIND THEM WE'D HAVE TO SEARCH THROUGH 200 MB OF EMAILS LOOKING FOR SOMETHING FORMATTED LIKE AN ADDRESS!



IT'S HOPELESS!

EVERYBODY STAND BACK.



I KNOW REGULAR EXPRESSIONS.



# Examples

AHVN 13: `756\.[0-9]{4}\.[0-9]{4}\.[0-9]{4}\.[0-9]{2}`

Matches: `756.1234.5678.90`

Does not match: `123.45.678.675`

US-SSN: `[0-9]{3}-[0-9]{2}-[0-9]{4}`

Matches: `166-11-4433`

Does not match: `11-55555555`

Email addresses: `[a-z]+@([a-z]+\.)+(ch|edu|com)`

Matches: `someone@unibas.ch`

Does not match: `someone@invalid.domain`

# Screening job candidates

“ [First name]! and pre/2 [last name] w/7  
bush or gore or republican! or democrat! or charg!  
or accus! or criticiz! or blam! or defend! or iran contra  
or clinton or spotted owl or florida recount or sex!  
or controversies! or fraud! or investigat! or bankrupt!  
or layoff! or downsiz! or PNTR or NAFTA or outsourc!  
or indict! or enron or kerry or iraq or wmd! or arrest!  
or intox! or fired or racis! or intox! or slur!  
or controversies! or abortion! or gay! or homosexual!  
or gun! or firearm! ”

— *LexisNexis search string used by Monica Goodling  
to illegally screen candidates for DOJ positions*



LexisNexis™

<http://www.justice.gov/oig/special/s0807/final.pdf>

# Regular expressions

- Characters in a regular expression are either regular characters (literal meaning) or metacharacters (special meaning).
- Generally, letters and numbers match themselves.
- Normally case sensitive, but can be set to ignore case.
- Careful with punctuation, most of it has special meanings.
- To match metacharacters literally, they need to be *escaped*, i.e. preceded by a backslash \.

# Matching string literals

Regular expression	Input string
<code>input</code>	This <code>input</code> string is short.
<code>15</code>	The due date is <code>15.12.</code>



# Matching string literals

Regular expression	Input string
<code>input</code>	This <code>input</code> string is short.
<code>15</code>	The due date is <code>15.12</code> .
<code>15.12</code>	The due date is <code>15.12</code> .

# Matching string literals

Regular expression	Input string
<code>input</code>	This <code>input</code> string is short.
<code>15</code>	The due date is <code>15.12</code> .
<code>15.12</code>	The due date is <code>15.12</code> .
but:	
<code>15.12</code>	The due date is <code>15712</code> .
match <code>.</code> literal:	
<code>15\.12</code>	The due date is <code>15712</code> .
<code>15\\.12</code>	The due date is <code>15.12</code> .

# Regex basics

Operation	Regular expression	Input string
concatenation	<code>foobar</code>	Matches <code>foobar</code> but not <code>foo</code> or <code>bar</code> .
disjunction	<code>this that</code>	Matches <code>this</code> or <code>that</code> .
closure (repetition)	<code>like.* apples</code> <code>like. apples</code>	<code>I like apples, Peter likes apples.</code> <code>Mary also likesALKFHEDL apples.</code> <code>Mary also likesALKFHEDL apples.</code>
grouping	<code>(He She) likes</code> <code>(He She).*(very )*much\.</code>	<code>She likes apples. He likes apples.</code> <code>She likes apples very very much.</code>

- More complicated patterns can be expressed via concatenation, disjunction, repetition and scope.
- Precedence in descending order.

# Quantifiers

Character	Matches
*	0 or more instances of preceding char
+	1 or more instances of preceding char
?	0 or one instance of preceding char
{m}	exactly m instances of preceding char
{m,n}	m through n instances of preceding char
{m,}	m or more instances of preceding char
{,n}	up to n instances of preceding char
?	add to a quantifier to match ungreedy

- Quantifiers match greedily by default (i.e. the longest string possible).

Ex: `^begin.*end` will match `'begin bla bla end bla end'`.

`^begin.*?end` will match `'begin bla bla end bla end'`.

# Groups, ranges and character classes

Character	Matches	Example RE	Matches
.	Any character, except \n	like.	likes like! like like
(a b)	a or b	(you me)	you or me
[ab]	Character range	202[01]	2019 2020 2021 2022
[a-z]	Character range	[A-Z][a-z]*	Capitalized words
[0-9]	Digit range	20[0-9]{2}	Years in the 21st century
[^ab]	Any character but (negation)	20[^0][01]	2000 2010 2020 2025 2031

- Quantifiers, ranges and other shortcuts improve expressiveness.

Ex: `[A-E]+` is shorthand for `(A|B|C|D|E)(A|B|C|D|E)*`.

- More character classes (sometimes) available.

Ex: `\w` for words (`[A-Za-z0-9_]`), or `\d` for digits (`[0-9]`), `\a`, `\s`, ...

# Anchors and other special characters

Character	Matches	Example RE	Matches
<code>^</code>	Beginning of a string	<code>^New</code>	New research in this field
<code>\$</code>	End of a string	<code>[A-Za-z]+!\$</code>	A breakthrough! Finally!
<code>\n</code>	Newline		
<code>\t</code>	Tab		
...	...		

- Strings can stretch multiple lines.
- Character encodings can sometimes cause problems. Stick to UTF-8.

# Regex syntaxes

- More elaborate regex syntaxes also support positive and negative lookahead/lookbehind, conditionals and group references.  
Ex: `^(?!.*word).*` matches lines not containing a word.
- Different syntaxes (basic, extended, perl, vim, ...) mostly similar with regard to basic features.
- Perl-compatible regular expressions (PCRE) is the de-facto standard.
- Most regex implementations feature switches to invert the search pattern, to ignore case, and more.

## Example: Valid RFC-822 email addresses

[illegible]



- Writing a regular expression is like writing a program.
  - Requires understanding the programming model.
  - Can be easier to write than read.
  - Can be difficult to debug.
- ➡ Break up problems into smaller pieces. Try not to do everything in one large regex. Comment liberally.

# Resources

- Pin a [cheat sheet](#) to your office wall.
- Regular expression tools help (e.g. [regex101.com](#)).
- Simple interactive tutorial: [regexone.com](#). Also [learn regex the easy way](#).
- For regular expression syntax specific to R, look up [this short tutorial](#).
- [This guide](#) also provides a more detailed overview of working with strings in R.

- Regexes are a powerful tool.
- Easy to grasp, complex to master.
- Using them in applications can be complex and error-prone.
- Regular expressions are not parsers.

*“Some people, when confronted with a problem, think ‘I know, I’ll use regular expressions.’ Now they have two problems.”*

*Emacs newsgroup*

Next lecture: Representing text as data

## References

---

# References



Chacon, S. and B. Straub (2014). *Pro Git*. Apress.



Fitzgerald, M. (2012). *Introducing Regular Expressions: Unraveling Regular Expressions, Step-by-Step*. 1. ed. Beijing: O'Reilly.



Matloff, N. (2011). *The Art of R Programming: A Tour of Statistical Software Design*. No Starch Press.



Shotts, W. E. (2019). *The Linux Command Line: A Complete Introduction*. Second edition. San Francisco: No Starch Press.