

Lesson 13 Intro to Text Mining

Lusine Zilfимian

May 11 (Monday), 2020

Contents

- Intro Glossary
- String Manipulation
- Stringr vs Base
- RegEx
- Intro to TM

Bibliography

- See on Github: Handling and Processing Strings
- Character Strings - Chapter 2
- String Manipulation + stringr - Chapter 3, 4
- RegEx - Chapter 5

- A **token** is a meaningful unit of text, such as a word, that we are interested in using for analysis, and **tokenization** is the process of splitting text into tokens.
- **Semantics** means the meaning and interpretation of words, signs, and sentence structure.
- In **semantic analysis** we will care about the place and order of the words

- String manipulation - series of functions used to extract information from text variables.
- **Regex** is set of commands used to match a family of text (alphanumeric, digits, words) to detect string sequences in a large text data.
- *Regular expressions* to do more complicated tasks such as extract email IDs or date from a set of text.
- String manipulation functions are not customized however we can customize regular expressions in any way we want.

String Manipulation

- We can create strings with either single quotes or double quotes.
- Unlike other languages, there is no difference in behavior.
- To include a literal single or double quote in a string you can use `\` to escape it.

```
(x <- "This is my first character!") # double quotes
```

```
## [1] "This is my first character!"
```

```
class(x) == typeof(x)
```

```
## [1] TRUE
```

```
(x <- 'This is my second string') # single quotes
```

```
## [1] "This is my second string"
```

```
(x <- "Why do we need the 'python'?")
```

```
## [1] "Why do we need the 'python'?"
```

```
(x <- 'Why do we need the "python"?')
```

Working with strings

- If you try to include single quotes within single quotes, or double quotes within double quotes, it will not work

Something went wrong

```
"Why call langauge "python" ? "  
'Why call langauge 'python' ? '
```

```
## Error: <text>:2:21: unexpected symbol  
## 1: # Something went wrong  
## 2: "Why call langauge "python  
##      ^
```

- \ escaping character

```
(x <- 'Why do we need the \'python\'?')
```

```
## [1] "Why do we need the 'python'?"
```

```
is.character(x)
```

```
## [1] TRUE
```

- Paste() and its arguments
- Takes or more R objects, convert them to character and then concatenates them to character (one or more)

```
paste("Life of", pi)
```

```
## [1] "Life of 3.14159265358979"
```

- Element wise concatenation of two vectors

```
x <- 1:26  
(a <- LETTERS[x[1:13]])
```

```
## [1] "A" "B" "C" "D" "E" "F" "G" "H" "I" "J" "K" "L" "M"
```

```
(b <- letters[x[14:26]])
```

```
## [1] "n" "o" "p" "q" "r" "s" "t" "u" "v" "w" "x" "y" "z"
```



```
(char1 <- paste(x, c(a, b)))
```

```
## [1] "1 A" "2 B" "3 C" "4 D" "5 E" "6 F" "7 G" "8 H" "9 I" "10  
## [11] "11 K" "12 L" "13 M" "14 n" "15 o" "16 p" "17 q" "18 r" "19 s" "20  
## [21] "21 u" "22 v" "23 w" "24 x" "25 y" "26 z"
```

```
length(char1)
```

```
## [1] 26
```

- sep = The element which separates every term.
- collapse = The element which separates every result.

```
paste(x, c(a, b), sep = " - is - ") # optional - Define your own separator
```

```
## [1] "1 - is - A" "2 - is - B" "3 - is - C" "4 - is - D" "5 - is - E  
## [6] "6 - is - F" "7 - is - G" "8 - is - H" "9 - is - I" "10 - is -  
## [11] "11 - is - K" "12 - is - L" "13 - is - M" "14 - is - n" "15 - is -  
## [16] "16 - is - p" "17 - is - q" "18 - is - r" "19 - is - s" "20 - is -  
## [21] "21 - is - u" "22 - is - v" "23 - is - w" "24 - is - x" "25 - is -  
## [26] "26 - is - z"
```

```
(char2 <- paste(x, c(a, b), sep = " - is - ", collapse = " , ")) # optional
```

Basic Sting manipulations (Base)

- `nchar()` counts the number of characters (including whitespaces)
- for single character

```
nchar("Here are 16 chars: ")
```

```
## [1] 19
```

```
char1
```

```
## [1] "1 A" "2 B" "3 C" "4 D" "5 E" "6 F" "7 G" "8 H" "9 I" "10  
## [11] "11 K" "12 L" "13 M" "14 n" "15 o" "16 p" "17 q" "18 r" "19 s" "20  
## [21] "21 u" "22 v" "23 w" "24 x" "25 y" "26 z"
```

```
length(char2)
```

```
## [1] 1
```

```
nchar(char1)
```

```
## [1] 3 3 3 3 3 3 3 3 3 3 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4
```

- We can convert all characters of our string to lower case with the `tolower` command
- and to upper case with `toupper`

```
toupper("Here are 16 chars: ")
```

```
## [1] "HERE ARE 16 CHARS: "
```

```
tolower(toupper("Here are 16 chars: "))
```

```
## [1] "here are 16 chars: "
```

```
chartr(old = "is", new = "is not", char2) # the size should be equal
```

```
## [1] "1 - is - A , 2 - is - B , 3 - is - C , 4 - is - D , 5 - is - E , 6 - is - F , 7 - is - G , 8 - is - H , 9 - is - I , 10 - is - J , 11 - is - K , 12 - is - L , 13 - is - M , 14 - is - N , 15 - is - O , 16 - is - P , 17 - is - Q , 18 - is - R , 19 - is - S , 20 - is - T , 21 - is - U , 22 - is - V , 23 - is - W , 24 - is - X , 25 - is - Y , 26 - is - Z , 27 - is - [ , 28 - is - \ , 29 - is - ] , 30 - is - ^ , 31 - is - _ , 32 - is - ` , 33 - is - ~ , 34 - is - space , 35 - is - ! , 36 - is - @ , 37 - is - # , 38 - is - $ , 39 - is - % , 40 - is - & , 41 - is - * , 42 - is - ( , 43 - is - ) , 44 - is - , , 45 - is - - , 46 - is - . , 47 - is - / , 48 - is - : , 49 - is - ; , 50 - is - < , 51 - is - = , 52 - is - > , 53 - is - ? , 54 - is - [ , 55 - is - \ , 56 - is - ] , 57 - is - ^ , 58 - is - _ , 59 - is - ` , 60 - is - ~ , 61 - is - space , 62 - is - ! , 63 - is - @ , 64 - is - # , 65 - is - $ , 66 - is - % , 67 - is - & , 68 - is - * , 69 - is - ( , 70 - is - ) , 71 - is - , , 72 - is - - , 73 - is - . , 74 - is - / , 75 - is - : , 76 - is - ; , 77 - is - < , 78 - is - = , 79 - is - > , 80 - is - ? , 81 - is - [ , 82 - is - \ , 83 - is - ] , 84 - is - ^ , 85 - is - _ , 86 - is - ` , 87 - is - ~ , 88 - is - space , 89 - is - ! , 90 - is - @ , 91 - is - # , 92 - is - $ , 93 - is - % , 94 - is - & , 95 - is - * , 96 - is - ( , 97 - is - ) , 98 - is - , , 99 - is - - , 100 - is - . , 101 - is - / , 102 - is - : , 103 - is - ; , 104 - is - < , 105 - is - = , 106 - is - > , 107 - is - ? , 108 - is - [ , 109 - is - \ , 110 - is - ] , 111 - is - ^ , 112 - is - _ , 113 - is - ` , 114 - is - ~ , 115 - is - space , 116 - is - ! , 117 - is - @ , 118 - is - # , 119 - is - $ , 120 - is - % , 121 - is - & , 122 - is - * , 123 - is - ( , 124 - is - ) , 125 - is - , , 126 - is - - , 127 - is - . , 128 - is - / , 129 - is - : , 130 - is - ; , 131 - is - < , 132 - is - = , 133 - is - > , 134 - is - ? , 135 - is - [ , 136 - is - \ , 137 - is - ] , 138 - is - ^ , 139 - is - _ , 140 - is - ` , 141 - is - ~ , 142 - is - space , 143 - is - ! , 144 - is - @ , 145 - is - # , 146 - is - $ , 147 - is - % , 148 - is - & , 149 - is - * , 150 - is - ( , 151 - is - ) , 152 - is - , , 153 - is - - , 154 - is - . , 155 - is - / , 156 - is - : , 157 - is - ; , 158 - is - < , 159 - is - = , 160 - is - > , 161 - is - ? , 162 - is - [ , 163 - is - \ , 164 - is - ] , 165 - is - ^ , 166 - is - _ , 167 - is - ` , 168 - is - ~ , 169 - is - space , 170 - is - ! , 171 - is - @ , 172 - is - # , 173 - is - $ , 174 - is - % , 175 - is - & , 176 - is - * , 177 - is - ( , 178 - is - ) , 179 - is - , , 180 - is - - , 181 - is - . , 182 - is - / , 183 - is - : , 184 - is - ; , 185 - is - < , 186 - is - = , 187 - is - > , 188 - is - ? , 189 - is - [ , 190 - is - \ , 191 - is - ] , 192 - is - ^ , 193 - is - _ , 194 - is - ` , 195 - is - ~ , 196 - is - space , 197 - is - ! , 198 - is - @ , 199 - is - # , 200 - is - $ , 201 - is - % , 202 - is - & , 203 - is - * , 204 - is - ( , 205 - is - ) , 206 - is - , , 207 - is - - , 208 - is - . , 209 - is - / , 210 - is - : , 211 - is - ; , 212 - is - < , 213 - is - = , 214 - is - > , 215 - is - ? , 216 - is - [ , 217 - is - \ , 218 - is - ] , 219 - is - ^ , 220 - is - _ , 221 - is - ` , 222 - is - ~ , 223 - is - space , 224 - is - ! , 225 - is - @ , 226 - is - # , 227 - is - $ , 228 - is - % , 229 - is - & , 230 - is - * , 231 - is - ( , 232 - is - ) , 233 - is - , , 234 - is - - , 235 - is - . , 236 - is - / , 237 - is - : , 238 - is - ; , 239 - is - < , 240 - is - = , 241 - is - > , 242 - is - ? , 243 - is - [ , 244 - is - \ , 245 - is - ] , 246 - is - ^ , 247 - is - _ , 248 - is - ` , 249 - is - ~ , 250 - is - space , 251 - is - ! , 252 - is - @ , 253 - is - # , 254 - is - $ , 255 - is - % , 256 - is - & , 257 - is - * , 258 - is - ( , 259 - is - ) , 260 - is - , , 261 - is - - , 262 - is - . , 263 - is - / , 264 - is - : , 265 - is - ; , 266 - is - < , 267 - is - = , 268 - is - > , 269 - is - ? , 270 - is - [ , 271 - is - \ , 272 - is - ] , 273 - is - ^ , 274 - is - _ , 275 - is - ` , 276 - is - ~ , 277 - is - space , 278 - is - ! , 279 - is - @ , 280 - is - # , 281 - is - $ , 282 - is - % , 283 - is - & , 284 - is - * , 285 - is - ( , 286 - is - ) , 287 - is - , , 288 - is - - , 289 - is - . , 290 - is - / , 291 - is - : , 292 - is - ; , 293 - is - < , 294 - is - = , 295 - is - > , 296 - is - ? , 297 - is - [ , 298 - is - \ , 299 - is - ] , 300 - is - ^ , 301 - is - _ , 302 - is - ` , 303 - is - ~ , 304 - is - space , 305 - is - ! , 306 - is - @ , 307 - is - # , 308 - is - $ , 309 - is - % , 310 - is - & , 311 - is - * , 312 - is - ( , 313 - is - ) , 314 - is - , , 315 - is - - , 316 - is - . , 317 - is - / , 318 - is - : , 319 - is - ; , 320 - is - < , 321 - is - = , 322 - is - > , 323 - is - ? , 324 - is - [ , 325 - is - \ , 326 - is - ] , 327 - is - ^ , 328 - is - _ , 329 - is - ` , 330 - is - ~ , 331 - is - space , 332 - is - ! , 333 - is - @ , 334 - is - # , 335 - is - $ , 336 - is - % , 337 - is - & , 338 - is - * , 339 - is - ( , 340 - is - ) , 341 - is - , , 342 - is - - , 343 - is - . , 344 - is - / , 345 - is - : , 346 - is - ; , 347 - is - < , 348 - is - = , 349 - is - > , 350 - is - ? , 351 - is - [ , 352 - is - \ , 353 - is - ] , 354 - is - ^ , 355 - is - _ , 356 - is - ` , 357 - is - ~ , 358 - is - space , 359 - is - ! , 360 - is - @ , 361 - is - # , 362 - is - $ , 363 - is - % , 364 - is - & , 365 - is - * , 366 - is - ( , 367 - is - ) , 368 - is - , , 369 - is - - , 370 - is - . , 371 - is - / , 372 - is - : , 373 - is - ; , 374 - is - < , 375 - is - = , 376 - is - > , 377 - is - ? , 378 - is - [ , 379 - is - \ , 380 - is - ] , 381 - is - ^ , 382 - is - _ , 383 - is - ` , 384 - is - ~ , 385 - is - space , 386 - is - ! , 387 - is - @ , 388 - is - # , 389 - is - $ , 390 - is - % , 391 - is - & , 392 - is - * , 393 - is - ( , 394 - is - ) , 395 - is - , , 396 - is - - , 397 - is - . , 398 - is - / , 399 - is - : , 400 - is - ; , 401 - is - < , 402 - is - = , 403 - is - > , 404 - is - ? , 405 - is - [ , 406 - is - \ , 407 - is - ] , 408 - is - ^ , 409 - is - _ , 410 - is - ` , 411 - is - ~ , 412 - is - space , 413 - is - ! , 414 - is - @ , 415 - is - # , 416 - is - $ , 417 - is - % , 418 - is - & , 419 - is - * , 420 - is - ( , 421 - is - ) , 422 - is - , , 423 - is - - , 424 - is - . , 425 - is - / , 426 - is - : , 427 - is - ; , 428 - is - < , 429 - is - = , 430 - is - > , 431 - is - ? , 432 - is - [ , 433 - is - \ , 434 - is - ] , 435 - is - ^ , 436 - is - _ , 437 - is - ` , 438 - is - ~ , 439 - is - space , 440 - is - ! , 441 - is - @ , 442 - is - # , 443 - is - $ , 444 - is - % , 445 - is - & , 446 - is - * , 447 - is - ( , 448 - is - ) , 449 - is - , , 450 - is - - , 451 - is - . , 452 - is - / , 453 - is - : , 454 - is - ; , 455 - is - < , 456 - is - = , 457 - is - > , 458 - is - ? , 459 - is - [ , 460 - is - \ , 461 - is - ] , 462 - is - ^ , 463 - is - _ , 464 - is - ` , 465 - is - ~ , 466 - is - space , 467 - is - ! , 468 - is - @ , 469 - is - # , 470 - is - $ , 471 - is - % , 472 - is - & , 473 - is - * , 474 - is - ( , 475 - is - ) , 476 - is - , , 477 - is - - , 478 - is - . , 479 - is - / , 480 - is - : , 481 - is - ; , 482 - is - < , 483 - is - = , 484 - is - > , 485 - is - ? , 486 - is - [ , 487 - is - \ , 488 - is - ] , 489 - is - ^ , 490 - is - _ , 491 - is - ` , 492 - is - ~ , 493 - is - space , 494 - is - ! , 495 - is - @ , 496 - is - # , 497 - is - $ , 498 - is - % , 499 - is - & , 500 - is - * , 501 - is - ( , 502 - is - ) , 503 - is - , , 504 - is - - , 505 - is - . , 506 - is - / , 507 - is - : , 508 - is - ; , 509 - is - < , 510 - is - = , 511 - is - > , 512 - is - ? , 513 - is - [ , 514 - is - \ , 515 - is - ] , 516 - is - ^ , 517 - is - _ , 518 - is - ` , 519 - is - ~ , 520 - is - space , 521 - is - ! , 522 - is - @ , 523 - is - # , 524 - is - $ , 525 - is - % , 526 - is - & , 527 - is - * , 528 - is - ( , 529 - is - ) , 530 - is - , , 531 - is - - , 532 - is - . , 533 - is - / , 534 - is - : , 535 - is - ; , 536 - is - < , 537 - is - = , 538 - is - > , 539 - is - ? , 540 - is - [ , 541 - is - \ , 542 - is - ] , 543 - is - ^ , 544 - is - _ , 545 - is - ` , 546 - is - ~ , 547 - is - space , 548 - is - ! , 549 - is - @ , 550 - is - # , 551 - is - $ , 552 - is - % , 553 - is - & , 554 - is - * , 555 - is - ( , 556 - is - ) , 557 - is - , , 558 - is - - , 559 - is - . , 560 - is - / , 561 - is - : , 562 - is - ; , 563 - is - < , 564 - is - = , 565 - is - > , 566 - is - ? , 567 - is - [ , 568 - is - \ , 569 - is - ] , 570 - is - ^ , 571 - is - _ , 572 - is - ` , 573 - is - ~ , 574 - is - space , 575 - is - ! , 576 - is - @ , 577 - is - # , 578 - is - $ , 579 - is - % , 580 - is - & , 581 - is - * , 582 - is - ( , 583 - is - ) , 584 - is - , , 585 - is - - , 586 - is - . , 587 - is - / , 588 - is - : , 589 - is - ; , 590 - is - < , 591 - is - = , 592 - is - > , 593 - is - ? , 594 - is - [ , 595 - is - \ , 596 - is - ] , 597 - is - ^ , 598 - is - _ , 599 - is - ` , 600 - is - ~ , 601 - is - space , 602 - is - ! , 603 - is - @ , 604 - is - # , 605 - is - $ , 606 - is - % , 607 - is - & , 608 - is - * , 609 - is - ( , 610 - is - ) , 611 - is - , , 612 - is - - , 613 - is - . , 614 - is - / , 615 - is - : , 616 - is - ; , 617 - is - < , 618 - is - = , 619 - is - > , 620 - is - ? , 621 - is - [ , 622 - is - \ , 623 - is - ] , 624 - is - ^ , 625 - is - _ , 626 - is - ` , 627 - is - ~ , 628 - is - space , 629 - is - ! , 630 - is - @ , 631 - is - # , 632 - is - $ , 633 - is - % , 634 - is - & , 635 - is - * , 636 - is - ( , 637 - is - ) , 638 - is - , , 639 - is - - , 640 - is - . , 641 - is - / , 642 - is - : , 643 - is - ; , 644 - is - < , 645 - is - = , 646 - is - > , 647 - is - ? , 648 - is - [ , 649 - is - \ , 650 - is - ] , 651 - is - ^ , 652 - is - _ , 653 - is - ` , 654 - is - ~ , 655 - is - space , 656 - is - ! , 657 - is - @ , 658 - is - # , 659 - is - $ , 660 - is - % , 661 - is - & , 662 - is - * , 663 - is - ( , 664 - is - ) , 665 - is - , , 666 - is - - , 667 - is - . , 668 - is - / , 669 - is - : , 670 - is - ; , 671 - is - < , 672 - is - = , 673 - is - > , 674 - is - ? , 675 - is - [ , 676 - is - \ , 677 - is - ] , 678 - is - ^ , 679 - is - _ , 680 - is - ` , 681 - is - ~ , 682 - is - space , 683 - is - ! , 684 - is - @ , 685 - is - # , 686 - is - $ , 687 - is - % , 688 - is - & , 689 - is - * , 690 - is - ( , 691 - is - ) , 692 - is - , , 693 - is - - , 694 - is - . , 695 - is - / , 696 - is - : , 697 - is - ; , 698 - is - < , 699 - is - = , 700 - is - > , 701 - is - ? , 702 - is - [ , 703 - is - \ , 704 - is - ] , 705 - is - ^ , 706 - is - _ , 707 - is - ` , 708 - is - ~ , 709 - is - space , 710 - is - ! , 711 - is - @ , 712 - is - # , 713 - is - $ , 714 - is - % , 715 - is - & , 716 - is - * , 717 - is - ( , 718 - is - ) , 719 - is - , , 720 - is - - , 721 - is - . , 722 - is - / , 723 - is - : , 724 - is - ; , 725 - is - < , 726 - is - = , 727 - is - > , 728 - is - ? , 729 - is - [ , 730 - is - \ , 731 - is - ] , 732 - is - ^ , 733 - is - _ , 734 - is - ` , 735 - is - ~ , 736 - is - space , 737 - is - ! , 738 - is - @ , 739 - is - # , 740 - is - $ , 741 - is - % , 742 - is - & , 743 - is - * , 744 - is - ( , 745 - is - ) , 746 - is - , , 747 - is - - , 748 - is - . , 749 - is - / , 750 - is - : , 751 - is - ; , 752 - is - < , 753 - is - = , 754 - is - > , 755 - is - ? , 756 - is - [ , 757 - is - \ , 758 - is - ] , 759 - is - ^ , 760 - is - _ , 761 - is - ` , 762 - is - ~ , 763 - is - space , 764 - is - ! , 765 - is - @ , 766 - is - # , 767 - is - $ , 768 - is - % , 769 - is - & , 770 - is - * , 771 - is - ( , 772 - is - ) , 773 - is - , , 774 - is - - , 775 - is - . , 776 - is - / , 777 - is - : , 778 - is - ; , 779 - is - < , 780 - is - = , 781 - is - > , 782 - is - ? , 783 - is - [ , 784 - is - \ , 785 - is - ] , 786 - is - ^ , 787 - is - _ , 788 - is - ` , 789 - is - ~ , 790 - is - space , 791 - is - ! , 792 - is - @ , 793 - is - # , 794 - is - $ , 795 - is - % , 796 - is - & , 797 - is - * , 798 - is - ( , 799 - is - ) , 800 - is - , , 801 - is - - , 802 - is - . , 803 - is - / , 804 - is - : , 805 - is - ; , 806 - is - < , 807 - is - = , 808 - is - > , 809 - is - ? , 810 - is - [ , 811 - is - \ , 812 - is - ] , 813 - is - ^ , 814 - is - _ , 815 - is - ` , 816 - is - ~ , 817 - is - space , 818 - is - ! , 819 - is - @ , 820 - is - # , 821 - is - $ , 822 - is - % , 823 - is - & , 824 - is - * , 825 - is - ( , 826 - is - ) , 827 - is - , , 828 - is - - , 829 - is - . , 830 - is - / , 831 - is - : , 832 - is - ; , 833 - is - < , 834 - is - = , 835 - is - > , 836 - is - ? , 837 - is - [ , 838 - is - \ , 839 - is - ] , 840 - is - ^ , 841 - is - _ , 842 - is - ` , 843 - is - ~ , 844 - is - space , 845 - is - ! , 846 - is - @ , 847 - is - # , 848 - is - $ , 849 - is - % , 850 - is - & , 851 - is - * , 852 - is - ( , 853 - is - ) , 854 - is - , , 855 - is - - , 856 - is - . , 857 - is - / , 858 - is - : , 859 - is - ; , 860 - is - < , 861 - is - = , 862 - is - > , 863 - is - ? , 864 - is - [ , 865 - is - \ , 866 - is - ] , 867 - is - ^ , 868 - is - _ , 869 - is - ` , 870 - is - ~ , 871 - is - space , 872 - is - ! , 873 - is - @ , 874 - is - # , 875 - is - $ , 876 - is - % , 877 - is - & , 878 - is - * , 879 - is - ( , 880 - is - ) , 881 - is - , , 882 - is - - , 883 - is - . , 884 - is - / , 885 - is - : , 886 - is - ; , 887 - is - < , 888 - is - = , 889 - is - > , 890 - is - ? , 891 - is - [ , 892 - is - \ , 893 - is - ] , 894 - is - ^ , 895 - is - _ , 896 - is - ` , 897 - is - ~ , 898 - is - space , 899 - is - ! , 900 - is - @ , 901 - is - # , 902 - is - $ , 903 - is - % , 904 - is - & , 905 - is - * , 906 - is - ( , 907 - is - ) , 908 - is - , , 909 - is - - , 910 - is - . , 911 - is - / , 912 - is - : , 913 - is - ; , 914 - is - < , 915 - is - = , 916 - is - > , 917 - is - ? , 918 - is - [ , 919 - is - \ , 920 - is - ] , 921 - is - ^ , 922 - is - _ , 923 - is - ` , 924 - is - ~ , 925 - is - space , 926 - is - ! , 927 - is - @ , 928 - is - # , 929 - is - $ , 930 - is - % , 931 - is - & , 932 - is - * , 933 - is - ( , 934 - is - ) , 935 - is - , , 936 - is - - , 937 - is - . , 938 - is - / , 939 - is - : , 940 - is - ; , 941 - is - < , 942 - is - = , 943 - is - > , 944 - is - ? , 945 - is - [ , 946 - is - \ , 947 - is - ] , 948 - is - ^ , 949 - is - _ , 950 - is - ` , 951 - is - ~ , 952 - is - space , 953 - is - ! , 954 - is - @ , 955 - is - # , 956 - is - $ , 957 - is - % , 958 - is - & , 959 - is - * , 960 - is - ( , 961 - is - ) , 962 - is - , , 963 - is - - , 964 - is - . , 965 - is - / , 966 - is - : , 967 - is - ; , 968 - is - < , 969 - is - = , 970 - is - > , 971 - is - ? , 972 - is - [ , 973 - is - \ , 974 - is - ] , 975 - is - ^ , 976 - is - _ , 977 - is - ` , 978 - is - ~ , 979 - is - space , 980 - is - ! , 981 - is - @ , 982 - is - # , 983 - is - $ , 984 - is - % , 985 - is - & , 986 - is - * , 987 - is - ( , 988 - is - ) , 989 - is - , , 990 - is - - , 991 - is - . , 992 - is - / , 993 - is - : , 994 - is - ; , 995 - is - < , 996 - is - = , 997 - is - > , 998 - is - ? , 999 - is - [ , 1000 - is - \ , 1001 - is - ] , 1002 - is - ^ , 1003 - is - _ , 1004 - is - ` , 1005 - is - ~ , 1006 - is - space , 1007 - is - ! , 1008 - is - @ , 1009 - is - # , 1010 - is - $ , 1011 - is - % , 1012 - is - & , 1013 - is - * , 1014 - is - ( , 1015 - is - ) , 1016 - is - , , 1017 - is - - , 1018 - is - . , 1019 - is - / , 1020 - is - : , 1021 - is - ; , 1022 - is - < , 1023 - is - = , 1024 - is - > , 1025 - is - ? , 1026 - is - [ , 1027 - is - \ , 1028 - is - ] , 1029 - is - ^ , 1030 - is - _ , 1031 - is - ` , 1032 - is - ~ , 1033 - is - space , 1034 - is - ! , 1035 - is - @ , 1036 - is - # , 1037 - is - $ , 1038 - is - % , 1039 - is - & , 1040 - is - * , 1041 - is - ( , 1042 - is - ) , 1043 - is - , , 1044 - is - - , 1045 - is - . , 1046 - is - / , 1047 - is - : , 1048 - is - ; , 1049 - is - < , 1050 - is - = , 1051 - is - > , 1052 - is - ? , 1053 - is - [ , 1054 - is - \ , 1055 - is - ] , 1056 - is - ^ , 1057 - is - _ , 1058 - is - ` , 1059 - is - ~ , 1060 - is - space , 1061 - is - ! , 1062 - is - @ , 1063 - is - # , 1064 - is - $ , 1065 - is - % , 1066 - is - & , 1067 - is - * , 1068 - is - ( , 1069 - is - ) , 1070 - is - , , 1071 - is - - , 1072 - is - . , 1073 - is - / , 1074 - is - : , 1075 - is - ; , 1076 - is - < , 1077 - is - = , 1078 - is - > , 1079 - is - ? , 1080 - is - [ , 1081 - is - \ , 1082 - is - ] , 1083 - is - ^ , 1084 - is - _ , 1085 - is - ` , 1086 - is - ~ , 1087 - is - space , 1088 - is - ! , 1089 - is - @ , 1090 - is - # , 1091 - is - $ , 1092 - is - % , 1093 - is - & , 1094 - is - * , 1095 - is - ( , 1096 - is - ) , 1097 - is - , , 1098 - is - - , 1099 - is - . , 1100 - is - / , 1101 - is - : , 1102 - is - ; , 1103 - is - < , 1104 - is - = , 1105 - is - > , 1106 - is - ? , 1107 - is - [ , 1108 - is - \ , 1109 - is - ] , 1110 - is - ^ , 1111 - is - _ , 1112 - is - ` , 1113 - is - ~ , 1114 - is - space , 1115 - is - ! , 1116 - is - @ , 1117 - is - # , 1118 - is - $ , 1119 - is - % , 1120 - is - & , 1121 - is - * , 1122 - is - ( , 1123 - is - ) , 1124 - is - , , 1125 - is - - , 1126 - is - . , 1127 - is - / , 1128 - is - : , 1129 - is - ; , 1130 - is - < , 1131 - is - = , 1132 - is - > , 1133 - is - ? , 1134 - is - [ , 1135 - is - \ , 1136 - is - ] , 1137 - is - ^ , 1138 - is - _ , 1139 - is - ` , 1140 - is - ~ , 1141 - is - space , 1142 - is - ! , 1143 - is - @ , 1144 - is - # , 1145 - is - $ , 1146 - is - % , 1147 - is - & , 1148 - is - * , 1149 - is - ( , 1150 - is - ) , 1151 - is - , , 1152 - is - - , 1153 - is - . , 1154 - is - / , 1155 - is - : , 1156 - is - ; , 1157 - is - < , 1158 - is - = , 1159 - is - > , 1160 - is - ? , 1161 - is - [ , 1162 - is - \ , 1163 - is - ] , 1164 - is - ^ , 1165 - is - _ , 1166 - is - ` , 1167 - is - ~ , 1168 - is - space , 1169 - is - ! , 1170 - is - @ , 1171 - is - # , 1172 - is - $ , 1173 - is - % , 1174 - is - & , 1175 - is - * , 1176 - is - ( , 1177 - is - ) , 1178 - is - , , 1179 - is - - , 1180 - is - . , 1181 - is - / , 1182 - is - : , 1183 - is - ; , 1184 - is - < , 1185 - is - = , 1186 - is - > , 1187 - is - ? , 1188 - is - [ , 1189 - is - \ , 1190 - is - ] , 1191 - is - ^ , 1192 - is - _ , 1193 - is - ` , 1194 - is - ~ , 1195 - is - space , 1196 - is - ! , 1197 - is - @ , 1198 - is - # , 1199 - is - $ , 1200 - is - % , 1201 - is - & , 1202 - is - * , 1203 - is - ( , 1204 - is - ) , 1205 - is - , , 1206 - is - - , 1207 - is - . , 1208 - is - / , 1209 - is - : , 1210 - is - ; , 1211 - is - < , 1212 - is - = , 1213 - is - > , 1214 - is - ? , 1215 - is - [ , 1216 - is - \ , 1217 - is - ] , 1218 - is - ^ , 1219 - is - _ , 1220 - is - ` , 1221 - is - ~ , 1222 - is - space , 1223 - is - ! , 1224 - is - @ , 1225 - is - # , 1226 - is - $ , 1227 - is - % , 1228 - is - & , 1229 - is - * , 1230 - is - ( , 1231 - is - ) , 1232 - is - , , 1233 - is - - , 1234 - is - . , 1235 - is - / , 1236 - is - : , 1237 - is - ; , 1238 - is - < , 1239 - is - = , 1240 - is - > , 1241 - is - ? , 1242 - is - [ , 1243 - is - \ , 1244 - is - ] , 1245 - is - ^ , 1246 - is - _ , 1247 - is - ` , 1248 - is - ~ , 1249 - is - space , 1250 - is - ! , 1251 - is - @ , 1252 - is - # , 1253 - is - $ , 1254 - is - % , 1255 - is - & , 1256 - is - * , 1257 - is - ( , 1258 - is - ) , 1259 - is - , , 1260 - is - - , 1261 - is - . , 1262 - is - / , 1263 - is - : , 1264 - is - ; , 1265 - is - < , 1266 - is - = , 1267 - is - > , 1268 - is - ? , 1269 - is - [ , 1270 - is - \ , 1271 - is - ] , 1272 - is - ^ , 1273 - is - _ , 1274 - is - ` , 1275 - is - ~ , 1276 - is - space , 1277 - is - ! , 1278 - is - @ , 1279 - is - # , 1280 - is - $ , 1281 - is - % , 1282 - is - & , 1283 - is - * , 1284 - is - ( , 1285 - is - ) , 1286 - is - , , 1287 - is - - , 1
```

Substring + replacement based on the positions

- `substr()` and `substring()`
- These function can be used to overwrite or replace a part of character string.

```
x <- "Lusine Zilfimian"  
substr(x, start = 8, stop = 7 + 5)
```

```
## [1] "Zilfi"
```

```
substring(x, first = 8) # last = 1000000L
```

```
## [1] "Zilfimian"
```

```
substr(x, start = 14, stop = 14) <- "y"  
x
```

```
## [1] "Lusine Zilfimyan"
```

Replacement and detection

```
sub("Lusine Zilfimian", pattern = "i", replacement = "y") # find and replace
```

```
## [1] "Lusyne Zilfimian"
```

```
gsub("Lusine Zilfimian", pattern = "i", replacement = "y")
```

```
## [1] "Lusyne Zylfymyan"
```

```
grep(pattern = "Lusine", x = c("Lusine Zilfimian", "David", "Mar")) #detect
```

```
## [1] 1
```

```
grepl(pattern = "Lusine", x = c("Lusine Zilfimian", "David", "Mar"))
```

```
## [1] TRUE FALSE FALSE
```

Abbreviations

- Abbreviate strings to at least minlength characters, such that they remain unique (if they were), unless strict = TRUE.

```
x <- c("Lusin ", " Lusine ", "Lusine Zilfimian", "Lusine Zilfi", "Lus - ",
(abb1 <- abbreviate(x, minlength = 4))
```

##	Lusin	Lusine	Lusine Zilfimian	Lusine Zilfi
##	"Lusin"	"Lusine"	"LsnZlfm"	"LusnZlf"
##	Lus -	Lusine 1995		
##	"Lus-"	"L199"		

```
(abb2 <- abbreviate(x, minlength = 4, strict = T))
```

##	Lusin	Lusine	Lusine Zilfimian	Lusine Zilfi
##	"Lusn"	"Lusn"	"LsnZ"	"LsnZ"
##	Lus -	Lusine 1995		
##	"Lus-"	"L199"		

Why is base not enough?

- The objects `NULL` and `character(0)` have zero length, yet when included inside `paste()` they are treated as an empty string:

```
paste("This", "is", NULL, character(0), "value")
```

```
## [1] "This is  value"
```

```
length(NULL)
```

```
## [1] 0
```

```
length(character(0))
```

```
## [1] 0
```

```
stringr::str_c("This", "is", NULL, character(0), "value", sep = " ")
```

```
## [1] "This is value"
```

- `str_c()` – zero argument are silently removed

- All functions start with str_
- All functions take a vector of strings as the first argument
- str_sub works the same way as substr() but allows for negative subsetting as well

```
library(stringr)
str_sub("Joe Satriani", start = 5)
```

```
## [1] "Satriani"
```

```
str_sub("Joe Satriani", start = 5, end = 5+4)
```

```
## [1] "Satri"
```

```
str_sub("Joe Satriani", start = -8, end = -4)
```

```
## [1] "Satri"
```


- `str_detect()` - For detecting whether a pattern is present (or absent) in a string vector. You get a TRUE if a match is detected in a string, FALSE otherwise
- `str_subset()` - Keep strings matching a pattern
- `str_count()` - Count the number of matches in a string
- `str_split()` - Similar to `strsplit()` to separate a character vector into a number of pieces
- `str_replace()` - For replacing the first occurrence of a matched pattern in a string
- `str_extract()` - For extracting a string containing a pattern

`detect` - boolean, `extract` - only pattern, `subset` - all the text.

Examples

```
(generators <- c("Action, Adventure, Comedy", "Comedy",  
  "Comedy, Drama, Drama, Romance", "Crime, Drama, History"))
```

```
## [1] "Action, Adventure, Comedy"      "Comedy"  
## [3] "Comedy, Drama, Drama, Romance"  "Crime, Drama, History"
```

```
str_detect(string = generators, pattern = "Drama")
```

```
## [1] FALSE FALSE TRUE TRUE
```

```
str_subset(string = generators, pattern = "Action")
```

```
## [1] "Action, Adventure, Comedy"
```

```
str <- c("123abd", "lz@yahoo.com", "Name567cd", "abc5.00", "lusinezilfimian@  
str_subset(str, pattern="@yahoo.com")
```

```
## [1] "lz@yahoo.com"      "lusinezilfimian@yahoo.com"
```

```
str_count(string = generators, pattern = "Drama")
```

Examples

```
str_replace(string = genres, pattern = ",", replacement = " &")
```

```
## [1] "Action & Adventure, Comedy"      "Comedy"  
## [3] "Comedy & Drama, Drama, Romance"  "Crime & Drama, History"
```

```
sub(x = genres, pattern = ",", replacement = " &") # from base
```

```
## [1] "Action & Adventure, Comedy"      "Comedy"  
## [3] "Comedy & Drama, Drama, Romance"  "Crime & Drama, History"
```

```
gsub(x = genres, pattern = ",", replacement = " &") # from base
```

```
## [1] "Action & Adventure & Comedy"      "Comedy"  
## [3] "Comedy & Drama & Drama & Romance"  "Crime & Drama & History"
```

```
str_replace_all(string = genres, pattern = ",", replacement = " &")
```

```
## [1] "Action & Adventure & Comedy"      "Comedy"  
## [3] "Comedy & Drama & Drama & Romance"  "Crime & Drama & History"
```

Examples

```
str_extract(string = genres, pattern = "Drama")
```

```
## [1] NA      NA      "Drama" "Drama"
```

```
str_extract_all(string = genres, pattern = "Drama", simplify = T)
```

```
##      [,1]      [,2]
## [1,] ""       ""
## [2,] ""       ""
## [3,] "Drama"  "Drama"
## [4,] "Drama"  ""
```

Let's consider the example

```
movies <- read.csv('movies3.csv', stringsAsFactors = F)
```

How many romcom movies are there?

```
movies$Comedy <- str_detect(string = movies$Genre, pattern = "Comedy")
movies$rom <- str_detect(string = movies$Genre, pattern = "Romance")
table(movies$Comedy, movies$rom)
```

```
##
##          FALSE TRUE
## FALSE  1478  202
## TRUE   886  319
```

Split string into pieces

```
(gen_m <- str_split(geners, ",", simplify = T))
```

```
##      [,1]      [,2]      [,3]      [,4]
## [1,] "Action" " Adventure" " Comedy" ""
## [2,] "Comedy" ""          ""          ""
## [3,] "Comedy" " Drama"    " Drama"    " Romance"
## [4,] "Crime"  " Drama"    " History"  ""
```

```
trimws("  Lus      ine  ") # trim whitespaces
```

```
## [1] "Lus      ine"
```

Regular Expressions

- Sequence of characters that define a search pattern
- For instance, one could have the following pattern: 2 digits, 2 letters and 4 digits
- ?regex - to see more

Form of Regular Expressions

- ① Literal characters which are characters that match themselves
- ② Metacharacters - any character that is not a literal character . ^ \$ + * ? | (), { }, [], -
- ③ Sequences
- ④ Character classes
- ⑤ Quantifiers
- ⑥ POSIX character classes
- If you want to match the metacharacter you need to use \ (escapes) before the metacharacter

Literal and Metacharacters

- Literal character - the simplest form of regular expressions are those that match a single character
- Metacharacters - special characters that have a reserved status:

```
str <-c("Lusine1995^", "Say$1995-02(", "David.", "M*77", "M&95", "1995.27")  
str_subset(string = str, pattern = "Lusine") # literal
```

```
## [1] "Lusine1995^"
```

```
str_subset(string = str, pattern = "Lusine|David|Say")
```

```
## [1] "Lusine1995^" "Say$1995-02(" "David."
```

```
str_subset(string = "money$", pattern = "money\\$")
```

```
## [1] "money$"
```

```
str_subset(string = "This is money1", pattern = "money$")
```

```
## character(0)
```


Metacharacters

```
str_subset(string = str, pattern = ".") # meta
```

```
## [1] "Lusine1995^" "Say$1995-02(" "David." "M*77"  
## [5] "M&95" "1995.27"
```

```
str_subset(string = "\n", pattern = ".") # meta
```

```
## character(0)
```

```
str_subset(string = " ", pattern = ".") # meta
```

```
## [1] " "
```

```
str_subset(string = "\"", pattern = ".") # meta
```

```
## character(0)
```

```
str_subset(string = str, pattern = "\\.") # literal
```

```
## [1] "David." "1995.27"
```

Metacharacters

```
str_subset(string = str, pattern = "^M") # meta
```

```
## [1] "M*77" "M&95"
```

```
str_subset(string = str, pattern = "\\^") # literal
```

```
## [1] "Lusine1995^"
```

```
str_subset(string = str, pattern = "\\$") # literal
```

```
## [1] "Say$1995-02("
```

Sequences

```
str_subset(string = c("123", "4L", "M"), pattern = "\\d") # [0-9]
```

```
## [1] "123" "4L"
```

```
str_subset(string = c("123", "4L", "M"), pattern = "\\D") # [^0-9]
```

```
## [1] "4L" "M"
```

```
str_subset(string = c(" ", "l ", "l"), pattern = "\\s")
```

```
## [1] " " "l "
```

```
str_subset(string = c(" ", "l ", "l"), pattern = "\\S")
```

```
## [1] "l " "l"
```

```
str_subset(string = c(" ", "l ", "l"), pattern = "\\w") # word
```

```
## [1] "l " "l"
```

```
str_subset(string = c(" ", "l ", "l"), pattern = "\\W") # non-word
```

Character classes

```
str_subset(string = str, pattern = "[0-9]")
```

```
## [1] "Lusine1995^" "Say$1995-02(" "M*77" "M&95"  
## [5] "1995.27"
```

```
str_subset(string = str, pattern = "[f-z]")
```

```
## [1] "Lusine1995^" "Say$1995-02(" "David."
```

```
str_subset(string = str, pattern = "[a-zA-Z]")
```

```
## [1] "Lusine1995^" "Say$1995-02(" "David." "M*77"  
## [5] "M&95"
```

```
str_subset(string = str, pattern = "[aeuio]")
```

```
## [1] "Lusine1995^" "Say$1995-02(" "David."
```

```
str_subset(string = c("AI", "ai"), pattern = "[AEIOU]")
```

```
## [1] "AI"
```

Special meaning of ^

- Not in specified character

```
str_subset(string = c("ou", "lk"), pattern = "[^aeuio]")
```

```
## [1] "lk"
```

```
str_subset(string = c("ou", "19", "-"), pattern = "[a-zA-Z0-9]")
```

```
## [1] "ou" "19"
```

```
str_subset(string = c("ou", "19", "-"), pattern = "[a-zA-Z]")
```

```
## [1] "ou"
```

```
str_subset(string = c("ou", "19", "-"), pattern = "[^0-9]")
```

```
## [1] "ou" "-"
```

```
str_subset(string = str, pattern = "[^1-9]")
```

```
## [1] "Lusine1995^" "Say$1995-02(" "David." "M*77"
```

```
## [5] "M&95" "1995.27"
```

Quantifiers

```
str_subset(string = c(" ", "1 ", "1"), pattern = "\\s+")
```

```
## [1] " " "1 "
```

```
str_subset(string = c(" ", "1 ", "1"), pattern = "\\s*")
```

```
## [1] " " "1 " "1"
```

```
str_subset(string = c(" ", "1 ", "1"), pattern = "\\s?")
```

```
## [1] " " "1 " "1"
```

```
str_subset(string = c(" ", "1 ", "1"), pattern = "\\s{3}")
```

```
## [1] "1 "
```

```
str_subset(string = c(" ", "1 ", "1"), pattern = "\\s{2,}")
```

```
## [1] "1 "
```

```
str_subset(string = c(" ", "1 ", "1"), pattern = "\\s{1,3}")
```

POSIX character classes

```
str_subset(string = str, pattern = "[:alnum:]") # [a-zA-Z0-9]
```

```
## [1] "Lusine1995^" "Say$1995-02(" "David." "M*77"  
## [5] "M&95" "1995.27"
```

```
str_subset(string = str, pattern = "[:alpha:]") # [a-zA-Z]
```

```
## [1] "Lusine1995^" "Say$1995-02(" "David." "M*77"  
## [5] "M&95"
```

```
str_subset(string = str, pattern = "[:digit:]") # [0-9]
```

```
## [1] "Lusine1995^" "Say$1995-02(" "M*77" "M&95"  
## [5] "1995.27"
```

```
str_subset(string = str, pattern = "[:lower:]") # [a-z]
```

```
## [1] "Lusine1995^" "Say$1995-02(" "David."
```

POSIX character classes

```
str_subset(string = str, pattern = "[:upper:]") # [A-Z]
```

```
## [1] "Lusine1995^" "Say$1995-02(" "David." "M*77"  
## [5] "M&95"
```

```
str_subset(string = str, pattern = "[:space:]") # [\f\n\r\t\v]
```

```
## character(0)
```

```
str_subset(string = str, pattern = "[:punct:]")
```

```
## [1] "Say$1995-02(" "David." "M*77" "M&95"  
## [5] "1995.27"
```

```
str_subset(string = c("grey", "gray"), pattern = "gr(e|a)y")
```

```
## [1] "grey" "gray"
```


How many wins and nominations are there?

```
movies$Awards[1:5]
```

```
## [1] "4 wins & 8 nominations." "7 wins & 17 nominations."  
## [3] "1 win & 2 nominations."  "9 wins & 28 nominations."  
## [5] "2 wins & 67 nominations."
```

```
movies$awards_num <- str_replace_all(movies$Awards, pattern="^[^0-9]",  
  replacement = " ")
```

```
movies$awards_num <- str_replace_all(movies$awards_num, pattern="\\s+",  
  replacement = " ")
```

How many wins and nominations are there?

```
movies$awards_num <- trimws(movies$awards_num)
movies$awards_num[1:20]
```

```
## [1] "4 8"      "7 17"      "1 2"      "9 28"      "2 67"      "1"
## [7] "1 87 171" "2 3"       "2 4"      "1 5 7"     "1 1"       "1 3"
## [13] "1 1 3"     "1 2"       "2 11 5"   "1 3 31"    "2 2"       "2 2"
## [19] "1 3"      "1 5"
```

```
x1 <- str_split(movies$awards_num, pattern = " ",simplify = T)
x1 <- apply(x1, 2, as.numeric)
x1 <- rowSums(x1, na.rm=T)
x1[1:20]
```

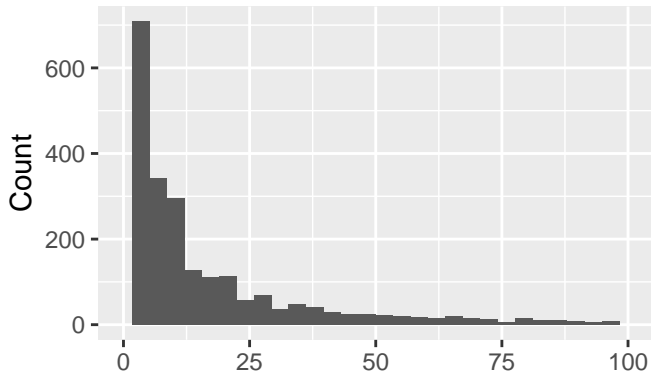
```
## [1] 12 24 3 37 69 1 259 5 6 13 2 4 5 3 18 35 4
## [18] 4 4 6
```

Visualization

```
movies$awards_num <- x1
ggplot(movies, aes(x = awards_num)) + geom_histogram() +
  xlim(c(0,100)) + labs(title = "Awards histogram (wins and nominations)",
    x = "Awards", y = "Count")
```

`stat_bin()` using `bins = 30`. Pick better value with `binwidth`.

Awards histogram (wins and nominati



Wins only

- Possible strings are ...

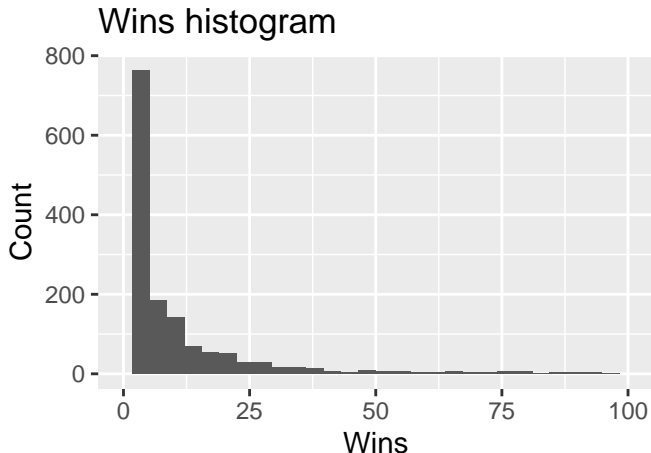
```
m <- movies$Awards[1:20]
m1 <- str_extract_all(movies$Awards, pattern="[0-9]+\\swin", simplify = T)
m1 <- str_remove_all(m1, pattern="[a-zA-Z]")
m1 <- as.numeric(m1)
movies$Wins <- m1
head(m1)
```

```
## [1] 4 7 1 9 2 NA
```

Visualization

```
ggplot(movies, aes(x = Wins)) + geom_histogram() +  
  labs(title = "Wins histogram", x = "Wins", y = "Count") +  
  xlim(c(0,100))
```

`stat_bin()` using `bins = 30`. Pick better value with `binwidth`.



Task

- Read the file Text.txt without warnings. Use regular expressions to solve the problems:
- Show only the first line. Eliminate the periods before and after the words, but not periods inside the values and at the end of the statement.
- Show only the second line. Convert the GDP of Armenia into proper numbers, create a data frame with variables year and GDP.

Solution

```
readLines("Text.txt", warn=FALSE)
```

```
## [1] "..The. ...economy. ...of ..Armenia ...grew ..by. 5.2% in ..2018. and  
## [2] "The GDP of Armenia was ?4,000,722.0 millions in 2012, ?4,555,638.2 m  
## [3] "Even in this small mountainous country, there are places where they  
## [4] "Armenian cuisine is wonderful and varied, with only kharavts and sh  
## [5] "The most popular tours in Armenia are tours to mountain lakes among
```

```
(text1 <- readLines("Text.txt", warn=FALSE)[1])
```

```
## [1] "..The. ...economy. ...of ..Armenia ...grew ..by. 5.2% in ..2018. and
```

```
words <- unlist(str_split(text1, pattern = " "))  
words <- str_replace_all(string = words, pattern = "^\\.\\.", replacement = "  
words <- str_replace_all(string = words, pattern = "\\..+$", replacement = "  
paste(words, sep = " ", collapse = " ")
```

```
## [1] "The economy of Armenia grew by 5.2% in 2018 and reached a nominal G
```

Solution

```
(text2 <- readLines("Text.txt", warn=FALSE)[2])
```

```
## [1] "The GDP of Armenia was ?4,000,722.0 millions in 2012, ?4,555,638.2 "
```

```
words2 <- unlist(str_split(text2, pattern = " "))
```

```
numbers <- str_subset(words2, pattern = "[0-9]")
```

```
numbers <- str_replace(numbers, pattern = "\\.[0-9]", replacement = "")
```

```
numbers <- str_replace_all(numbers, pattern = "^\\.|$,|\\.\\.|[,]", replacement = "")
```

```
df = data.frame(matrix(numbers, ncol = 2, byrow = T))
```

```
colnames(df)= c("gdp", "year")
```

```
df
```

```
##      gdp year
```

```
## 1 4000722 2012
```

```
## 2 4555638 2013
```

```
## 3 4828626 2014
```

```
## 4 5032089 2015
```


Text Mining

- We will use the tm text mining package.

```
library(tm)
```

- Load the data into a Corpus (a collection of documents) which is the main data structure used by tm
- In order to create a VCorpus using tm, we need to pass a "Source" object as a paramter to the 'VCorpus
- First you need to create a VectorSource - A vector source interprets each element of the vector x as a document
- VectorSource is for only character vectors.

```
text <- c("The best guitar models are Gibson and Fender!",  
  "My guitars are the best guitars.",  
  "Are these guitars yours?")  
vs <- VectorSource(text)  
corpus <- VCorpus(vs)
```

- use `[[]]` to access single element in the corpus

```
inspect(corpus[[3]]) # like print
```

```
## <<PlainTextDocument>>
## Metadata:  7
## Content:  chars: 24
##
## Are these guitars yours?
```

- Access content directly for the second document

```
corpus[[2]][1]
```

```
## $content
## [1] "My guitars are the best guitars."
```

```
corpus[[2]][2]
```

```
## $meta
##   author      : character(0)
##   timestamp: 2020-05-20 19:41:38
##   description : character(0)
##   heading     : character(0)
##   id          : 2
##   language    : en
##   origin      : character(0)
```

- Metadata is used to annotate text documents or whole corpora with additional information. The easiest way to accomplish this with tm is to use the meta() function

```
meta(corpus, tag = 'language')
```

```
## $`1`
## [1] "en"
##
## $`2`
## [1] "en"
##
```

- Create new metadata tag class for the first document and assign Guitar to it
- Same for IDs

```
meta(corpus[[1]], tag="class") <- "Guitar"  
corpus[[1]][2]
```

```
## $meta  
##   author      : character(0)  
##   timestamp: 2020-05-20 19:41:38  
##   description : character(0)  
##   heading     : character(0)  
##   id          : 1  
##   language    : en  
##   origin      : character(0)  
##   class       : Guitar
```

```
meta(corpus, tag = "ids") <- c("DOC1", "DOC2", "DOC3")  
meta(corpus, "ids")
```

```
##   ids  
## 1 DOC1  
## 2 DOC2  
## 3 DOC3
```

Document-term matrix or term-document matrix

- A matrix that describes the frequency of terms that occur in a collection of documents.
- In a document-term matrix, rows correspond to documents in the collection and columns correspond to

```
tdm <- TermDocumentMatrix(corpus)
inspect(tdm)
```

```
## <<TermDocumentMatrix (terms: 12, documents: 3)>>
## Non-/sparse entries: 17/19
## Sparsity           : 53%
## Maximal term length: 8
## Weighting          : term frequency (tf)
## Sample            :
##
##      Docs
## Terms  1 2 3
##  and   1 0 0
##  are   1 1 1
##  best  1 1 0
##  fender! 1 0 0
##  gibson 1 0 0
```

- Document term matrix is the transposed Term document matrix

```
dtm <- DocumentTermMatrix(corpus, control = list(weightTfIdf))  
inspect(dtm)
```

```
## <<DocumentTermMatrix (documents: 3, terms: 12)>>  
## Non-/sparse entries: 17/19  
## Sparsity           : 53%  
## Maximal term length: 8  
## Weighting          : term frequency (tf)  
## Sample            :  
##      Terms  
## Docs and are best fender! gibson guitar guitars guitars. models the  
##    1  1  1  1      1      1      1      0      0      1  1  
##    2  0  1  1      0      0      0      1      1      0  1  
##    3  0  1  0      0      0      0      1      0      0  0
```

- Sparsity is defined as a percentage of element with value 0 in total number of elements in document term matrix

```
mm <- as.matrix(dtm)
# Sparsity
sum(mm == 0)/(3*12)
```

```
## [1] 0.5277778
```

Problems

- Guitar; Guitars, Guitar.
- One way would be to remove all punctuations before creating corpus, using regex

```
text2 <- str_remove_all(text, pattern = "[:punct:]")  
text2
```

```
## [1] "The best guitar models are Gibson and Fender"  
## [2] "My guitars are the best guitars"  
## [3] "Are these guitars yours"
```


- Or alternatively use `tm_map()`.
- `tm_map()` takes corpus and applies on it a function, in this case `removePunctuation`

```
corpus <- VCorpus(VectorSource(text))
corpus <- tm_map(corpus, removePunctuation)
dtm <- DocumentTermMatrix(corpus)
inspect(dtm)
```

```
## <<DocumentTermMatrix (documents: 3, terms: 11)>>
## Non-/sparse entries: 16/17
## Sparsity           : 52%
## Maximal term length: 7
## Weighting          : term frequency (tf)
## Sample            :
##      Terms
## Docs and are best fender gibson guitar guitars models the these
##   1  1  1  1  1  1  1  0  1  1  0
##   2  0  1  1  0  0  0  2  0  1  0
##   3  0  1  0  0  0  0  1  0  0  1
```

- Term frequency-inverse document frequency - numerical statistic that is intended to reflect how important a word is to a document in a collection or corpus
- The tf-idf value increases proportionally to the number of times a word appears in the document and is offset by the frequency of the word in the corpus, which helps to adjust for the fact that some words appear more frequently in general.
- TF - Term Frequency, which measures how frequently a term occurs in a document
- $TF(t) = \frac{\text{Number of times term } t \text{ appears in a document}}{\text{Total number of terms in the document}}$
- Since every document is different in length, it is possible that a term would appear much more times in long documents than shorter ones. Thus, the term frequency is often divided by the document length (the total number of terms in the document) as a way of normalization.

- Tf – the importance of the term in the separate document.
- Inverse Document Frequency measures how important a term is.
- While computing TF, all terms are considered equally important. However it is known that certain terms, such as “is”, “of”, and “that”, may appear a lot of times but have little importance. Thus we need to weigh down the frequent terms while scale up the rare ones, by computing the following:
 - $IDF(t) = \log_2\left(\frac{\text{Total number of documents}}{\text{Number of documents with term in it}}\right)$
 - IDF decrease the weight of frequent words.
 - Could you say which words have the high tf-idf weight?
 - The smaller the weight, the more common the term.

```
dtm1 <- DocumentTermMatrix(corpus, control = list(weighting = weightTfIdf))
as.matrix(dtm1)
```

```
##      Terms
## Docs      and are      best      fender      gibson      guitar      guitars
##   1 0.1981203    0 0.07312031 0.1981203 0.1981203 0.1981203 0.0000000
##   2 0.0000000    0 0.11699250 0.0000000 0.0000000 0.0000000 0.2339850
##   3 0.0000000    0 0.00000000 0.0000000 0.0000000 0.0000000 0.1462406
##      Terms
## Docs      models      the      these      yours
##   1 0.1981203 0.07312031 0.0000000 0.0000000
##   2 0.0000000 0.11699250 0.0000000 0.0000000
##   3 0.0000000 0.00000000 0.3962406 0.3962406
```

- for guitar

```
(1/8)*log(3/1, base=2)
```

```
## [1] 0.1981203
```

300,000 songs and their lyrics

- Let's consider the example from the file lyrics.csv

```
lyrics <- read.csv("lyrics.csv", stringsAsFactors = F)
summary(lyrics$year)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      67     2006     2008     2009     2014     2038
```

Data Cleaning

```
lyrics %>%  
  group_by(year) %>%  
  summarise(count = n()) %>%  
  arrange(desc(year))
```

```
## # A tibble: 52 x 2  
##   year count  
##   <int> <int>  
## 1  2038    10  
## 2  2016 35042  
## 3  2015 27159  
## 4  2014 26393  
## 5  2013 16245  
## 6  2012 14581  
## 7  2011 11387  
## 8  2010 11374  
## 9  2009 11333  
## 10 2008 20375  
## # ... with 42 more rows
```

Data Cleaning

```
lyrics %>%  
  group_by(year) %>%  
  summarise(count = n()) %>%  
  arrange(year)
```

```
## # A tibble: 52 x 2  
##   year count  
##   <int> <int>  
## 1     67     1  
## 2    112     4  
## 3    702     1  
## 4   1968     1  
## 5   1970   421  
## 6   1971   480  
## 7   1972   496  
## 8   1973   504  
## 9   1974   395  
## 10  1975   321  
## # ... with 42 more rows
```

Data Cleaning

```
lyrics %>%  
  group_by(artist) %>%  
  summarise(n_songs=n()) %>%  
  arrange(desc(n_songs))
```

```
## # A tibble: 17,088 x 2  
##   artist          n_songs  
##   <chr>          <int>  
## 1 dolly-parton      755  
## 2 american-idol    700  
## 3 elton-john       680  
## 4 b-b-king         667  
## 5 chris-brown      655  
## 6 eddy-arnold      628  
## 7 barbra-streisand 624  
## 8 ella-fitzgerald  623  
## 9 bob-dylan        614  
## 10 bee-gees        599  
## # ... with 17,078 more rows
```


Data Cleaning

```
lyrics$char_num <- nchar(lyrics$lyrics)
summary(lyrics$char_num)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      0.0     0.0   744.0   844.1  1212.0 42155.0
```

```
table(lyrics$genre)
```

```
##
##      Country      Electronic      Folk      Hip-Hop      Indie
##      16449      15446      3095      30644      5373
##      Jazz      Metal Not Available      Other      Pop
##      15659      26900      27984      22394      46371
##      R&B      Rock
##      5560      123402
```

```
sum(lyrics$char_num==0)
```

```
## [1] 89449
```

Data Cleaning

```
lyrics <- lyrics %>%  
  filter(char_num>100 & genre != "Not Available" & year > 1968 & year <= 20  
dim(lyrics)
```

```
## [1] 221669      7
```

```
table(lyrics$genre)
```

```
##  
##      Country Electronic      Folk      Hip-Hop      Indie      Jazz  
##      13710      7027      2006      22168      2888      6915  
##      Metal      Other      Pop      R&B      Rock  
##      21250      4838      37770      3198      99899
```

```
summary(lyrics$char_num)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.  
##      101      681      986     1184     1419    42155
```

Data Cleaning

```
b <- "[Intro: BeyoncÃfÂ©]"
iconv(b, to="ASCII", sub="")
```

```
## [1] "[Intro: Beyonc]"
```

```
lyrics$lyrics <- iconv(lyrics$lyrics, to="ASCII", sub="")
```

Let's analyze lyrics of Beyonce

```
beyonce <- lyrics[lyrics$artist == 'beyonce-knowles',]
beyonce_vs <- VectorSource(beyonce$lyrics)
beyonce_corpus <- VCorpus(beyonce_vs)
```

Stopwords

```
stopwords("english")
```

##	[1]	"i"	"me"	"my"	"myself"	"we"
##	[6]	"our"	"ours"	"ourselves"	"you"	"your"
##	[11]	"yours"	"yourself"	"yourselves"	"he"	"him"
##	[16]	"his"	"himself"	"she"	"her"	"hers"
##	[21]	"herself"	"it"	"its"	"itself"	"they"
##	[26]	"them"	"their"	"theirs"	"themselves"	"what"
##	[31]	"which"	"who"	"whom"	"this"	"that"
##	[36]	"these"	"those"	"am"	"is"	"are"
##	[41]	"was"	"were"	"be"	"been"	"being"
##	[46]	"have"	"has"	"had"	"having"	"do"
##	[51]	"does"	"did"	"doing"	"would"	"should"
##	[56]	"could"	"ought"	"i'm"	"you're"	"he's"
##	[61]	"she's"	"it's"	"we're"	"they're"	"i've"
##	[66]	"you've"	"we've"	"they've"	"i'd"	"you'd"
##	[71]	"he'd"	"she'd"	"we'd"	"they'd"	"i'll"
##	[76]	"you'll"	"he'll"	"she'll"	"we'll"	"they'll"
##	[81]	"isn't"	"aren't"	"wasn't"	"weren't"	"hasn't"
##	[86]	"haven't"	"hadn't"	"doesn't"	"don't"	"didn't"
##	[91]	"won't"	"shan't"	"can't"	"couldn't"	"mustn't"

Stemming

```
love <- c("love", "loving", "lovingly", "loved", "lover", "lovely", "love",  
         "game", "gaming", "games")  
stemDocument(love)
```

```
## [1] "love"  "love"  "love"  "love"  "lover" "love"  "love"  "game"  
## [9] "game"  "game"
```

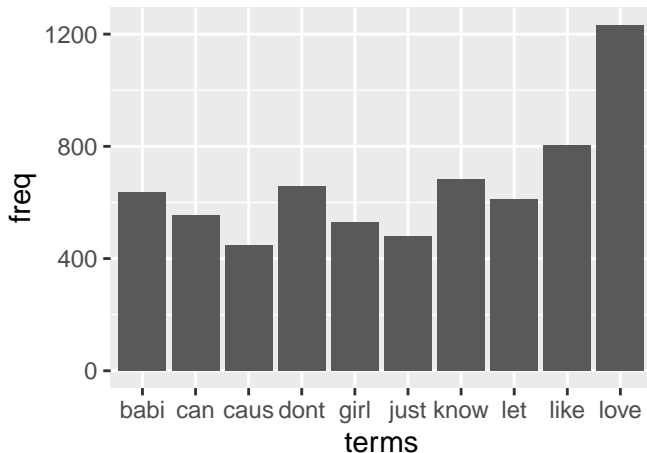
Analysis

```
beyonce_tdm <- TermDocumentMatrix(beyonce_corpus,  
                                   control=list(removeNumbers = T,removePunctuation = T,  
                                                stopwords = T, stemming = T))  
  
dtm_mat <- as.matrix(beyonce_tdm)  
freqs <- rowSums(dtm_mat)  
df_freq <- data.frame(terms=rownames(dtm_mat),  
                      freq = freqs, stringsAsFactors = F)  
df_freq <- df_freq[order(df_freq$freq, decreasing = T),]  
head(df_freq)
```

```
##      terms freq  
## love  love 1233  
## like  like  803  
## know  know  684  
## dont  dont  659  
## babi  babi  638  
## let   let   610
```

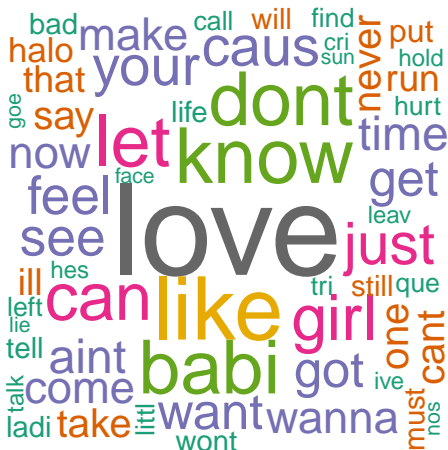
Analysis

```
df_top10 <- df_freq[1:10,]  
ggplot(df_top10, aes(x=terms, y=freq)) + geom_bar(stat = 'identity')
```



Wordcloud

```
set.seed(1)
wordcloud(words = df_freq$terms, freq = df_freq$freq, min.freq = 10,
max.words = 200, random.order = FALSE, colors=brewer.pal(8, "Dark2"))
```



Tf-Idf

```
beyonce_tdm <- TermDocumentMatrix(beyonce_corpus,  
  control = list(removeNumbers = T, removePunctuation = T,  
    stopwords = T, stemming = T, weighting = weightTfIdf))  
tdm_mat <- as.matrix(beyonce_tdm)  
freqs <- rowSums(tdm_mat)  
df_freq <- data.frame(terms=rownames(tdm_mat),  
  freq = freqs)  
df_freq <- df_freq[order(df_freq$freq, decreasing = T),]  
head(df_freq)
```

```
##      terms      freq  
## halo  halo 5.295742  
## love  love 4.083825  
## girl  girl 3.326538  
## run   run  3.293408  
## let   let  3.269829  
## babi  babi 2.828966
```

Tf-Idf

```
set.seed(27)
wordcloud(words = df_freq$terms, freq = df_freq$freq, min.freq = 1,
          max.words = 90, random.order = FALSE, colors = brewer.pal(8, "Dark2"))
```



Task

- Extract the phone numbers in one column and the names in another:
- Combine them in the data frame

```
phones <- c("Joe 027-789663",  
            "Jimi 99656565",  
            "Adri2 099-65-1995 GUITARIS")
```