

Homework 03 Data Cleaning, Normal Distributions

Due by 11:59pm, Saturday, February 11, 2023

S&DS 230

(1) More on List Manipulation (18 points - 3 points each).

(1.1) Make an object called `myList` that contains the following elements (in order):

- The integers 1 through 10
- A matrix with the integers 1 through 25 that has five rows, filled by row
- A list that contains
 - The text “latte”
 - A vector with the text “taco” and “nan”
 - A vector with the integers 1 through 7

You should be able to make this object in a single line of code.

Use `[]`, `[[]]`, `[,]` notation to answer parts b) through f).

(1.2) Make an object called `ans1` that is the fourth row of the matrix contained in `myList`.

(1.3) Make an object called `ans2` that is the sum of the 5th column of the matrix contained in `myList`.

(1.4) Make an object called `ans3` that is the sum of EACH row of the matrix contained in `myList` (use the `apply()` function or check out `rowSums()`).

(1.5) Make an object called `ans4` that is whichever single element of `myList` that you’d like to consume (yes, comedians, it has to be a food ...).

(1.6) Make an object called `ans5` that is the third element of the third element of `myList` converted to characters.

Get the results of each of your objects you created above (i.e. get them to show up in your knitted file by typing their names or putting the code line that creates each object in parentheses).

```
myList <- list(c(1:10), matrix(c(1:25), nrow = 5, byrow = TRUE),
              list(str = "latte", vec = c("taco", "nan"), vec = c(1:7)))
myList
## [[1]]
## [1] 1 2 3 4 5 6 7 8 9 10
##
```

```
## [[2]]
##      [,1] [,2] [,3] [,4] [,5]
## [1,]    1    2    3    4    5
## [2,]    6    7    8    9   10
## [3,]   11   12   13   14   15
## [4,]   16   17   18   19   20
## [5,]   21   22   23   24   25
##
## [[3]]
## [[3]]$str
## [1] "latte"
##
## [[3]]$vec
## [1] "taco" "nan"
##
## [[3]]$vec
## [1] 1 2 3 4 5 6 7

ans1 <- myList[[2]][4, ]
ans1

## [1] 16 17 18 19 20

ans2 <- sum(myList[[2]][ ,5])
ans2

## [1] 75

ans3 <- rowSums(myList[[2]][c(1,2,3,4,5), ])
ans3

## [1] 15 40 65 90 115

ans4 <- myList[[3]][[1]]
ans4

## [1] "latte"

ans5 <- as.character(myList[[3]][[3]])
ans5

## [1] "1" "2" "3" "4" "5" "6" "7"
```

(2) Normal Quantile Plots and the Binomial Distribution (20 points, 3 points each, part (2.5) is 5 points).

You may recall from your Intro Statistics course that a binomial distribution looks like a normal distribution if $np > 10$ and $n(1-p) > 10$ (i.e. as long as the average number of successes and failures are both larger than 10). Recall that n is the number of trials, and p is the probability of success for each Bernoulli trial. *As an example, flip a coin 30 times, count the number of heads. $n=30$, $p=.5$, $np = 15 > 10$ and $n(1-p) = 15 > 10$, so the distribution should be approximately normal).*

You are going to make six normal quantile plots that simulate 100 random observations from binomial distributions with $p = .2$ and various values of n .

(2.1) Install the `car` package. This will allow you use the `qqPlot()` function. Load this package.

(2.2) Make a vector called `vec` that is powers of 10 for powers 0 through 5. The one caveat, is that you need to use the `**` operator which reads as 'to the power of' (i.e. $2^{**}3$ is 8).

(2.3) Use the `par()` function to set up your plot region to show 6 plots on a page. The `par` argument you want is `mfrow = c(2,3)` which sets your plot region to have 2 rows and 3 columns.

(2.4) Use the `rbinom()` function to generate 5 random binomial observations, each with 20 trials, and with $p=0.8$. You may need to type `?rbinom` to get the syntax for this function. Store the result in an object called `vec2`.

(2.5) Write a loop that repeatedly creates a normal quantile plot for 100 random samples each from a binomial distribution with $p=0.2$ and n equal to the 6 values stored in `vec`. A few plot details : * Use the `qqPlot` function. * Make the graph points red solid dots (`pch = 19`). * Make the boundary lines blue (use `col.lines`) * Make a main graph title that pastes the text "100 Binomial Samples, N =" to the corresponding value from `vec`.

```
library(car)

## Loading required package: carData

vec <- 10 ** c(0:5)
vec

## [1] 1e+00 1e+01 1e+02 1e+03 1e+04 1e+05

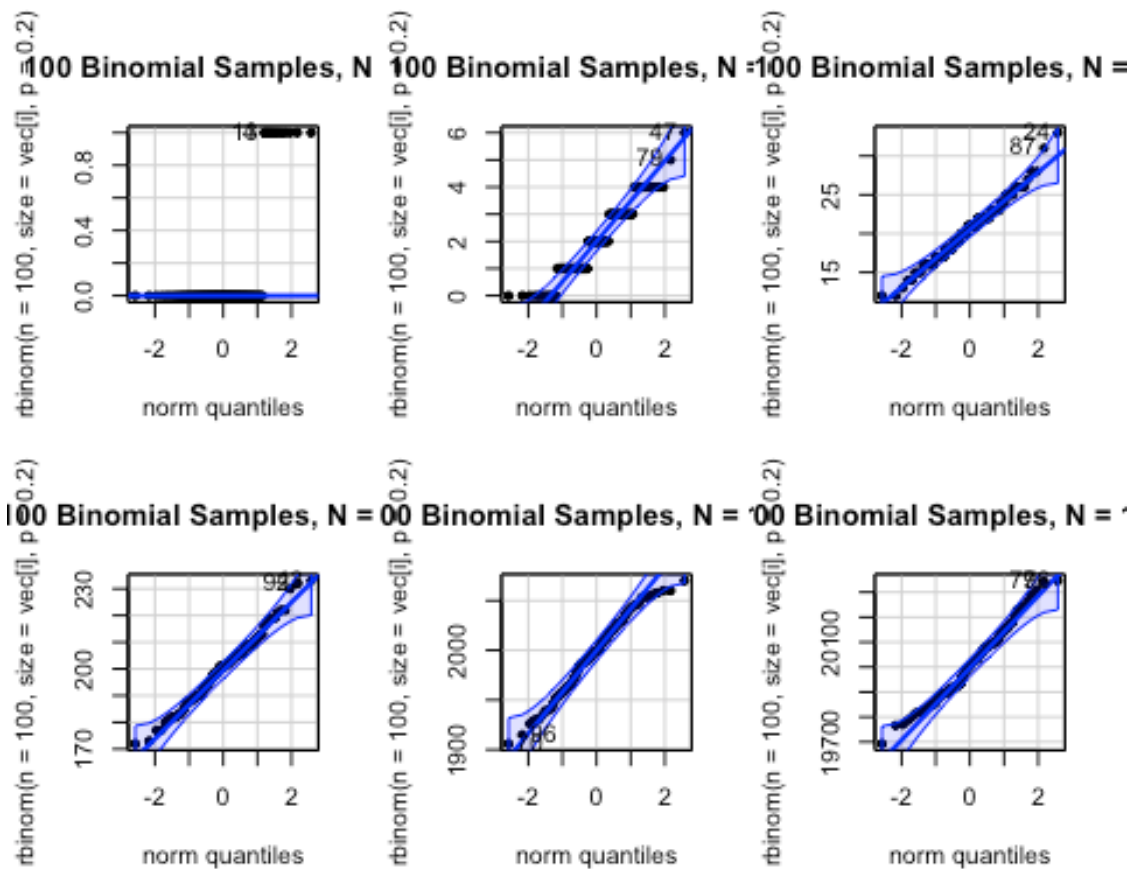
par(mfrow = c(2, 3))

vec2 <- rbinom(n = 5, size = 20, prob = 0.8)
vec2

## [1] 16 16 19 17 18

#n : 100 random samples
#size = number of trial
#prob = prob of success on each trial

for (i in 1:length(vec)) {
  qqPlot(rbinom(n = 100, size = vec[i], p = 0.2), pch = 19, col.lines =
'blue',
        main = paste(c("100 Binomial Samples, N =", vec[i]), collapse = "
"))
}
```



(2.6) Take a look at the normal quantile plots. For what value of n do the graphs seem to be approximately normally distributed? Is this consistent with what you expect? Write two complete sentences to answer these questions.

N values above 1000 seems to be approximately normally distributed. Yes, it is consistent with what was expected. According to the Central Limit Theorem, the sample mean of any distribution will become approximately normal if the sample size is sufficiently large. After the 3rd plot we can see a straight line which depicts normal distribution on qqPlot.

(3) Favorite food and Data Cleaning (62 points. Parts 3.2 through 3.5, 2 pts each, other values listed below).

This is data generated by you - thank you! I simply asked "What is your favorite food?". You can get the data [HERE](#).

Your goal is similar to what we did with the question "What animal would you like to be?" in Class 5 : clean this variable, make a barplot, and discuss the results.

(3.1) (1 pt) Read in the data to a new object called food.

(3.2) Get a sense of the dataset - dimensions, variable names, look at the first few rows.

(3.3) Convert food to a single vector that is just the first column (literally, replace food with food\$Food).

(3.4) Show the sorted unique values of food. Calculate how many unique values exist in food.

```
food <- read.csv("http://reuningscherer.net/S&DS230/data/food_230.csv")

dim(food)
## [1] 317  1

names(food)
## [1] "Food"

head(food)

##           Food
## 1 french fries.
## 2         Chicken
## 3 butter chicken
## 4         Ginger
## 5         Dal Bhaat
## 6         Thai food

tail(food)

##           Food
## 312          ham
## 313 I love Korean cuisine.
## 314       Fried chicken
## 315          steak
## 316          Steak
## 317         chicken

food <- food$Food

sort(unique(food))

## [1] "Albanian Food "
## [2] "all kinds of delicious food"
## [3] "amaretto dark chocolate"
## [4] "any type of cheese"
## [5] "Arepa"
## [6] "Artichoke"
## [7] "Asian cuisine"
## [8] "asian food"
## [9] "Asian food"
## [10] "Bagel"
## [11] "baguettes"
## [12] "Bananas"
## [13] "Blue Point Oyster"
## [14] "brazilian chorizo pizzas"
## [15] "Brazilian food."
```

```
## [16] "Bread"
## [17] "Burgers"
## [18] "burritos"
## [19] "Burritos"
## [20] "butter chicken"
## [21] "Cajun Fries"
## [22] "cake"
## [23] "Ceviche"
## [24] "cheese"
## [25] "Cheese"
## [26] "Cheeseburgers"
## [27] "cheez its"
## [28] "chicken"
## [29] "Chicken"
## [30] "Chicken Malai Kabab"
## [31] "Chicken parmesan and penne alla vodka"
## [32] "chicken tenders"
## [33] "Chicken Tikka and Naan"
## [34] "Chicken Tikka Masala"
## [35] "Chicken wings"
## [36] "chinese"
## [37] "Chinese food"
## [38] "Chinese Food"
## [39] "Chinese food "
## [40] "Chinese food is my favorite."
## [41] "chinese hotpot"
## [42] "Chipotle"
## [43] "chocolate"
## [44] "Chocolate\n\n \n\nA cat"
## [45] "chocolate chip cookies"
## [46] "Cinnamon Rolls"
## [47] "comfort food"
## [48] "Cookies"
## [49] "Corn"
## [50] "Cottage Pie"
## [51] "Crepes"
## [52] "Curries of all sorts"
## [53] "Curry"
## [54] "Curry vindaloo"
## [55] "Dal Bhaat"
## [56] "Dandan noodles"
## [57] "Delicious "
## [58] "Dessert"
## [59] "Donuts!"
## [60] "Dried mangos"
## [61] "dumplings"
## [62] "Empanadas"
## [63] "enchiladas"
## [64] "Escargot in garlic butter"
## [65] "Ethiopian Cuisine"
```

```

## [66] "farofa"
## [67] "Farofa "
## [68] "Fish and chips"
## [69] "Fish tacos"
## [70] "Five guys' fries"
## [71] "Flautas"
## [72] "Freeze Dried"
## [73] "french fries"
## [74] "french fries."
## [75] "French onion soup"
## [76] "fried chicken"
## [77] "Fried chicken"
## [78] "Fried fish"
## [79] "Fried okra"
## [80] "Fried Rice"
## [81] "fruit"
## [82] "Fruit"
## [83] "Fruits"
## [84] "Ginger"
## [85] "Gizzard (chicken) and vegetables (greens), Igbo dish"
## [86] "gnocchi"
## [87] "Good pizza"
## [88] "Grilled chicken breast"
## [89] "guacamole"
## [90] "guacamole "
## [91] "Gyros"
## [92] "ham"
## [93] "Hamburgers"
## [94] "Hibachi"
## [95] "Hot dog"
## [96] "hot pot"
## [97] "Hot pot"
## [98] "hotpot"
## [99] "I love Asian food, especially Chinese food, Thai food, and
Sushi. \n\n "
## [100] "I love Korean cuisine."
## [101] "i love mexican food"
## [102] "ice cream"
## [103] "Ice cream"
## [104] "Ice Cream"
## [105] "ICE CREAM"
## [106] "Ice cream!!!"
## [107] "Ice cream. "
## [108] "indian"
## [109] "Indian"
## [110] "Indian food"
## [111] "Indian Food"
## [112] "Indian Food "
## [113] "Indian food is great."
## [114] "Indonesian food"

```

```
## [115] "Instant ramen"
## [116] "italian"
## [117] "Italian"
## [118] "italian food"
## [119] "Italian food"
## [120] "Italian Food"
## [121] "italian food - pasta"
## [122] "Italian food!"
## [123] "Japanese "
## [124] "Japanese food"
## [125] "Jollof Rice and Chicken"
## [126] "Jollof Rice with chicken "
## [127] "Kelewele, it's an African dish"
## [128] "Korean"
## [129] "Korean Barbeque"
## [130] "Korean BBQ"
## [131] "korean bbq!"
## [132] "korean bbq!!"
## [133] "korean food"
## [134] "Korean food"
## [135] "Korean Food"
## [136] "korean food "
## [137] "kosher Steak "
## [138] "lahmacun (turkish pizza)"
## [139] "Lasagna"
## [140] "Lasagna."
## [141] "Lasagne"
## [142] "Lebanese Food"
## [143] "mac and cheese"
## [144] "Macaroni and Cheese"
## [145] "mango"
## [146] "Meat"
## [147] "mediterranean "
## [148] "Mediterranean "
## [149] "Mediterranean food "
## [150] "Mexican"
## [151] "Mexican food"
## [152] "Mexican Food"
## [153] "Multigrain pancake"
## [154] "My favorite kind of food is pastries."
## [155] "nectarines"
## [156] "noodles"
## [157] "Noodles"
## [158] "noodles with broth"
## [159] "Noodles."
## [160] "palak paneer"
## [161] "pasta"
## [162] "Pasta"
## [163] "Patty Melt"
## [164] "Persian and Mediterranean"
```



```
## [165] "Pho"
## [166] "pizza"
## [167] "Pizza"
## [168] "Pizza!!!"
## [169] "platanos"
## [170] "Poke Bowl"
## [171] "Postickers"
## [172] "probably either casual diner fare, Italian, or shellfish"
## [173] "Puerto Rican food"
## [174] "ramen"
## [175] "Ramen"
## [176] "ramen\n\n "
## [177] "Ribz"
## [178] "rice with curry"
## [179] "Roast dinner "
## [180] "salmon"
## [181] "Salmon"
## [182] "salty"
## [183] "seafood"
## [184] "shahi paneer"
## [185] "sharp cheddar cheese"
## [186] "Shrimp curry"
## [187] "soup"
## [188] "Soup"
## [189] "South Asian Rice Dishes for example Biryani, Pulao, Mandi, Tahri "
## [190] "Spaghetti "
## [191] "Spicy"
## [192] "Spicy and flavorful food"
## [193] "spicy food "
## [194] "Spicy Hotpot"
## [195] "spicy tofu"
## [196] "steak"
## [197] "Steak"
## [198] "strawberries"
## [199] "Strawberries "
## [200] "sundried tomatoes"
## [201] "sushi"
## [202] "Sushi"
## [203] "Sushi with an overwhelming amount of raw salmon"
## [204] "Swedish meatballs"
## [205] "Sweet chili Doritos"
## [206] "sweet potato"
## [207] "Tagine"
## [208] "Thai"
## [209] "thai food"
## [210] "Thai food"
## [211] "Thai Food"
## [212] "Thai food "
## [213] "Thai food is my favorite kind of food "
## [214] "Thai specifically beef pad see-ew and some rice."
```

```
## [215] "Tofu"
## [216] "udon noodle"
## [217] "Vietnamese food"
## [218] "Vietnamese spring rolls"

length(unique(food))

## [1] 218
```

(3.5) Write a couple of sentences about what data cleaning issues you notice among the unique values of food.

In the dataset there are some multiple responses out of which we will have to select one option of the favorite food. There are few food items which are misspelled which would also need some cleaning. Some food items are written as plurals/singulars. A lot of food items have additional words which can be cleaned. There are unique food items which can be combined together in a specific cuisine. Additionally, I noticed that there are different varieties of chicken dishes which can be clubbed together. Some foods are same but they are not unique due to the issues of lowercase and uppercase. Some food names have stopwords and regular expressions which can be cleaned to get a decent number of food items to analyse. Some examples include: "Thai food" written in different ways, different types of curries, sushi written in different ways.

Before proceeding to data cleaning, a quick reminder example of how to remove text before or after a particular word using the regular expression `.*` Note that `.` stands for 'any character' (other than a new line), and `*` stands for '0 or more times'.

Example: Find " have ", delete this AND everything preceding or following:

```
#Deletes " have " and everything preceding
gsub(".* have ", "", "Cats have personality")

#Deletes " have" and everything following
gsub(" have .*", "", "Cats have personality")
```

(3.6) Cleaning Part I (8 pts): Clean the data using the following steps (in order):

- Convert data to lower case
- Find " or " and remove this and anything that follows.
- Find " and " and remove this and anything that follows.
- Find " food" and remove this AND anything that follows.
- Find " cuisine" and remove this AND anything that follows.
- Remove all special characters and punctuation.
- Remove trailing spaces at the end of text (use the `trimws()` function)

At each step, you'll probably want to check what unique values of food are left to make sure your functions are working correctly. By the time you finish, you should have 156 unique levels.

Your final two lines of code should again show the sorted unique values of food and the current number of unique values.

```
food <- tolower(food)
food <- gsub(" or .*", "", food)
food <- gsub(" and .*", "", food)
food <- gsub(" food.*", "", food)
food <- gsub(" cuisine.*", "", food)
food <- gsub("[^0-9A-Za-z//' ]", "", food)
food <- gsub("'", "", food)
food <- gsub("[[:punct:]]", "", food)
food <- trimws(food, which = "right", whitespace = "[ \\t\\r\\n]")

sort(unique(food))

## [1] "albanian"
## [2] "all kinds of delicious"
## [3] "amaretto dark chocolate"
## [4] "any type of cheese"
## [5] "arepa"
## [6] "artichoke"
## [7] "asian"
## [8] "bagel"
## [9] "baguettes"
## [10] "bananas"
## [11] "blue point oyster"
## [12] "brazilian"
## [13] "brazilian chorizo pizzas"
## [14] "bread"
## [15] "burgers"
## [16] "burritos"
## [17] "butter chicken"
## [18] "cajun fries"
## [19] "cake"
## [20] "ceviche"
## [21] "cheese"
## [22] "cheeseburgers"
## [23] "cheez its"
## [24] "chicken"
## [25] "chicken malai kabab"
## [26] "chicken parmesan"
## [27] "chicken tenders"
## [28] "chicken tikka"
## [29] "chicken tikka masala"
## [30] "chicken wings"
## [31] "chinese"
## [32] "chinese hotpot"
## [33] "chipotle"
## [34] "chocolate"
## [35] "chocolate chip cookies"
```

```
## [36] "chocolatea cat"
## [37] "cinnamon rolls"
## [38] "comfort"
## [39] "cookies"
## [40] "corn"
## [41] "cottage pie"
## [42] "crepes"
## [43] "curries of all sorts"
## [44] "curry"
## [45] "curry vindaloo"
## [46] "dal bhaat"
## [47] "dandan noodles"
## [48] "delicious"
## [49] "dessert"
## [50] "donuts"
## [51] "dried mangos"
## [52] "dumplings"
## [53] "empanadas"
## [54] "enchiladas"
## [55] "escargot in garlic butter"
## [56] "ethiopian"
## [57] "farofa"
## [58] "fish"
## [59] "fish tacos"
## [60] "five guys fries"
## [61] "flautas"
## [62] "freeze dried"
## [63] "french fries"
## [64] "french onion soup"
## [65] "fried chicken"
## [66] "fried fish"
## [67] "fried okra"
## [68] "fried rice"
## [69] "fruit"
## [70] "fruits"
## [71] "ginger"
## [72] "gizzard chicken"
## [73] "gnocchi"
## [74] "good pizza"
## [75] "grilled chicken breast"
## [76] "guacamole"
## [77] "gyros"
## [78] "ham"
## [79] "hamburgers"
## [80] "hibachi"
## [81] "hot dog"
## [82] "hot pot"
## [83] "hotpot"
## [84] "i love asian"
## [85] "i love korean"
```

```
## [86] "i love mexican"
## [87] "ice cream"
## [88] "indian"
## [89] "indonesian"
## [90] "instant ramen"
## [91] "italian"
## [92] "japanese"
## [93] "jollof rice"
## [94] "jollof rice with chicken"
## [95] "kelewele its an african dish"
## [96] "korean"
## [97] "korean barbeque"
## [98] "korean bbq"
## [99] "kosher steak"
## [100] "lahmacun turkish pizza"
## [101] "lasagna"
## [102] "lasagne"
## [103] "lebanese"
## [104] "mac"
## [105] "macaroni"
## [106] "mango"
## [107] "meat"
## [108] "mediterranean"
## [109] "mexican"
## [110] "multigrain pancake"
## [111] "my favorite kind of"
## [112] "nectarines"
## [113] "noodles"
## [114] "noodles with broth"
## [115] "palak paneer"
## [116] "pasta"
## [117] "patty melt"
## [118] "persian"
## [119] "pho"
## [120] "pizza"
## [121] "platanos"
## [122] "poke bowl"
## [123] "postickers"
## [124] "probably either casual diner fare italian"
## [125] "puerto rican"
## [126] "ramen"
## [127] "ribs"
## [128] "rice with curry"
## [129] "roast dinner"
## [130] "salmon"
## [131] "salty"
## [132] "seafood"
## [133] "shahi paneer"
## [134] "sharp cheddar cheese"
## [135] "shrimp curry"
```

```
## [136] "soup"
## [137] "south asian rice dishes for example biryani pulao mandi tahri"
## [138] "spaghetti"
## [139] "spicy"
## [140] "spicy hotpot"
## [141] "spicy tofu"
## [142] "steak"
## [143] "strawberries"
## [144] "sundried tomatoes"
## [145] "sushi"
## [146] "sushi with an overwhelming amount of raw salmon"
## [147] "swedish meatballs"
## [148] "sweet chili doritos"
## [149] "sweet potato"
## [150] "tagine"
## [151] "thai"
## [152] "thai specifically beef pad seeew"
## [153] "tofu"
## [154] "udon noodle"
## [155] "vietnamese"
## [156] "vietnamese spring rolls"

length(unique(food))

## [1] 156
```

(3.7) Cleaning Part II (10 pts): A few quick random cleaning items:

Clean up the following types of food (in order) - one line of code per type of food. In each case, deal with misspellings, modifiers ("shrimp curry" vs just "curry"), two words ('hot pot' instead of 'hotpot'), plurals, etc.

- hotpot
- curry
- lasagna
- noodles
- cookies
- chocolate
- cheese
- steak
- sushi
- fries (french, cajun, five guys' all just call 'fries')
- ramen
- tofu
- burgers (of any kind)
- soup
- anything containing 'delicious' just call 'delicious'

When you're finished, you should have 130 unique values.

Your final two lines of code should again show the sorted unique values of food and the current number of unique values.

```
food <- gsub("hot pot|chinese hotpot|spicy hotpot", "hotpot", food)
food <- gsub(".*curr.*", "curry", food)
food <- gsub("lasagne", "lasagna", food)
food <- gsub("udon noodle|noodles with broth|dandan noodles", "noodles", food)
food <- gsub(".*cookies.*", "cookies", food)
food <- gsub(".*chocolat.*", "chocolate", food)
food <- gsub(".*cheese.*", "cheese", food)
food <- gsub(".*steak.*", "steak", food)
food <- gsub("french fries|cajun fries|five guys fries", "fries", food)
food <- gsub(".*sushi.*", "sushi", food)
food <- gsub(".*ramen.*", "ramen", food)
food <- gsub(".*tofu.*", "tofu", food)
food <- gsub(".*burger.*", "burger", food)
food <- gsub(".*soup.*", "soup", food)
food <- gsub(".*delicious.*", "delicious", food)
```

```
sort(unique(food))
```

```
## [1] "albanian"
## [2] "arepa"
## [3] "artichoke"
## [4] "asian"
## [5] "bagel"
## [6] "baguettes"
## [7] "bananas"
## [8] "blue point oyster"
## [9] "brazilian"
## [10] "brazilian chorizo pizzas"
## [11] "bread"
## [12] "burger"
## [13] "burritos"
## [14] "butter chicken"
## [15] "cake"
## [16] "ceviche"
## [17] "cheese"
## [18] "cheez its"
## [19] "chicken"
## [20] "chicken malai kabab"
## [21] "chicken parmesan"
## [22] "chicken tenders"
## [23] "chicken tikka"
## [24] "chicken tikka masala"
## [25] "chicken wings"
## [26] "chinese"
```

```
## [27] "chipotle"
## [28] "chocolate"
## [29] "cinnamon rolls"
## [30] "comfort"
## [31] "cookies"
## [32] "corn"
## [33] "cottage pie"
## [34] "crepes"
## [35] "curry"
## [36] "dal bhaat"
## [37] "delicious"
## [38] "dessert"
## [39] "donuts"
## [40] "dried mangos"
## [41] "dumplings"
## [42] "empanadas"
## [43] "enchiladas"
## [44] "escargot in garlic butter"
## [45] "ethiopian"
## [46] "farofa"
## [47] "fish"
## [48] "fish tacos"
## [49] "flautas"
## [50] "freeze dried"
## [51] "fried chicken"
## [52] "fried fish"
## [53] "fried okra"
## [54] "fried rice"
## [55] "fries"
## [56] "fruit"
## [57] "fruits"
## [58] "ginger"
## [59] "gizzard chicken"
## [60] "gnocchi"
## [61] "good pizza"
## [62] "grilled chicken breast"
## [63] "guacamole"
## [64] "gyros"
## [65] "ham"
## [66] "hibachi"
## [67] "hot dog"
## [68] "hotpot"
## [69] "i love asian"
## [70] "i love korean"
## [71] "i love mexican"
## [72] "ice cream"
## [73] "indian"
## [74] "indonesian"
## [75] "italian"
## [76] "japanese"
```



```
## [77] "jollof rice"
## [78] "jollof rice with chicken"
## [79] "kelewele its an african dish"
## [80] "korean"
## [81] "korean barbeque"
## [82] "korean bbq"
## [83] "lahmacun turkish pizza"
## [84] "lasagna"
## [85] "lebanese"
## [86] "mac"
## [87] "macaroni"
## [88] "mango"
## [89] "meat"
## [90] "mediterranean"
## [91] "mexican"
## [92] "multigrain pancake"
## [93] "my favorite kind of"
## [94] "nectarines"
## [95] "noodles"
## [96] "palak paneer"
## [97] "pasta"
## [98] "patty melt"
## [99] "persian"
## [100] "pho"
## [101] "pizza"
## [102] "platanos"
## [103] "poke bowl"
## [104] "postickers"
## [105] "probably either casual diner fare italian"
## [106] "puerto rican"
## [107] "ramen"
## [108] "ribs"
## [109] "roast dinner"
## [110] "salmon"
## [111] "salty"
## [112] "seafood"
## [113] "shahi paneer"
## [114] "soup"
## [115] "south asian rice dishes for example biryani pulao mandi tahri"
## [116] "spaghetti"
## [117] "spicy"
## [118] "steak"
## [119] "strawberries"
## [120] "sundried tomatoes"
## [121] "sushi"
## [122] "swedish meatballs"
## [123] "sweet chili doritos"
## [124] "sweet potato"
## [125] "tagine"
## [126] "thai"
```

```
## [127] "thai specifically beef pad seeew"
## [128] "tofu"
## [129] "vietnamese"
## [130] "vietnamese spring rolls"

length(unique(food))

## [1] 130
```

(3.8) Cleaning Part III (8 pts): Cleaning types of cuisine.

Clean up the following types of cuisine (in order) - in this case, you'll want to make a vector called `searchvec` that contains the types of cuisine. Then create a loop following the example in Class 5 to replace all the modifiers for each cuisine type so that you ultimately end up with cleaned up versions of each cuisine type. Use not more than 5 lines of code.

The cuisine types (in order) are * asian * chinese * vietnamese * italian * indian * thai * mexican * brazilian * korean

(there are other types of cuisine, but they don't require cleaning).

When you're finished, you should have 120 unique values.

Your final two lines of code should again show the sorted unique values of food and the current number of unique values.

```
searchvec <- c("asian", "chinese", "vietnamese", "italian", "indian", "thai",
               "mexican", "brazilian", "korean")

for (i in 1:length(searchvec)){
  food <- gsub(paste0(".*", searchvec[i], ".*"), searchvec[i], food)
}

sort(unique(food))

## [1] "albanian"
## [3] "artichoke"
## [5] "bagel"
## [7] "bananas"
## [9] "brazilian"
## [11] "burger"
## [13] "butter chicken"
## [15] "ceviche"
## [17] "cheez its"
## [19] "chicken malai kabab"
## [21] "chicken tenders"
## [23] "chicken tikka masala"
## [25] "chinese"
## [27] "chocolate"
## [29] "comfort"
## [31] "corn"
## [33] "arepa"
## [35] "asian"
## [37] "baguettes"
## [39] "blue point oyster"
## [41] "bread"
## [43] "burritos"
## [45] "cake"
## [47] "cheese"
## [49] "chicken"
## [51] "chicken parmesan"
## [53] "chicken tikka"
## [55] "chicken wings"
## [57] "chipotle"
## [59] "cinnamon rolls"
## [61] "cookies"
## [63] "cottage pie"
```

```

## [33] "crepes"
## [35] "dal bhaat"
## [37] "dessert"
## [39] "dried mangos"
## [41] "empanadas"
## [43] "escargot in garlic butter"
## [45] "farofa"
## [47] "fish tacos"
## [49] "freeze dried"
## [51] "fried fish"
## [53] "fried rice"
## [55] "fruit"
## [57] "ginger"
## [59] "gnocchi"
## [61] "grilled chicken breast"
## [63] "gyros"
## [65] "hibachi"
## [67] "hotpot"
## [69] "indian"
## [71] "italian"
## [73] "jollof rice"
## [75] "kelewele its an african dish"
## [77] "lahmacun turkish pizza"
## [79] "lebanese"
## [81] "macaroni"
## [83] "meat"
## [85] "mexican"
## [87] "my favorite kind of"
## [89] "noodles"
## [91] "pasta"
## [93] "persian"
## [95] "pizza"
## [97] "poke bowl"
## [99] "puerto rican"
## [101] "ribs"
## [103] "salmon"
## [105] "seafood"
## [107] "soup"
## [109] "spicy"
## [111] "strawberries"
## [113] "sushi"
## [115] "sweet chili doritos"
## [117] "tagine"
## [119] "tofu"

"curry"
"delicious"
"donuts"
"dumplings"
"enchiladas"
"ethiopian"
"fish"
"flautas"
"fried chicken"
"fried okra"
"fries"
"fruits"
"gizzard chicken"
"good pizza"
"guacamole"
"ham"
"hot dog"
"ice cream"
"indonesian"
"japanese"
"jollof rice with chicken"
"korean"
"lasagna"
"mac"
"mango"
"mediterranean"
"multigrain pancake"
"nectarines"
"palak paneer"
"patty melt"
"pho"
"platanos"
"postickers"
"ramen"
"roast dinner"
"salty"
"shahi paneer"
"spaghetti"
"steak"
"sundried tomatoes"
"swedish meatballs"
"sweet potato"
"thai"
"vietnamese"

length(unique(food))

## [1] 120

```

(3.9) (15 pts) Following the example from Class 05, display a dataframe of the sorted tabular results of food to see how many individuals prefer each kind of food.

From here on, the decisions of how to clean and combine categories are yours! Any food that currently has a count of 3 or more should remain (you can add to these categories - for example, you could add 'lasagna' to 'italian' or to 'pasta'). All other levels should be recoded or incorporated into a 'miscellaneous' food category. Points awarded based on thoughtfulness, effort, and quality/preciseness of your code.

Include your code below, and add comments where appropriate to describe the choices you make. You should have no more than 40 levels by the time you finish.

Display a dataframe of the sorted tabular results of food to see how many individuals prefer each kind of food AGAIN after you've finished your coding.

```
table1 <- data.frame(sort(table(food), decreasing = T))
table1
```

##	food	Freq
## 1	sushi	24
## 2	pizza	17
## 3	korean	14
## 4	thai	14
## 5	chinese	13
## 6	ice cream	10
## 7	mexican	10
## 8	indian	9
## 9	italian	9
## 10	ramen	9
## 11	noodles	8
## 12	pasta	8
## 13	steak	8
## 14	chocolate	6
## 15	hotpot	6
## 16	asian	5
## 17	cheese	5
## 18	curry	5
## 19	soup	5
## 20	fries	4
## 21	burger	3
## 22	fruit	3
## 23	japanese	3
## 24	lasagna	3
## 25	mediterranean	3
## 26	pho	3
## 27	spicy	3
## 28	brazilian	2
## 29	bread	2
## 30	burritos	2
## 31	chicken	2

## 32	chicken wings	2
## 33	cookies	2
## 34	delicious	2
## 35	farofa	2
## 36	fried chicken	2
## 37	guacamole	2
## 38	salmon	2
## 39	strawberries	2
## 40	tofu	2
## 41	vietnamese	2
## 42	albanian	1
## 43	arepa	1
## 44	artichoke	1
## 45	bagel	1
## 46	baguettes	1
## 47	bananas	1
## 48	blue point oyster	1
## 49	butter chicken	1
## 50	cake	1
## 51	ceviche	1
## 52	cheez its	1
## 53	chicken malai kabab	1
## 54	chicken parmesan	1
## 55	chicken tenders	1
## 56	chicken tikka	1
## 57	chicken tikka masala	1
## 58	chipotle	1
## 59	cinnamon rolls	1
## 60	comfort	1
## 61	corn	1
## 62	cottage pie	1
## 63	crepes	1
## 64	dal bhaat	1
## 65	dessert	1
## 66	donuts	1
## 67	dried mangos	1
## 68	dumplings	1
## 69	empanadas	1
## 70	enchiladas	1
## 71	escargot in garlic butter	1
## 72	ethiopian	1
## 73	fish	1
## 74	fish tacos	1
## 75	flautas	1
## 76	freeze dried	1
## 77	fried fish	1
## 78	fried okra	1
## 79	fried rice	1
## 80	fruits	1
## 81	ginger	1

```

## 82          gizzard chicken    1
## 83              gnocchi       1
## 84          good pizza        1
## 85    grilled chicken breast  1
## 86              gyros         1
## 87              ham          1
## 88              hibachi       1
## 89              hot dog       1
## 90              indonesian    1
## 91              jollof rice   1
## 92    jollof rice with chicken 1
## 93    kelewele its an african dish 1
## 94      lahmacun turkish pizza 1
## 95              lebanese      1
## 96              mac           1
## 97              macaroni      1
## 98              mango        1
## 99              meat         1
## 100     multigrain pancake     1
## 101     my favorite kind of   1
## 102              nectarines   1
## 103      palak paneer         1
## 104      patty melt          1
## 105              persian      1
## 106              platanos     1
## 107      poke bowl           1
## 108      postickers          1
## 109      puerto rican        1
## 110              ribs        1
## 111      roast dinner        1
## 112              salty       1
## 113              seafood      1
## 114      shahi paneer        1
## 115              spaghetti    1
## 116      sundried tomatoes   1
## 117      swedish meatballs   1
## 118      sweet chili doritos  1
## 119      sweet potato        1
## 120              tagine       1

```

all the food items with a count of 2 or less were combined together in a group wherever applicable i.e recoded into a similar group or put into already existing category or added to misc. food category.

#Fruits

```

food <- gsub("bananas|mango|nectarines|platanos|strawberries|dried
mangos|fruits", "fruit", food)

```

#Indian Cuisine

```

food <- gsub("chicken tikka|palak paneer|dal bhaat|butter chicken|shahi
paneer|chicken malai kabab|fried okra", "indian", food)

food <- gsub("indian masala", "indian", food) # combined extra Indian food
items which were left in the previous

#seafood
food <- gsub("fish|salmon|blue point oyster|fried fish", "seafood", food)

#dessert
food <- gsub("cinnamon rolls|crepes|multigrain pancake|cookies|donuts|cake",
"dessert", food)

#Pizza
food <- gsub("lahmacun turkish pizza|good pizza", "pizza", food)

#Pasta
food <- gsub("macaroni|mac|spaghetti", "pasta", food)

#Mexican dishes to be combined under Mexican Cuisine
food <- gsub("seafood tacos|flautas|burritos|enchiladas|chipotle", "mexican",
food)

#meat
food <- gsub("chicken tenders|ham|chicken wings|fried chicken|swedish
meatballs|chicken|chicken parmesan|cottage pie|gizzard chicken|ribs|roast
dinner", "meat", food)

food <- gsub("grilled meat breast", "meat", food) #somehow grilled meat
breast wasnt getting replaced in the previous code therefore created a new
code for it

#African dishes to be combined under African Cuisine
food <- gsub("tagine|jollof rice|jollof rice with chicken|kelewele its an
african dish", "african", food)

food <- gsub("african with meat", "african", food) #somehow african with meat
wasnt getting replaced in the previous code therefore created a new code for
it

#Chinese dishes to be combined under Chinese Cuisine
food <- gsub("fried rice|dumplings|postickers", "chinese", food)

#Bread products to be combined under Bread Category
food <- gsub("baguettes|bread|bagel", "bread", food)

#Burger/Sandwich types products to be combined under Burger Category
food <- gsub("hot dog|patty melt", "burger", food)

```

```

#Vegetables
food <- gsub("sweet potato|sundried
tomatoes|artichoke|corn|guacamole", "vegetable", food)

#farofa under Brazilian cuisine category
food <- gsub("farofa", "brazilian", food)

#Other Cuisine
food <- gsub("vietnamese|albanian|ethiopian|persian|puerto
rican|indonesian|persian|lebanese", "other_cuisine", food)

#Japanese dishes to be combined under Japanese Cuisine
food <- gsub("hibachi", "japanese", food)

#Italian dishes to be combined under Italian Cuisine
food <- gsub("gnocchi", "italian", food)

#snacks
food <- gsub("cheez its|sweet chili doritos", "snacks", food)

#mediterranean foo items
food <- gsub("gyros", "mediterranean", food)

sort(unique(food))

## [1] "african"
## [3] "asian"
## [5] "bread"
## [7] "ceviche"
## [9] "chinese"
## [11] "comfort"
## [13] "delicious"
## [15] "empanadas"
## [17] "freeze dried"
## [19] "fruit"
## [21] "hotpot"
## [23] "indian"
## [25] "japanese"
## [27] "lasagna"
## [29] "mediterranean"
## [31] "my favorite kind of"
## [33] "other_cuisine"
## [35] "pho"
## [37] "poke bowl"
## [39] "salty"
## [41] "snacks"
## [43] "spicy"

"arepa"
"brazilian"
"burger"
"cheese"
"chocolate"
"curry"
"dessert"
"escargot in garlic butter"
"fries"
"ginger"
"ice cream"
"italian"
"korean"
"meat"
"mexican"
"noodles"
"pasta"
"pizza"
"ramen"
"seafood"
"soup"
"steak"

```



```

## [45] "sushi"
## [47] "tofu"

length(unique(food))

## [1] 48

table2 <- data.frame(sort(table(food), decreasing = T))
table2

##           food Freq
## 1          sushi  24
## 2          pizza  19
## 3         indian  17
## 4         chinese  16
## 5          meat  16
## 6         mexican  16
## 7          korean  14
## 8           thai  14
## 9          fruit  11
## 10         pasta  11
## 11        ice cream  10
## 12         italian  10
## 13          ramen   9
## 14         dessert   8
## 15         noodles   8
## 16    other_cuisine   8
## 17          steak   8
## 18        chocolate   6
## 19         hotpot   6
## 20         seafood   6
## 21        vegetable   6
## 22          asian   5
## 23         burger   5
## 24          cheese   5
## 25          curry   5
## 26          soup   5
## 27         african   4
## 28        brazilian   4
## 29          bread   4
## 30          fries   4
## 31         japanese   4
## 32    mediterranean   4
## 33         lasagna   3
## 34           pho   3
## 35          spicy   3
## 36        delicious   2
## 37          snacks   2
## 38          tofu   2
## 39          arepa   1
## 40         ceviche   1

```

```
## 41          comfort      1
## 42          empanadas    1
## 43 escargot in garlic butter 1
## 44          freeze dried 1
## 45          ginger      1
## 46      my favorite kind of 1
## 47          poke bowl    1
## 48          salty       1

#categorizing the row after 36 into misc food category
for (i in 36:49){
  food <- gsub(paste0("^", as.character(table2[i, 1])), "misc. food", food)
}

table3 <- data.frame(sort(table(food), decreasing = T))
table3
```

##	food	Freq
## 1	sushi	24
## 2	pizza	19
## 3	indian	17
## 4	chinese	16
## 5	meat	16
## 6	mexican	16
## 7	misc. food	16
## 8	korean	14
## 9	thai	14
## 10	fruit	11
## 11	pasta	11
## 12	ice cream	10
## 13	italian	10
## 14	ramen	9
## 15	dessert	8
## 16	noodles	8
## 17	other_cuisine	8
## 18	steak	8
## 19	chocolate	6
## 20	hotpot	6
## 21	seafood	6
## 22	vegetable	6
## 23	asian	5
## 24	burger	5
## 25	cheese	5
## 26	curry	5
## 27	soup	5
## 28	african	4
## 29	brazilian	4
## 30	bread	4
## 31	fries	4
## 32	japanese	4

```
## 33 mediterranean    4
## 34      lasagna     3
## 35        pho       3
## 36        spicy     3

sort(unique(food))

## [1] "african"      "asian"        "brazilian"    "bread"
## [5] "burger"       "cheese"       "chinese"      "chocolate"
## [9] "curry"        "dessert"      "fries"        "fruit"
## [13] "hotpot"       "ice cream"    "indian"       "italian"
## [17] "japanese"     "korean"       "lasagna"      "meat"
## [21] "mediterranean" "mexican"     "misc. food"   "noodles"
## [25] "other_cuisine" "pasta"       "pho"          "pizza"
## [29] "ramen"        "seafood"     "soup"         "spicy"
## [33] "steak"        "sushi"       "thai"         "vegetable"

length(unique(food))

## [1] 36
```

(3.10) (8 pts) Final steps and a plot: You'll want to CAREFULLY follow the example in the code at the end of Class 05.

- Use the `toTitleCase()` function from the package `tools` to convert food to title case.
- Make an object called `finaltab` that is a table of your final vector `food`.
- Calculate percents, rounded to the nearest integer, for each food type. Save this as an object called `percents`.
- Change the names of `finaltab` to include a space and then the percents followed by a "%" in curved parentheses.
- Make a horizontal barplot of your final plot. Choose a nice bar color, adjust the left margins as necessary, give a main title and label the horizontal axis.

```
library(tools)
food <- toTitleCase(food)

finaltab <- table(food)

percents <- round(finaltab/sum(finaltab)*100, 1)
percents

## food
##      African      Asian      Brazilian      Bread      Burger
##          1.3          1.6          1.3          1.3          1.6
##      Cheese      Chinese      Chocolate      Curry      Dessert
##          1.6          5.0          1.9          1.6          2.5
##      Fries      Fruit      Hotpot      Ice Cream      Indian
##          1.3          3.5          1.9          3.2          5.4
##      Italian      Japanese      Korean      Lasagna      Meat
##          3.2          1.3          4.4          0.9          5.0
```

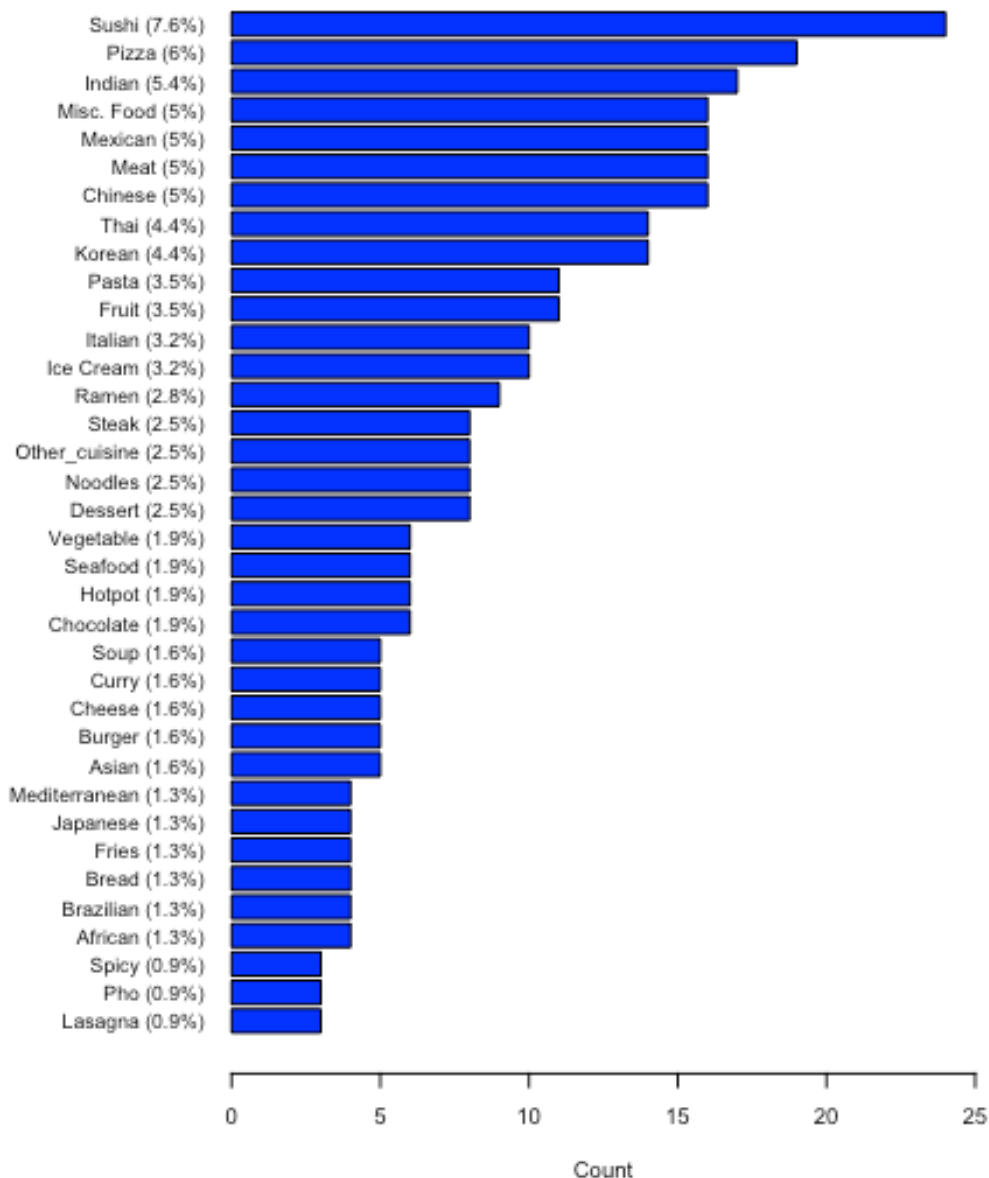
## Mediterranean	Mexican	Misc. Food	Noodles	Other_cuisine
## 1.3	5.0	5.0	2.5	2.5
## Pasta	Pho	Pizza	Ramen	Seafood
## 3.5	0.9	6.0	2.8	1.9
## Soup	Spicy	Steak	Sushi	Thai
## 1.6	0.9	2.5	7.6	4.4
## Vegetable				
## 1.9				

```
names(finaltab) <- paste0(names(finaltab), " (", percents, "%)")
```

```
par(mar = c(5, 12, 4, 2), cex = 0.6)
```

```
barplot(sort(finaltab),
        horiz = T,
        las = 1, col = "blue",
        main = "Favorite Foods Among Students of Class S&DS 230",
        xlab = "Count",
        cex.names = 0.9,
        xlim = c(0,25))
```

Favorite Foods Among Students of Class S&DS 230



(3.11) (3 pts) In no more than three sentences, discuss your process and results. Be sure to mention how many unique values of 'food' you started and ended with. Any surprises?

the dataset had 218 unique food items. After doing the suggested steps in the assignment such as converting to lowercase, removing regular expressions, stopwords, combining the same cuisine into one I was left with 120 unique food items. To get to a level of less than 40 food items I re-coded all the seafood items into one category, since there was a category of fruit, I included other fruit into one category. Similarly I did for pizza, Mexican dishes, pizza, desserts. There were various other cuisines with a count 1 so I re-coded them into one category as

other_cuisine. In the end I was left with 36 unique food items. It was interesting to know that top 5 food liked by the students are sushi, pizza, Indian cuisine, Mexican, Chinese and meat/chicken items.