

# Regular Expressions

Greg Ridgeway      Ruth Moyer

2026-02-10

## Table of contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Finding patterns in text with grep()</b>	<b>2</b>
2.1	Find letters and numbers . . . . .	3
2.1.1	Exercises . . . . .	7
2.2	More symbols that help with grepping . . . . .	7
2.2.1	Carets ^ . . . . .	7
2.2.2	Dollar Signs \$ . . . . .	8
2.2.3	Plus Sign + . . . . .	8
2.2.4	Exercises . . . . .	9
2.2.5	Parentheses () . . . . .	9
2.2.6	Vertical Bar   . . . . .	9
2.2.7	Question Mark ? . . . . .	11
2.2.8	Boundaries \b . . . . .	12
2.2.9	Exercises . . . . .	12
<b>3</b>	<b>Finding and replacing patterns in text with gsub()</b>	<b>12</b>
3.1	Basic find and replace examples . . . . .	12
3.1.1	Exercises . . . . .	14
3.2	Back-Referencing . . . . .	14
3.2.1	Exercise . . . . .	16
3.3	Lookaheads . . . . .	17
3.3.1	Exercises . . . . .	19
<b>4</b>	<b>Introduction to Webscraping: A Practical Application of Regular Expressions</b>	<b>19</b>
4.1	The 10,000 most common English words . . . . .	19
4.1.1	Exercises . . . . .	22
<b>5</b>	<b>Solutions to the exercises</b>	<b>23</b>

## 1 Introduction

A regular expression is a sequence of characters that defines a search pattern. Sometimes we use regular expressions to help us find a pattern in text that we want, like the Find functionality in Word. Other times, we use a regular expression to help us find and *replace* a piece of text, like the Find and Replace functionality in Word. The two main R functions that we will learn about are `grep()` and `gsub()`. You may use them in verb form (e.g., “I need to grep all the gun cases,” “I was grepping like crazy to find all the opioid cases,” “I first had to gsub the commas with semicolons”). Regular expressions are available in numerous other software packages, so everything you learn here will port over to using regular expressions in Linux, Stata, Python, Java, ArcGIS, and many others.

We have already used regular expressions in previous work that we have done. We wanted to select protests in Philadelphia and ran

```
dataProtest |>  
  filter(grepl("Philadelphia, PA", Location))
```

When working with NIBRS we cleaned up the column names that had multiple periods in them.

```
fmtNames <- fmt$Description |>  
  strsplit(split = " - ", fixed = TRUE) |>  
  sapply(head, n=1) |>  
  make.names() |>  
  gsub("\\\\.+\\\\.|\\$)", "\\\\1", x=_)
```

To estimate the number of sexual assault victimizations from the NCVS we selected all the crime labels that had “rape” or “sex” in them.

```
dataExt |>  
  filter(grepl("rape|[Ss]ex", V4529)) |>  
  distinct(V4529)
```

All of this will become clear in these notes.

## 2 Finding patterns in text with `grep()`

First we will learn about `grep()`, which was first developed in the early 1970s. `grep()` searches data for lines that match a given expression. The name `grep` stands for “globally search a regular expression and print.” Some tangible examples will help us see how `grep()` works.

Run the following line that provides a collection of text elements with a diverse range of capitalization, letters, punctuation, and other features. We will use this for practice and testing.

```
dataText <- c("3718 Locust Walk",
            "4 Privet Drive",
            "10880 Malibu Point",
            "221B Baker St.",
            "19104-6286",
            "20015",
            "90291",
            "90210",
            "(215) 573-9097",
            "215-573-9097",
            "2155739097",
            "Kendall Roy",
            "Roman Roy",
            "Siobhan Roy",
            "Roy Wood, Jr.",
            "Royal Caribbean",
            "Casino Royale",
            "two children",
            "2 children",
            "twins",
            "Philadelphia",
            "Philly",
            "Phila",
            "Dr. Phil",
            "$23456",
            "$10000",
            "$60,000")
```

## 2.1 Find letters and numbers

Let's find all of the items that have the letter "a" in dataText.

```
grep("a", dataText)
```

```
[1] 1 3 4 12 13 14 16 17 21 23
```

This shows the indices of which elements of `dataText` have an “a” in them. Sure enough the first element, “3718 Locust Walk”, has an “a” and you can confirm that the other indices `grep()` returned also indicate elements of `dataText` that have an “a”.

For most of these notes we are going to set `value=TRUE`, which will return the actual elements matched rather than their indices. This will make it easier to check that `grep()` is finding the elements that we are expecting to match.

```
grep("a", dataText, value=TRUE)
```

```
[1] "3718 Locust Walk"    "10880 Malibu Point" "221B Baker St."
[4] "Kendall Roy"        "Roman Roy"           "Siobhan Roy"
[7] "Royal Caribbean"   "Casino Royale"      "Philadelphia"
[10] "Phila"
```

As you can see, `grep()` uses the following general syntax: `grep(“what we’re searching for”, the text, and an optional value=TRUE if we want to return text)`. If we put `value=FALSE` (which is the default if we do not set `value`), we will just receive the index in the text where we can find what we are searching for.

Let’s try another example. Instead of a letter, let’s try to find a number such as “1”. Specifically, which items in `dataText` have a “1” in them?

```
grep("1", dataText, value=TRUE)
```

```
[1] "3718 Locust Walk"    "10880 Malibu Point" "221B Baker St."
[4] "19104-6286"          "20015"                 "90291"
[7] "90210"                "(215) 573-9097"     "215-573-9097"
[10] "2155739097"         "$10000"
```

We can also search for multiple characters, instead of individual characters. We place our list of desired characters within square brackets [ ]. For example, let’s find items that contain numbers.

```
grep("[0123456789]", dataText, value=TRUE)
```

```
[1] "3718 Locust Walk"    "4 Privet Drive"      "10880 Malibu Point"
[4] "221B Baker St."     "19104-6286"        "20015"
[7] "90291"               "90210"              "(215) 573-9097"
[10] "215-573-9097"       "2155739097"       "2 children"
[13] "$23456"              "$10000"            "$60,000"
```

If we wanted items that contain an odd number, we could do this.

```
grep("[13579]", dataText, value=TRUE)
```

```
[1] "3718 Locust Walk"      "10880 Malibu Point"  "221B Baker St."
[4] "19104-6286"           "20015"                  "90291"
[7] "90210"                 "(215) 573-9097"    "215-573-9097"
[10] "2155739097"          "$23456"                "$10000"
```

So the [] in a regular expression means “match any of these characters.” We can also place square brackets next to each other. For example, let’s say we wanted to find an item in `dataText` that has four adjacent numbers.

```
grep("[0-9][0-9][0-9][0-9]", dataText, value=TRUE)
```

```
[1] "3718 Locust Walk"      "10880 Malibu Point"  "19104-6286"
[4] "20015"                 "90291"                  "90210"
[7] "(215) 573-9097"       "215-573-9097"    "2155739097"
[10] "$23456"                "$10000"
```

Note that we can use the shorthand 0-9 to mean any number between 0 and 9, including 0 and 9. This regular expression says “find a number, followed by a number, followed by another number, followed by another number.” Alternatively we can use {4} to mean “match four of the previous character”.

```
grep("[0-9]{4}", dataText, value=TRUE)
```

```
[1] "3718 Locust Walk"      "10880 Malibu Point"  "19104-6286"
[4] "20015"                 "90291"                  "90210"
[7] "(215) 573-9097"       "215-573-9097"    "2155739097"
[10] "$23456"                "$10000"
```

{n} means the preceding item will be matched exactly n times. Note that this also matches text that has five or more numbers in a row, since if they have five numbers in a row, then they will also have four numbers in a row. Later we will learn about how to find exactly four numbers in a row.

It is not used often, but we can also use the squiggly brackets to make a range.

```
grep("[0-9]{5,10}", dataText, value=TRUE)
```

```
[1] "10880 Malibu Point" "19104-6286"      "20015"  
[4] "90291"              "90210"            "2155739097"  
[7] "$23456"             "$10000"
```

Thus `{n,m}` will find something matched between n and m times. We have used examples above with numbers, but you can apply this syntax to letters.

```
grep("[a-zA-Z]{5}", dataText, value=TRUE)
```

```
[1] "3718 Locust Walk"    "4 Privet Drive"    "10880 Malibu Point"  
[4] "221B Baker St."     "Kendall Roy"       "Roman Roy"  
[7] "Siobhan Roy"        "Royal Caribbean"  "Casino Royale"  
[10] "two children"       "2 children"       "twins"  
[13] "Philadelphia"       "Philly"           "Phila"
```

Note how we used the shorthand `a-z` to mean any lowercase letter and `A-Z` to mean any uppercase letter. `grep()` will match a different set of elements if we just search for lower case letters.

```
grep("[a-z]{5}", dataText, value=TRUE)
```

```
[1] "3718 Locust Walk"    "4 Privet Drive"    "10880 Malibu Point"  
[4] "Kendall Roy"        "Siobhan Roy"       "Royal Caribbean"  
[7] "Casino Royale"       "two children"      "2 children"  
[10] "twins"              "Philadelphia"     "Philly"
```

Searches using `grep()` are case-sensitive. For example, let's find all items that contain capital letters.

```
grep("[A-Z]", dataText, value=TRUE)
```

```
[1] "3718 Locust Walk"    "4 Privet Drive"    "10880 Malibu Point"  
[4] "221B Baker St."     "Kendall Roy"       "Roman Roy"  
[7] "Siobhan Roy"        "Roy Wood, Jr."    "Royal Caribbean"  
[10] "Casino Royale"      "Philadelphia"     "Philly"  
[13] "Phila"              "Dr. Phil"
```

### 2.1.1 Exercises

1. Find text with vowels
2. Find text with a number immediately followed by a letter
3. What is the difference between `grep("Roy", dataText, value=TRUE)` and `grep("[Roy]", dataText, value=TRUE)`?

## 2.2 More symbols that help with grepping

So far you have seen that we use [] to match any character listed between the [], the - to specify a range to match like 0-9, a-z, and A-Z, and the {} to match the previous character multiple times.

In addition to these, there are several more symbols or sequences of symbols that are useful. The number of symbols on the keyboard are fairly limited compared to the numerous combinations of patterns of text we might wish to find. As a result you will see that some of these symbols are used in very different ways depending on the context.

### 2.2.1 Carets ^

Let's look at the caret ^ symbol. When we put a ^ within square brackets, it means "not". Let's try to find text in `dataText` that has something that is not a letter or a space immediately followed by a letter.

```
grep("[^A-Za-z] [A-Za-z]", dataText, value=TRUE)
```

```
[1] "221B Baker St."
```

"221B Baker St." has a character that is not a letter (1) immediately followed by a letter (B).

Outside of square brackets the ^ means something *completely* different. When we put the ^ outside of square bracket, it means "the beginning of the text". For example, the following regular expression matches text where the first character is either an upper-case or a lower-case letter.

```
grep("^ [A-Za-z]", dataText, value=TRUE)
```

```
[1] "Kendall Roy"      "Roman Roy"       "Siobhan Roy"      "Roy Wood, Jr."
[5] "Royal Caribbean" "Casino Royale"   "two children"    "twins"
[9] "Philadelphia"     "Philly"          "Phila"           "Dr. Phil"
```

### 2.2.2 Dollar Signs \$

While we use carets to signal the beginning of the text, the dollar sign signals the end of the text. For example, to search for all items that end with either an upper-case or lower-case letter, we would do the following:

```
grep("[A-Za-z]$", dataText, value=TRUE)
```

```
[1] "3718 Locust Walk"    "4 Privet Drive"      "10880 Malibu Point"  
[4] "Kendall Roy"        "Roman Roy"          "Siobhan Roy"  
[7] "Royal Caribbean"   "Casino Royale"     "two children"  
[10] "2 children"         "twins"              "Philadelphia"  
[13] "Philly"              "Phila"                "Dr. Phil"
```

Note how all of these have a letter as the last character. We can be more specific and ask for text that end with the letter “y” or end with a 7.

```
grep("y$", dataText, value=TRUE)
```

```
[1] "Kendall Roy" "Roman Roy" "Siobhan Roy" "Philly"
```

```
grep("7$", dataText, value=TRUE)
```

```
[1] "(215) 573-9097" "215-573-9097" "2155739097"
```

### 2.2.3 Plus Sign +

The + means “at least one of the previous.” For example, suppose we wanted to find items in our list that have numbers and those numbers are followed by letters.

```
grep("[0-9]+[A-Za-z]+", dataText, value=TRUE)
```

```
[1] "221B Baker St."
```

Or search for text that starts with some numbers, then has a space, followed by some letters, and then ends.

```
grep("^[0-9]+ [A-Za-z]+$", dataText, value=TRUE)
```

```
[1] "2 children"
```

#### 2.2.4 Exercises

4. Find the Roy children. That is, pick out Kendall, Roman, and Siobhan Roy.
5. Find text that starts with a number and ends with a letter. (Hint: You will probably need to use ^, \$, and +).

#### 2.2.5 Parentheses ()

Parentheses group together characters as words. For example, suppose we wanted to find the word “two”.

```
grep("(two)", dataText, value=TRUE)
```

```
[1] "two children"
```

On its own the parentheses are no different than just running the regular expression “two” with no parentheses. However, in combination with | and ? it gets more interesting and powerful.

#### 2.2.6 Vertical Bar |

The vertical bar functions as an “or.” Suppose we wanted to get both the phrases “two children” and “2 children”, allowing the number to be spelled out or not. We would use parentheses as well as a vertical bar.

```
grep("(two|2) children", dataText, value=TRUE)
```

```
[1] "two children" "2 children"
```

Let’s find all Los Angeles and Washington DC ZIP codes.

```
grep("(902|200)[0-9][0-9]", dataText, value=TRUE)
```

```
[1] "20015" "90291" "90210"
```

Here's how we can find words that have exactly four characters. To be exactly four characters, on either side of them there needs to be a space or the beginning or end of the line.

```
grep("( |^)[A-Za-z]{4}(|$)", dataText, value=TRUE)
```

```
[1] "3718 Locust Walk" "Dr. Phil"
```

( |<sup>^</sup>) says that we start with either a space or the beginning of the text. ( |\$) says that at the end there must be a space or the end of the text. In the middle there has to be exactly four letters. Note that this did not select "Phil." since it has a period following it. We'll need a more general regular expression to select that one as well.

Let's try to find phone numbers in `dataText`. The problem is that phone numbers have three different formats in our data. The easiest phone number pattern to find is the one with 10 digits in a row.

```
grep("[0-9]{10}", dataText, value=TRUE)
```

```
[1] "2155739097"
```

It is also not too hard to find the hyphenated phone number pattern.

```
grep("[0-9]{3}-[0-9]{3}-[0-9]{4}", dataText, value=TRUE)
```

```
[1] "215-573-9097"
```

The phone number format with parentheses needs a little more caution. We pointed out earlier that parentheses have special meaning in regular expressions. In fact, there are several "special characters" in regular expressions, \ ^ \$ {} [] () . \* + ? | -. If you actually want to search for the symbol itself you need to "protect" it with a backslash \\. However, the \\ is a special character for R too and R will think that something special is coming next after the \\. So to tell R "no really, I really want a backslash here" you have to protect the backslash too. So to look for those phone numbers with parentheses we use a regular expression like this.

```
grep("\\\\([0-9]{3}\\\\) [0-9]{3}-[0-9]{4}", dataText, value=TRUE)
```

```
[1] "(215) 573-9097"
```

There is a workaround to the extra backslashes, which can be overwhelming in some regular expressions. An R raw string literal prevents R from interpreting the backslashes as special characters. Wrap your text inside `r"( ... )"` and you can do away with double backslashes.

```
# an R string literal requiring only single backslashes  
r"\\([0-9]{3}\\) [0-9]{3}-[0-9]{4})"
```

```
[1] "\\([0-9]{3}\\) [0-9]{3}-[0-9]{4}"
```

```
# regular expression, now with 50% fewer backslashes!  
grep(r"\\([0-9]{3}\\) [0-9]{3}-[0-9]{4})", dataText, value=TRUE)
```

```
[1] "(215) 573-9097"
```

Now, let's put all of these patterns together and use `|` to search for any phone number in any of the three formats.

```
grep("[0-9]{10}|[0-9]{3}-[0-9]{3}-[0-9]{4}|\\([0-9]{3}\\) [0-9]{3}-[0-9]{4}",  
      dataText, value=TRUE)
```

```
[1] "(215) 573-9097" "215-573-9097" "2155739097"
```

## 2.2.7 Question Mark ?

The question mark indicates optional text. To illustrate,

```
grep("Philadelphia", dataText, value=TRUE)
```

```
[1] "Philadelphia"
```

```
grep("Philadelphia?", dataText, value=TRUE)
```

```
[1] "Philadelphia" "Philly" "Phila" "Dr. Phil"
```

The first one is no different from searching for the word “Philadelphia,” but the second one says that the “adelphia” part is optional. Note that this regular expression also picked up “Dr. Phil”. Again we will need a more careful regular expression to avoid matching this name and only select abbreviations and nicknames for Philadelphia.

### 2.2.8 Boundaries \b

\b will try to find boundaries around words. This includes spaces, punctuation, and the beginning and end of the text. So another way to find all text with four letter words, including “Phil.”, is

```
grep("\\b[A-Za-z]{4}\\b", dataText, value=TRUE)
```

```
[1] "3718 Locust Walk" "Roy Wood, Jr."      "Dr. Phil"
```

Remember, the \ is a special character for R. To “protect” the backslash we put another backslash in front of it. If we did not include the \b in this regular expression, then we would have also matched words with five or more letters too.

### 2.2.9 Exercises

Write regular expressions to find the following items in `dataText`.

6. Find commas
7. Find ZIP codes
8. Find those with six letter word
9. Find mentions of Philadelphia
10. Find capitalized words
11. Find the addresses
12. Find a shorter way to get phone numbers using ?

## 3 Finding and replacing patterns in text with `gsub()`

`gsub()` stands for “global substitution.” It is useful for automating the editing of text, including deleting characters from text or extracting only parts of text.

### 3.1 Basic find and replace examples

Let’s remove all numbers from our `dataText` list.

```
gsub("[0-9]","",dataText)
```

```
[1] "Locust Walk"      "Privet Drive"      "Malibu Point"      "B Baker St."  
[5] "-"               ""                  ""                  ""  
[9] "()"              "--"                ""                  "Kendall Roy"  
[13] "Roman Roy"      "Siobhan Roy"      "Roy Wood, Jr."    "Royal Caribbean"  
[17] "Casino Royale"  "two children"    "children"        "twins"  
[21] "Philadelphia"   "Philly"           "Phila"            "Dr. Phil"  
[25] "$"               "$"                 "$,"
```

or remove all the Roy last names

```
gsub(" Roy$","",dataText)
```

```
[1] "3718 Locust Walk"  "4 Privet Drive"    "10880 Malibu Point"  
[4] "221B Baker St."   "19104-6286"       "20015"  
[7] "90291"             "90210"            "(215) 573-9097"  
[10] "215-573-9097"    "2155739097"     "Kendall"  
[13] "Roman"             "Siobhan"          "Roy Wood, Jr."  
[16] "Royal Caribbean"  "Casino Royale"   "two children"  
[19] "2 children"        "twins"            "Philadelphia"  
[22] "Philly"            "Phila"            "Dr. Phil"  
[25] "$23456"           "$10000"          "$60,000"
```

or turn all punctuation to "X"

```
gsub("[.)(-$]","X", dataText)
```

```
[1] "3718 Locust Walk"  "4 Privet Drive"    "10880 Malibu Point"  
[4] "221B Baker StX"   "19104X6286"      "20015"  
[7] "90291"             "90210"            "X215X 573X9097"  
[10] "215X573X9097"    "2155739097"     "Kendall Roy"  
[13] "Roman Roy"        "Siobhan Roy"      "Roy Wood, JrX"  
[16] "Royal Caribbean"  "Casino Royale"   "two children"  
[19] "2 children"        "twins"            "Philadelphia"  
[22] "Philly"            "Phila"            "DrX Phil"  
[25] "X23456"           "X10000"          "X60,000"
```

or remove HTML tags from text (very very handy!!!).

```
gsub("<[^>]+>", "",  
  "<td><b><a href=\"/movie/Dont-Breathe-2-(2021)#tab=box-office\">Don't Breathe 2</a></b>")
```

[1] "Don't Breathe 2"

This last one we will use a lot. The regular expression says start with a <, then match several characters that are non >, followed by a >. Anything matching this pattern, like <td>, will get removed.

The first argument of the `gsub()` function is the “find pattern,” what we are asking `gsub()` to find. The second argument is a “replacement pattern” that will replace whatever `gsub()` matched with the find pattern. Lastly, like `grep()`, we need to tell `gsub()` the name of the data object in which it should look for the pattern.

With `gsub()` we will be using a lot of the same grammar, such as carets, brackets, and dashes, that we used with `grep()`. To illustrate, let’s remove anything that is not a number from text.

```
gsub("[^0-9]", "", dataText)
```

[1] "3718" "4" "10880" "221" "191046286"  
[6] "20015" "90291" "90210" "2155739097" "2155739097"  
[11] "2155739097" "" "" "" ""  
[16] "" "" "" "2" ""  
[21] "" "" "" "" "23456"  
[26] "10000" "60000"

### 3.1.1 Exercises

13. Remove the all the commas. Note that we probably just want to remove them from the numbers, but we do not quite yet have the tools to avoid removing the comma from Roy Wood, Jr.

## 3.2 Back-Referencing

A bit more complicated aspect of `gsub()` is “backreferencing.” Any part of the regular expression in the find pattern that is wrapped in parentheses gets stored in a “register.” You can have multiple pairs of parentheses and, therefore, save different parts of what gets matched in the find pattern. You can then recall what `gsub()` stored in the registers in the replacement

pattern, using `\$1` for what was stored in the first register, `\$2` for what was stored in the second register, and so on.

For example, let's delete the "plus four" part of the ZIP codes in our data.

```
gsub("( [0-9]{5})-[0-9]{4}", "\$1", dataText)
```

```
[1] "3718 Locust Walk"    "4 Privet Drive"      "10880 Malibu Point"  
[4] "221B Baker St."     "19104"                 "20015"  
[7] "90291"                "90210"                 "(215) 573-9097"  
[10] "215-573-9097"       "2155739097"        "Kendall Roy"  
[13] "Roman Roy"          "Siobhan Roy"         "Roy Wood, Jr."  
[16] "Royal Caribbean"    "Casino Royale"        "two children"  
[19] "2 children"          "twins"                  "Philadelphia"  
[22] "Philly"               "Phila"                  "Dr. Phil"  
[25] "$23456"                "$10000"                "$60,000"
```

We wrapped the parentheses around the first five digits, so only those first five digits get stored in register 1. Then the replacement pattern replaces everything that the find pattern matched (the entirety of the five digits, the hyphen, and the plus four digits) with the contents of register 1, containing just the first five digits.

Let's use `gsub()` to just keep the first two "words" in each element of `dataText`.

```
gsub("^( [^ ]+ [^ ]+) .*$", "\$1", dataText)
```

```
[1] "3718 Locust"        "4 Privet"           "10880 Malibu"      "221B Baker"  
[5] "19104-6286"        "20015"              "90291"            "90210"  
[9] "(215) 573-9097"    "215-573-9097"     "2155739097"     "Kendall Roy"  
[13] "Roman Roy"         "Siobhan Roy"        "Roy Wood,"        "Royal Caribbean"  
[17] "Casino Royale"     "two children"       "2 children"       "twins"  
[21] "Philadelphia"      "Philly"              "Phila"             "Dr. Phil"  
[25] "$23456"                "$10000"              "$60,000"
```

This regular expression says: start at the beginning of the text, find a bunch of characters that are not spaces (remember that the `+` means one or more of the previous), then find a space, then find another bunch of characters that are not spaces, followed by another space, followed by anything else until the end of the text. The `.` is a wild card that matches any one character. The `*` is like the `+` but it means *zero* or more of the previous character (the `+` matches *one* or more of the previous character). Note how we have used the parentheses. They are wrapped around those first two words. Those words get stored in register 1 and the replacement pattern just recalls whatever got stored in register 1.

An alternative strategy is to use `\w`, which means a “word character” (any numbers of letters) and `\s`, which means any “whitespace character” (spaces, tabs, line feeds, non-breaking space).

```
gsub("^(\\w+\\s\\w+)\\s.*$","\\1",dataText)
```

```
[1] "3718 Locust"      "4 Privet"        "10880 Malibu"    "221B Baker"  
[5] "19104-6286"      "20015"          "90291"          "90210"  
[9] "(215) 573-9097"  "215-573-9097"   "2155739097"   "Kendall Roy"  
[13] "Roman Roy"       "Siobhan Roy"     "Roy Wood, Jr."  "Royal Caribbean"  
[17] "Casino Royale"   "two children"    "2 children"     "twins"  
[21] "Philadelphia"    "Philly"          "Phila"          "Dr. Phil"  
[25] "$23456"          "$10000"         "$60,000"
```

When working with `grep()`, we wrote a regular expression to find all the phone numbers in all the various formats. Now let’s use `gsub()` to standardize all the phone numbers to have the hyphenated format, like 123-456-7890. We use the same regular expression you figured out in the exercises when using `?`  in `grep()` to find phone numbers, but we insert pairs of parentheses to capture the phone number digits.

```
gsub("^\\(\\d{3}\\)\\s|-\\d{3}-\\d{4}$","\\1-\\2-\\3",dataText)
```

```
[1] "3718 Locust Walk"  "4 Privet Drive"   "10880 Malibu Point"  
[4] "221B Baker St."   "19104-6286"      "20015"  
[7] "90291"             "90210"          "215-573-9097"  
[10] "215-573-9097"     "215-573-9097"   "Kendall Roy"  
[13] "Roman Roy"        "Siobhan Roy"     "Roy Wood, Jr."  
[16] "Royal Caribbean"  "Casino Royale"   "two children"  
[19] "2 children"        "twins"          "Philadelphia"  
[22] "Philly"            "Phila"          "Dr. Phil"  
[25] "$23456"           "$10000"         "$60,000"
```

Note that register 2 captures whatever text matches the optional text `(\\) |-`. That’s why you don’t see `\\2` in the replacement pattern.

### 3.2.1 Exercise

14. Add commas to these numbers `c("9453", "2332", "4645", "1234")`. That is, make these numbers look like 9,453 and 2,332 and so on. Fill in the find pattern and replacement pattern in `gsub("", "", c("9453", "2332", "4645", "1234"))`

### 3.3 Lookaheads

Although not as commonly used as a backreference, a “lookahead” can be helpful to find what comes next or, more generally, to `gsub()` items that are a bit more complicated. Let’s say you wanted to check that every q is followed by a u. If it’s not, then insert the u after the q.

Maybe you would do something like this

```
gsub("q[^u]", "qu", c("quiet", "quick", "quagmire", "quixotic"))
```

```
[1] "quiet"    "quick"     "quagmire" "quixotic"
```

As you can see, it doesn’t quite do what we wanted. It problematically drops the “i” from “quixotic.” Remember that the replacement pattern will overwrite whatever matches the find pattern. This find pattern will match the “qi,” the q *and* the adjacent character that is not a u.

A “lookahead” just peeks at the next character to see what it is, but does not consider it part of the match. In parentheses we signal a lookahead with ? and then to ask for a character that is not a “u” we use !u.

```
gsub("q(?![u])", "qu", c("quiet", "quick", "quagmire", "quixotic"), perl=TRUE)
```

```
[1] "quiet"    "quick"     "quagmire" "quixotic"
```

Now we have “quixotic” spelled properly as `gsub()` looked ahead to check for a “u” but did not consider it part of the match to be replaced. Note that we have set `perl=TRUE`. There is not a single regular expression standard. Lookaheads (and there are lookbehinds too) are not part of the POSIX 1003.2 standard that R uses by default. However, you can ask R to use perl-style regular expressions that do support lookaheads (and lookbehinds) by simply setting `perl=TRUE`.

Let’s say we want to add commas to large numbers to display like 1,234,567 instead of 1234567. Ready for this one?

```
gsub("(?<!\\.)([0-9])(?=([0-9]{3})+($|\\.))", "\\\1,",  
      c("$1234567890", "234", "1234", "12345", "12345.6789"),  
      perl=TRUE)
```

```
[1] "$1,234,567,890" "234"           "1,234"          "12,345"  
[5] "12,345.6789"
```

First note the replacement text, `\1`. This is going to take whatever is in register 1 and produce it again followed by a `,`. The challenge for the matching part is to write a regular expression that identifies those digits where a comma is needed to its right. There are three parts to the regular expression.

- `(?<!\\.):` This is a “lookbehind” that first makes sure there are no decimals to the left of the digit
- `([0-9]):` Match a single digit and store it in register 1
- `(?=([0-9]{3})+($|\\.)):` This is a “lookahead” checking that there are blocks of 3 digits to the right until either the end or a decimal

This gets you a handy regular expression to produce nicer looking numbers. Just so you know, R also provides some handy tools for formatting numbers.

```
prettyNum(c(1234567890, 234, 1234, 12345, 12345.6789), big.mark=",")
```

```
[1] "1,234,567,890" "234"           "1,234"           "12,345"  
[5] "12,345.68"
```

Here's how lookaheads are going to be very important for us in working with data. We often get datasets in comma-separated value format, typically with a “.csv” extension on the file name. The R function `read.csv()` can read in data in this format. Problems can occur when values in the dataset have commas inside of them.

Here's some example text that could be problematic.

```
text.w.quotes <- 'white,male,"loitering,narcotics",19,"knife,gun"'  
cat(text.w.quotes, "\n")
```

```
white,male,"loitering,narcotics",19,"knife,gun"
```

Some of the commas in this text are separating values for race, sex, arrest charge, age, and recovered contraband. However, there are other commas inside the quotes listing the arrest charges and the contraband items. Fortunately, `read.csv()` is smart enough to keep quoted items together as one data element, as long as the parameter `quote = "\"\""`, which is the default. Other functions are not so kind. Later we will use `dbWriteTable()` to build an SQL database and it will think that all the commas separate the values. So it will think there is a separate `"loitering` data element and then a `narcotics`. And the same for `"knife` and `gun"`.

So here is a very handy regular expression using lookaheads that changes commas that are not inside quotes to semicolons. You can also choose a stranger symbol, like `|` or `@`. This regular

expression looks for a , and then it uses a lookahead for the remainder of the line. It will match that comma if there are no quotes for the rest of the line, [^\"], or if there are only pairs of quotes each with non-quote characters inside of them, \"[^\"]\*\".

```
gsub(",(?=([^"]|\"[^"]*\")*$)", ";", text.w.quotes, perl=TRUE) |>  
cat()
```

```
white;male;"loitering,narcotics";19;"knife,gun"
```

As you can see, commas inside the quotes are preserved and those outside have been transformed into ;. Now we would be able to tell a function like `dbWriteTable()` that the data elements are separated by ; and it will read in the data properly.

### 3.3.1 Exercises

Make the following changes to `dataText`.

15. Change Dr. Phil to Dr. Phillip
16. Spell out Philadelphia
17. Keep just the first word or first number
18. Keep only area codes of phone numbers

## 4 Introduction to Webscraping: A Practical Application of Regular Expressions

As you have already seen, regular expressions are an extremely valuable tool when working with data. In fact, we are going to learn about webscraping next. Webscraping enables us to extract data from a website by searching the underlying HTML code for the website and extracting the desired data. To do webscraping, you are going to rely heavily on regular expressions.

### 4.1 The 10,000 most common English words

Suppose we want a list of the 10,000 most common words in English. [wiktionary.org](http://wiktionary.org) provides this list. Have a look at the webpage.

The `scan()` function in R can read data from a text file, but if given a URL, it will download the HTML code for that page.

```
words <- scan("https://en.wiktionary.org/wiki/Wiktionary:Frequency_lists/PG/2006/04/1-10000"
               what="", sep="\n")
```

If instead you got a message like “Read 50706 items” then you are in luck. You just used R to scrape a webpage. You might get a number different from 50706 and it might change if you run the script another day. Websites make updates to headers, footers, and banners and these alterations change the number of lines of HTML code required to generate the page.

The first several lines of HTML code in `words` all contain HTML code to set up the page.

```
words[1:5]
```

```
[1] "<!DOCTYPE html>"
[2] "<html class=\"client-nojs vector-feature-language-in-header-enabled vector-
feature-language-in-main-page-header-disabled vector-feature-page-tools-
pinned-disabled vector-feature-toc-pinned-clientpref-1 vector-feature-main-
menu-pinned-disabled vector-feature-limited-width-clientpref-1 vector-feature-
limited-width-content-enabled vector-feature-custom-font-size-clientpref-
1 vector-feature-appearance-pinned-clientpref-1 skin-theme-clientpref-day vector-
sticky-header-enabled vector-toc-available\" lang=\"en\" dir=\"ltr\">"
[3] "<head>"
[4] "<meta charset=\"UTF-8\">"
[5] "<title>Wiktionary:Frequency lists/PG/2006/04/1-10000 - Wiktionary, the free dictionary</title>
```

There is nothing of interest to us here. We just want the part with the 10,000 most common words. Let’s look further down the page.

```
words[490+0:9]
```

```
[1] "<td><a href=\"/wiki/and\" title=\"and\">and</a></td>"
[2] "<td>29944184"
[3] "</td></tr>"
[4] "<tr>"
[5] "<td>4</td>"
[6] "<td><a href=\"/wiki/to\" title=\"to\">to</a></td>"
[7] "<td>25956096"
[8] "</td></tr>"
[9] "<tr>"
[10] "<td>5</td>"
```

Here we start finding some words! Note that every line that has one of the words we are looking for has `title=""`. We can use that phrase to find lines that have the 10,000 words on them. Use `grep()` to find those lines and print out the first 10 of them to see if this will work for us.

```
words <- grep("title=\"\"", words, value=TRUE)
words[23:30]
```

```
[1] "\t\t\t\t\t<div id=\"mw-content-text\" class=\"mw-body-content\"><div class=\"mw-content-ltr mw-parser-output\" lang=\"en\" dir=\"ltr\"><div class=\"mw-heading mw-heading3\"><h3 id=\"1_-_10000\" data-mw-thread-id=\"h-1_-_10000\"><span data-mw-comment-start=\"\" id=\"h-1_-_10000\"></span>1 - 10000<span data-mw-comment-end=\"h-1_-_10000\"></span></h3><span class=\"mw-editsection\"><span class=\"mw-editsection-bracket\">[</span><a href=\"/w/index.php?title=Wiktionary:Frequency_lists/PG/2000&action=edit&section=1\" title=\"Edit section: 1 - 10000\"><span>edit</span></a></span>-<span class=\"mw-editsection-bracket\">]</span></span></div>"
[2] "<td><a href=\"/wiki/the\" title=\"the\">the</a></td>"
[3] "<td><a href=\"/wiki/of\" title=\"of\">of</a></td>"
[4] "<td><a href=\"/wiki/and\" title=\"and\">and</a></td>"
[5] "<td><a href=\"/wiki/to\" title=\"to\">to</a></td>"
[6] "<td><a href=\"/wiki/in\" title=\"in\">in</a></td>"
[7] "<td><a href=\"/wiki/i\" title=\"i\">i</a></td>"
[8] "<td><a href=\"/wiki/that\" title=\"that\">that</a></td>"
```

While the first 24 lines do have the phrase `title=""`, they are not the ones with the 10,000 words. But starting at line 25 we start seeing the most common words: the, of, and, to. So let's cut the first 24 lines from `words` and keep the next 10,000 lines. `title=""` shows up more in the HTML code in the footer after the 10,000th word.

```
words <- words[-(1:24)]
words <- words[1:10000]
head(words) # look at the first several rows
```

```
[1] "<td><a href=\"/wiki/of\" title=\"of\">of</a></td>"
[2] "<td><a href=\"/wiki/and\" title=\"and\">and</a></td>"
[3] "<td><a href=\"/wiki/to\" title=\"to\">to</a></td>"
[4] "<td><a href=\"/wiki/in\" title=\"in\">in</a></td>"
[5] "<td><a href=\"/wiki/i\" title=\"i\">i</a></td>"
[6] "<td><a href=\"/wiki/that\" title=\"that\">that</a></td>"
```

```
tail(words) # look at the last several rows
```

```
[1] "<td><a href=\"/wiki/repressed\" title=\"repressed\">repressed</a></td>"  
[2] "<td><a href=\"/wiki/cooperation\" title=\"cooperation\">cooperation</a></td>"  
[3] "<td><a href=\"/wiki/sequel\" title=\"sequel\">sequel</a></td>"  
[4] "<td><a href=\"/wiki/wench\" title=\"wench\">wench</a></td>"  
[5] "<td><a href=\"/wiki/calves\" title=\"calves\">calves</a></td>"  
[6] NA
```

These lines have the text we need, but there are still a lot of HTML tags, all that HTML code wrapped in < >. <td> is the HTML tag marking a cell in a table and <a> is the HTML tag for creating hyperlinks to other pages. Each of these also has ending tags </td> and </a>. Fortunately for us they are all contained within the < and > characters. So let's use `gsub()` to remove all the HTML tags.

```
words <- gsub("<[^>]*>", "", words)
```

This regular expression looks for a <, then a bunch of characters that are not a >, followed by a > and replaces them with nothing. Now `words` contains just our list of 10,000 most common words. Here are the first 50 of them.

```
words[1:50]
```

```
[1] "of"      "and"     "to"       "in"       "i"        "that"     "was"      "his"      "he"  
[10] "it"      "with"    "is"       "for"      "as"       "had"      "you"      "not"      "be"  
[19] "her"     "on"      "at"       "by"       "which"    "have"     "or"       "from"    "this"  
[28] "him"     "but"     "all"      "she"      "they"     "were"     "my"       "are"      "me"  
[37] "one"     "their"   "so"       "an"       "said"     "them"     "we"       "who"      "would"  
[46] "been"    "will"    "no"       "when"    "there"
```

#### 4.1.1 Exercises

19. We were told “i before e except after c, or when sounded like a as in neighbor or weigh”. Is that true?
20. Find words with punctuation
21. Find words that do not have aeiou in the first four characters

## 5 Solutions to the exercises

1. Find text with vowels

```
grep("[AEIOUaeiou]", dataText, value=TRUE)
```

```
[1] "3718 Locust Walk"      "4 Privet Drive"      "10880 Malibu Point"  
[4] "221B Baker St."       "Kendall Roy"        "Roman Roy"  
[7] "Siobhan Roy"         "Roy Wood, Jr."       "Royal Caribbean"  
[10] "Casino Royale"       "two children"       "2 children"  
[13] "twins"                "Philadelphia"     "Philly"  
[16] "Phila"                 "Dr. Phil"
```

2. Find text with a number immediately followed by a letter

```
grep("[0-9][A-Za-z]", dataText, value=TRUE)
```

```
[1] "221B Baker St."
```

3. What is the difference between `grep("Roy", dataText, value=TRUE)` and `grep("[Roy]", dataText, value=TRUE)`? Roy will only match text that is exactly “Roy”, while [Roy] will match text that has any one of the characters “R”, “o”, or “y”.

```
grep("Roy", dataText, value=TRUE)
```

```
[1] "Kendall Roy"          "Roman Roy"          "Siobhan Roy"          "Roy Wood, Jr."  
[5] "Royal Caribbean"     "Casino Royale"
```

```
grep("[Roy]", dataText, value=TRUE)
```

```
[1] "3718 Locust Walk"      "10880 Malibu Point" "Kendall Roy"  
[4] "Roman Roy"            "Siobhan Roy"        "Roy Wood, Jr."  
[7] "Royal Caribbean"     "Casino Royale"      "two children"  
[10] "Philly"
```

4. Find the Roy children

```
grep("Roy$", dataText, value = TRUE)
```

```
[1] "Kendall Roy" "Roman Roy" "Siobhan Roy"
```

5. Find text that starts with a number and ends with a letter

```
grep("^ [0-9] [0-z ]+[A-Za-z] $", dataText, value=TRUE)
```

```
[1] "3718 Locust Walk" "4 Privet Drive" "10880 Malibu Point"  
[4] "2 children"
```

6. Find commas

```
grep(",", dataText, value=TRUE)
```

```
[1] "Roy Wood, Jr." "$60,000"
```

Remember that not every regular expression is complicated. Sometimes you just need to search for something specific and it requires no fancy regular expression.

7. Find ZIP codes

```
grep("^ [0-9]{5}$|^ [0-9]{5}-", dataText, value=TRUE)
```

```
[1] "19104-6286" "20015" "90291" "90210"
```

8. Find those with six letter word

```
grep("\\b[A-Za-z]{6}\\b", dataText, value=TRUE)
```

```
[1] "3718 Locust Walk" "4 Privet Drive" "10880 Malibu Point"  
[4] "Casino Royale" "Philly"
```

9. Find mentions of Philadelphia

```
grep("Phil(adelphia|ly|a)", dataText, value=TRUE)
```

```
[1] "Philadelphia" "Philly" "Phila"
```

or

```
# "Phil" followed by something that is not whitespace
grep("Phil[^\\s]", dataText, value=TRUE)
```

```
[1] "Philadelphia" "Philly" "Phila"
```

10. Find capitalized words

```
grep("\\b[A-Z]", dataText, value=TRUE)
```

```
[1] "3718 Locust Walk" "4 Privet Drive" "10880 Malibu Point"
[4] "221B Baker St." "Kendall Roy" "Roman Roy"
[7] "Siobhan Roy" "Roy Wood, Jr." "Royal Caribbean"
[10] "Casino Royale" "Philadelphia" "Philly"
[13] "Phila" "Dr. Phil"
```

11. Find the addresses

```
grep("[0-9]+[A-Z]? [A-Za-z ]+ (Drive|Walk|Point|St)", dataText, value=TRUE)
```

```
[1] "3718 Locust Walk" "4 Privet Drive" "10880 Malibu Point"
[4] "221B Baker St."
```

12. Find a shorter way to get phone numbers using ?

```
grep("\\(?[0-9]{3}\\) |-| ?[0-9]{3}(-| )?[0-9]{4}", dataText, value=TRUE)
```

```
[1] "(215) 573-9097" "215-573-9097" "2155739097"
```

13. Remove the all the commas

```
gsub(", ", "", dataText)
```

```
[1] "3718 Locust Walk" "4 Privet Drive" "10880 Malibu Point"
[4] "221B Baker St." "19104-6286" "20015"
[7] "90291" "90210" "(215) 573-9097"
[10] "215-573-9097" "2155739097" "Kendall Roy"
[13] "Roman Roy" "Siobhan Roy" "Roy Wood Jr."
[16] "Royal Caribbean" "Casino Royale" "two children"
[19] "2 children" "twins" "Philadelphia"
[22] "Philly" "Phila" "Dr. Phil"
[25] "$23456" "$10000" "$60000"
```

14. Add commas to these numbers `c("9453", "2332", "4645", "1234")`.

```
gsub("(\\d)(\\d\\d\\d)", "\\\\$1,\\\$2", c("9453", "2332", "4645", "1234"))
```

```
[1] "9,453" "2,332" "4,645" "1,234"
```

or

```
gsub("^([0-9])", "\\\\$1,", c("9453", "2332", "4645", "1234"))
```

```
[1] "9,453" "2,332" "4,645" "1,234"
```

15. Change Dr. Phil to Dr. Phillip

```
gsub("Phil$", "Phillip", dataText)
```

```
[1] "3718 Locust Walk"      "4 Privet Drive"      "10880 Malibu Point"  
[4] "221B Baker St."       "19104-6286"        "20015"  
[7] "90291"                 "90210"           "(215) 573-9097"  
[10] "215-573-9097"         "2155739097"      "Kendall Roy"  
[13] "Roman Roy"            "Siobhan Roy"       "Roy Wood, Jr."  
[16] "Royal Caribbean"     "Casino Royale"     "two children"  
[19] "2 children"            "twins"            "Philadelphia"  
[22] "Philly"                "Phila"             "Dr. Phillip"  
[25] "$23456"                "$10000"           "$60,000"
```

16. Spell out Philadelphia

```
gsub(("Phil(adelphia|ly|a)", "Philadelphia", dataText)
```

```
[1] "3718 Locust Walk"      "4 Privet Drive"      "10880 Malibu Point"  
[4] "221B Baker St."       "19104-6286"        "20015"  
[7] "90291"                 "90210"           "(215) 573-9097"  
[10] "215-573-9097"         "2155739097"      "Kendall Roy"  
[13] "Roman Roy"            "Siobhan Roy"       "Roy Wood, Jr."  
[16] "Royal Caribbean"     "Casino Royale"     "two children"  
[19] "2 children"            "twins"            "Philadelphia"  
[22] "Philadelphia"         "Philadelphia"      "Dr. Phil"  
[25] "$23456"                "$10000"           "$60,000"
```

or

```
gsub("Phil[^\b]+", "Philadelphia", dataText)
```

```
[1] "3718 Locust Walk"      "4 Privet Drive"      "10880 Malibu Point"  
[4] "221B Baker St."       "19104-6286"        "20015"  
[7] "90291"                 "90210"             "(215) 573-9097"  
[10] "215-573-9097"        "2155739097"      "Kendall Roy"  
[13] "Roman Roy"           "Siobhan Roy"       "Roy Wood, Jr."  
[16] "Royal Caribbean"     "Casino Royale"    "two children"  
[19] "2 children"           "twins"            "Philadelphia"  
[22] "Philadelphia"         "Philadelphia"     "Dr. Phil"  
[25] "$23456"                "$10000"           "$60,000"
```

17. Keep just the first word or first number

```
gsub("^( [^ ]+) .*$", "\\\1", dataText)
```

```
[1] "3718"          "4"           "10880"        "221B"        "19104-6286"  
[6] "20015"         "90291"        "90210"        "(215)"      "215-573-  
9097"  
[11] "2155739097"   "Kendall"      "Roman"        "Siobhan"    "Roy"  
[16] "Royal"         "Casino"       "two"          "2"          "twins"  
[21] "Philadelphia"  "Philly"       "Phila"        "Dr."        "$23456"  
[26] "$10000"        "$60,000"
```

or

```
gsub(" .*", "", dataText)
```

```
[1] "3718"          "4"           "10880"        "221B"        "19104-6286"  
[6] "20015"         "90291"        "90210"        "(215)"      "215-573-  
9097"  
[11] "2155739097"   "Kendall"      "Roman"        "Siobhan"    "Roy"  
[16] "Royal"         "Casino"       "two"          "2"          "twins"  
[21] "Philadelphia"  "Philly"       "Phila"        "Dr."        "$23456"  
[26] "$10000"        "$60,000"
```

18. Keep only area codes of phone numbers

```
gsub("^\\(?(\\d{3})\\)\\ |-| )?\\d{3}-?\\d{4}", "\\1", dataText)
```

```
[1] "3718 Locust Walk"      "4 Privet Drive"      "10880 Malibu Point"  
[4] "221B Baker St."       "19104-6286"        "20015"  
[7] "90291"                "90210"            "215"  
[10] "215"                  "215"              "Kendall Roy"  
[13] "Roman Roy"           "Siobhan Roy"       "Roy Wood, Jr."  
[16] "Royal Caribbean"     "Casino Royale"    "two children"  
[19] "2 children"           "twins"            "Philadelphia"  
[22] "Philly"               "Phila"             "Dr. Phil"  
[25] "$23456"               "$10000"           "$60,000"
```

19. We were told “i before e except after c, or when sounded like a as in neighbor or weigh”.  
Is that true?

```
grep("[^c]ei", words, value=TRUE)
```

```
[1] "their"          "being"          "neither"         "seeing"         "foreign"  
[6] "weight"          "seized"          "height"          "reign"          "sovereign"  
[11] "beings"          "wherein"         "veil"            "therein"        "leisure"  
[16] "seize"           "neighborhood"    "heir"            "theirs"          "neighbors"  
[21] "veins"           "heights"          "neighbour"       "neighbouring"   "weighed"  
[26] "neighbor"        "reigned"          "sovereignty"     "foreigners"     "neighboring"  
[31] "reins"            "seizing"          "weigh"            "vein"           "Keith"  
[36] "sovereigns"       "leisurely"        "foreigner"       "weird"          "deity"  
[41] "weighing"         "freight"          "deities"          "rein"           "reigns"  
[46] "weighty"          "Raleigh"          "albeit"          "heightened"     "Seine"  
[51] "well-being"        "forfeit"          "Marseilles"      "twenty-eight"   "herein"
```

```
grep("cie", words, value=TRUE)
```

```
[1] "ancient"          "society"          "sufficient"       "species"  
[5] "science"           "conscience"        "sufficiently"     "scientific"  
[9] "fancied"           "fancies"          "efficient"        "efficiency"  
[13] "tendencies"        "conscientious"    "insufficient"    "deficient"  
[17] "deficiency"
```

There are a lot of words with “ei” where the letter before the “ei” ” is not a “c”. Also, there are a lot of words that have “ie” immediately following a “c”.

20. Find words with punctuation

```
grep("[^a-zA-Z0-9]", words, value=TRUE)
```

```
[1] "I'll"          "can't"        "I've"         "didn't"  
[5] "an'"           "won't"        "man's"        "c."  
[9] "that's"        "etc."         "o'"          "I'd"  
[13] "v."            "father's"      "wouldn't"     "he's"  
[17] "couldn't"      "isn't"        "'em"         "king's"  
[21] "there's"       "ain't"        "you'll"       "god's"  
[25] "mother+'s"    "wasn't"       "doesn't"     "haven't"  
[29] "you've"        "woman's"      "we'll"        "she's"  
[33] "'tis"          "you'd"        "what's"       "cf."  
[37] "hadn't"        "th'"          "pp."          "rest."  
[41] "o'er"          "u.s."         "he'd"         "they're"  
[45] "we're"          "day's"        "girl's"       "he'll"  
[49] "shouldn't"     "husband's"    "twenty-five"  "men's"  
[53] "other's"        "i.e."         "brother's"    "wife's"  
[57] "we've"          "well-known"   "lady's"       "twenty-four"  
[61] "ma'am"          "enemy's"      "imp."         "hasn't"  
[65] "majesty's"      "viz."         "mustn't"      "old-fashioned"  
[69] "friend's"       "so-called"    "re-use"       "aren't"  
[73] "they'll"        "needn't"      "pub."         "middle-aged"  
[77] "she'd"          "t'"          "of."          "oe."  
[81] "prof."          "she'll"       "shan't"       "ne'er"  
[85] "as."            "they'd"       "we'd"         "here's"  
[89] "they've"        "adj."         "weren't"      "sq."  
[93] "gen."           "cong."        "col."         "brother-in-law"  
[97] "son-in-law"      "e.g."         "ch."          "e-mail"  
[101] "twenty-one"     "twenty-two"    "it'll"        "ft."  
[105] "self-control"   "anglo-saxon"  "ff."          "forty-five"  
[109] "'ll"            "Icel."        "good-looking" "p.m."  
[113] "e'er"           "twenty-three" "thirty-six"   "a.m."  
[117] "father-in-law"  "e'en"        "seventy-five" "'cause"  
[121] "well-being"     "Mass."        "thirty-two"   "self-"  
[125] "twenty-eight"   "twenty-six"   "sister-in-law" "three-quarters"  
[129] "'n'"
```

or

```
grep("[[:punct:]]", words, value=TRUE)
```

```
[1] "I'll"          "can't"         "I've"          "didn't"  
[5] "an'"           "won't"          "man's"          "c."  
[9] "that's"        "etc."           "o'"             "I'd"  
[13] "v."            "father's"       "wouldn't"      "he's"  
[17] "couldn't"      "isn't"          "'em"           "king's"  
[21] "there's"       "ain't"          "you'll"        "god's"  
[25] "mother+'s"     "wasn't"         "doesn't"       "haven't"  
[29] "you've"        "woman's"        "we'll"          "she's"  
[33] "'tis"          "you'd"          "what's"         "cf."  
[37] "hadn't"        "th'"            "pp."            "rest."  
[41] "o'er"          "u.s."           "he'd"           "they're"  
[45] "we're"          "day's"          "girl's"         "he'll"  
[49] "shouldn't"     "husband's"      "twenty-five"    "men's"  
[53] "other's"        "i.e."           "brother's"     "wife's"  
[57] "we've"          "well-known"     "lady's"         "twenty-four"  
[61] "ma'am"          "enemy's"         "imp."           "hasn't"  
[65] "majesty's"      "viz."           "mustn't"        "old-fashioned"  
[69] "friend's"       "so-called"      "re-use"         "aren't"  
[73] "they'll"        "needn't"        "pub."           "middle-aged"  
[77] "she'd"          "t'"              "of."             "oe."  
[81] "prof."          "she'll"          "shan't"         "ne'er"  
[85] "as."             "they'd"          "we'd"           "here's"  
[89] "they've"        "adj."            "weren't"        "sq."  
[93] "gen."           "cong."          "col."           "brother-in-law"  
[97] "son-in-law"      "e.g."           "ch."            "e-mail"  
[101] "twenty-one"     "twenty-two"      "it'll"          "ft."  
[105] "self-control"   "anglo-saxon"    "ff."            "forty-five"  
[109] "'ll"            "Icel."          "good-looking"   "p.m."  
[113] "e'er"           "twenty-three"   "thirty-six"     "a.m."  
[117] "father-in-law"  "e'en"           "seventy-five"   "'cause"  
[121] "well-being"     "Mass."          "thirty-two"     "self-"  
[125] "twenty-eight"   "twenty-six"     "sister-in-law"  "three-quarters"  
[129] "'n'"
```

[**:punct:**] is a special set containing all punctuation characters. or

```
grep("\W", words, value=TRUE)
```

```
[1] "I'll"          "can't"         "I've"          "didn't"
```

[5]	"an'"	"won't"	"man's"	"c."
[9]	"that's"	"etc."	"o'"	"I'd"
[13]	"v."	"father's"	"wouldn't"	"he's"
[17]	"couldn't"	"isn't"	"'em"	"king's"
[21]	"there's"	"ain't"	"you'll"	"god's"
[25]	"mother+'s"	"wasn't"	"doesn't"	"haven't"
[29]	"you've"	"woman's"	"we'll"	"she's"
[33]	"'tis"	"you'd"	"what's"	"cf."
[37]	"hadn't"	"th'"	"pp."	"rest."
[41]	"o'er"	"u.s."	"he'd"	"they're"
[45]	"we're"	"day's"	"girl's"	"he'll"
[49]	"shouldn't"	"husband's"	"twenty-five"	"men's"
[53]	"other's"	"i.e."	"brother's"	"wife's"
[57]	"we've"	"well-known"	"lady's"	"twenty-four"
[61]	"ma'am"	"enemy's"	"imp."	"hasn't"
[65]	"majesty's"	"viz."	"mustn't"	"old-fashioned"
[69]	"friend's"	"so-called"	"re-use"	"aren't"
[73]	"they'll"	"needn't"	"pub."	"middle-aged"
[77]	"she'd"	"t'"	"of."	"oe."
[81]	"prof."	"she'll"	"shan't"	"ne'er"
[85]	"as."	"they'd"	"we'd"	"here's"
[89]	"they've"	"adj."	"weren't"	"sq."
[93]	"gen."	"cong."	"col."	"brother-in-law"
[97]	"son-in-law"	"e.g."	"ch."	"e-mail"
[101]	"twenty-one"	"twenty-two"	"it'll"	"ft."
[105]	"self-control"	"anglo-saxon"	"ff."	"forty-five"
[109]	"ll"	"Icel."	"good-looking"	"p.m."
[113]	"e'er"	"twenty-three"	"thirty-six"	"a.m."
[117]	"father-in-law"	"e'en"	"seventy-five"	"'cause"
[121]	"well-being"	"Mass."	"thirty-two"	"self-"
[125]	"twenty-eight"	"twenty-six"	"sister-in-law"	"three-quarters"
[129]	"'n'"			

\W matches any character that is not a number or a letter.

21. Find words that do not have aeiou in the first four characters

```
grep("^[^aeiouAEIOU]{4}", words, value=TRUE)
```

[1]	"system"	"physical"	"style"	"sympathy"
[5]	"mysterious"	"mystery"	"physician"	"thyself"
[9]	"sympathetic"	"mysteries"	"systems"	"Dryden"

[13]	"symptoms"	"symbol"	"crystal"	"10th"
[17]	"hymn"	"mystic"	"sympathies"	"Sydney"
[21]	"12th"	"syllable"	"11th"	"rhyme"
[25]	"symbols"	"physically"	"rhythm"	"systematic"
[29]	"psychology"	"psychological"	"p.m."	"mystical"
[33]	"mythology"	"myth"	"syllables"	"hysterical"