

Problem Set 5 Solutions

CSCI 5a: Programming in R

Let's clear the global computing environment:

```
rm( list = ls() )
```

```
load( "Problem Set 5 R Objects.Rdata" )
```

```
ls()
```

```
## [1] "problem.1.data" "problem.2.data" "problem.3.data" "problem.5.data"  
## [5] "problem.8.data"
```

Problem 1: Report Missing Values

In this problem, we will examine the vector `problem.1.data` for missing data.

Part (a)

Does the vector contain any NA values? Write R code to determine this, and then report your answer using one or two sentences of text.

Solution

```
any( is.na( problem.1.data ) )
```

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

At least one element in `problem.1.data` has the special value NA.

Part (b)

How many elements of `problem.1.data` have the value NA? Report your result using a `cat()` statement.

Solution

```
na.count <-  
  sum( is.na( problem.1.data ) )  
  
cat(  
  "Number of elements with NA:",  
  na.count  
)
```

```
## Number of elements with NA: 4
```

Four elements of `problem.1.data` have the special value NA.

Part (c)

What proportion of the elements of `problem.1.data` have the value NA? Report your result using a `cat()` statement, displaying this value using 2 decimal places.

Solution

```
na.proportion <-  
  mean( is.na( problem.1.data ) )  
  
cat(  
  "Proportion of elements that are NA:",  
  formatC(  
    na.proportion,  
    format = "f",  
    digits = 2  
  )  
)
```

```
## Proportion of elements that are NA: 0.22
```

Approximately 22% of the elements of `problem.1.data` have the special value NA.

Part (d)

Determine which locations of `problem.1.data` contain NA values, and save these locations in a variable. Then use a `for` loop to print out a simple report of these values:

- Each location should be reported on a single line.
- Enumerate the items, so that each line is numbered.
- Make sure that you have correct spacing in your answer.
- Each location should be reported on a single line.
- Enumerate the items, so that each line is numbered.
- Make sure that you have correct spacing in your answer.

Solution

```
na.locations <-  
  which( is.na( problem.1.data ) )  
  
for( index in 1:length(na.locations) ) {  
  
  cat(  
    index,
```

```
        ". Location: ",
        na.locations[ index ],
        "\n",
        sep = ""
    )
}
```

```
## 1. Location: 5
## 2. Location: 8
## 3. Location: 11
## 4. Location: 12
```

End of problem 1

Problem 2: Report and Replace -9 Values

Some people use the numeric value -9 to represent missing data instead of NA. I think this is a bad idea, but people do it nonetheless. The best practice here is that if you encounter a dataset that uses -9 values to represent missing data you should immediately convert these to NA values, and in this problem we're going to explore this idea.

Part (a): Sample mean before cleaning

To see why it's a bad idea to represent missing data by using -9 instead of NA, let's start by calculating the sample mean of `problem.2.data`. Report your result using a `cat()` statement, displaying this value with 2 decimal places.

Solution

```
cat(
  "Sample mean (before cleaning): ",
  formatC(
    mean( problem.2.data ),
    format = "f",
    digits = 2
  ),
  "\n",
  sep = ""
)
```

```
## Sample mean (before cleaning): -0.20
```

Part (b): Detecting -9 values

Does the vector `problem.2.data` contain any -9 values? Write R code to determine this, and then report your answer using one or two sentences of text.

Solution

```
any( problem.2.data == -9 )
```

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

Thus, the vector `problem.2.data` contains at least one -9 value.

Part (c): Counting -9 elements

How many elements of `problem.2.data` have the value -9? (Hint: use a vectorized comparison operation.) Report your result using a `cat()` statement.

Solution

```
cat(
  "Number of -9 values:",
  sum( problem.2.data == -9 )
)
```

```
## Number of -9 values: 7
```

Part (d): Proportion of elements equal to -9

What proportion of the elements of `problem.2.data` have the value -9? Report your result using a `cat()` statement, displaying this value using 2 decimal places.

Solution

```
cat(
  "Proportion of -9 values:",
  formatC(
    mean( problem.2.data == -9 ),
    format = "f",
    digits = 2
  )
)
```

```
## Proportion of -9 values: 0.08
```

Thus, approximately 8% of the elements in `problem.2.data` have the value -9.

Part (e): Locations of -9 elements

Determine which locations of `problem.2.data` contain -9 values, and save these locations in a variable. Then use a `for` loop to print out a simple report of these values:

- Each location should be reported on a single line.
- Enumerate the items, so that each line is numbered.
- Make sure that you have correct spacing in your answer.

Solution

```
minus.9.locations <-
  which( problem.2.data == -9 )

for( index in 1:length(minus.9.locations) ) {

  cat(
    index,
    ". Location: ",
    minus.9.locations[ index ],
    "\n",
    sep = ""
  )
}
```

```
## 1. Location: 15
## 2. Location: 21
## 3. Location: 22
## 4. Location: 39
## 5. Location: 58
## 6. Location: 62
## 7. Location: 85
```

Part (f): Replacing -9 with NA

Assign the value `NA` to the locations in `problem.2.data` that currently have a -9 value.

There are a couple of ways to do this:

- Use logical indexing on the left-hand side of an assignment operator to select all the elements of the vector that are equal to -9, and assign `NA` to these elements.
- Use positive integer indexing along with the vector of locations from part (e) on the left-hand side of an assignment operator to select all the elements of the vector that are equal to -9, and assign `NA` to these elements.

You only need one line of code to assign `NA` to all the elements that are equal to -9.

There's nothing to report here, but write your code clearly so that the TAs can understand it.

Solution

First, let's do this with logical indexing:

```
problem.2.data[ problem.2.data == -9 ] <-  
  NA
```

Now let's do it with positive integer indexing:

```
problem.2.data[ minus.9.locations ] <-  
  NA
```

Part (g): Sample mean after cleaning

Calculate the sample mean of the non-missing values of the corrected version of `problem.2.data` that you created in part (f). Report your result using a `cat()` statement, displaying this value with 2 decimal places.

Write a single sentence comparing your answer for this part with your answer for part (a). Now do you see why it's a bad idea to use -9 to represent a missing value?

Solution

```
cat(  
  "Sample mean (after cleaning): ",  
  formatC(  
    mean( problem.2.data, na.rm = TRUE ),  
    format = "f",  
    digits = 2  
  ),  
  "\n",  
  sep = ""  
)
```

```
## Sample mean (after cleaning): 0.52
```

Before we repaired the -9 values, the sample mean was -0.20, but once we fixed these the sample mean was 0.52.

End of problem 2

Problem 3: Report and Replace Outliers

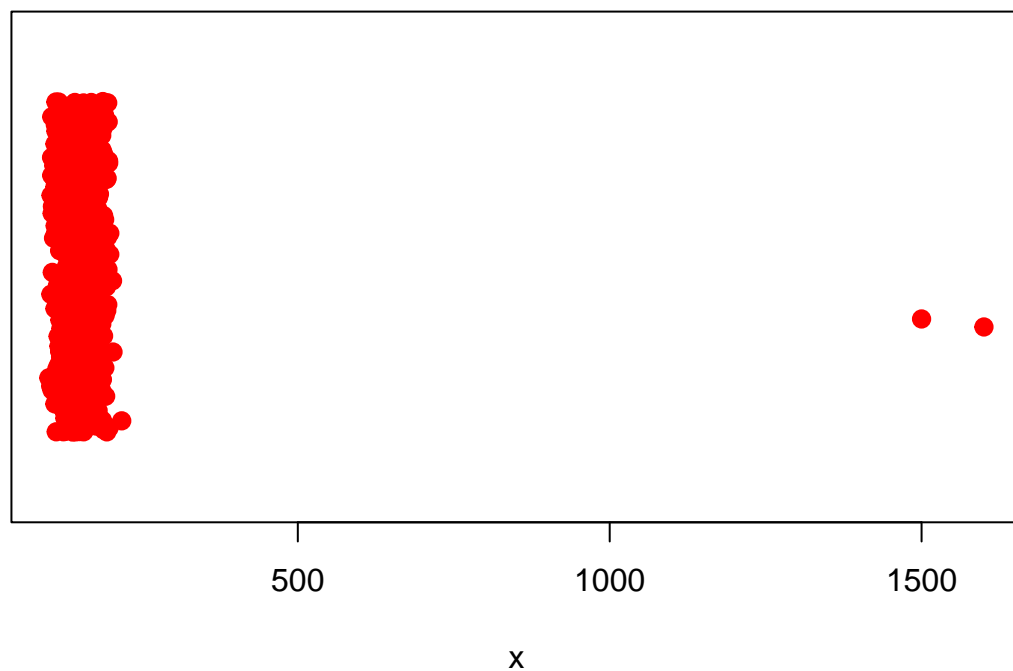
Part (a): Stripchart

Construct a stripchart of the values in `problem.3.data`. Do you think that there are any outliers in this data? Explain your answer with one or two sentences.

Solution

```
stripchart(  
  problem.3.data,  
  ylim = c(0, 2),  
  main = "Stripchart of problem.3.data",  
  xlab = "x",  
  method = "jitter",  
  jitter = 0.7,  
  pch = 19,  
  cex = 1.2,  
  col = "red"  
)
```

Stripchart of problem.3.data



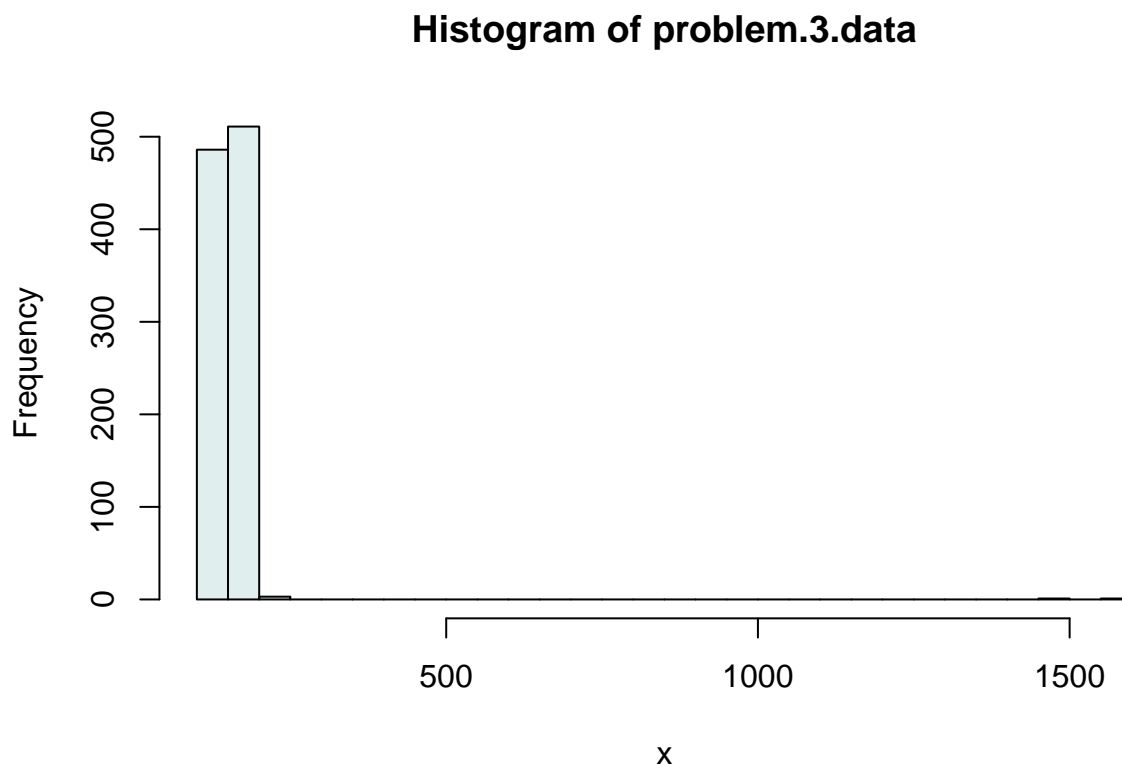
It looks like there are two outliers in this data.

Part (b): Histogram

Construct a histogram of the values in `problem.3.data`. Do you think that there are any outliers in this data? Explain your answer with one or two sentences.

Solution

```
hist(  
  problem.3.data,  
  main = "Histogram of problem.3.data",  
  xlab = "x",  
  ylab = "Frequency",  
  col = "azure2",  
  breaks = 50  
)
```



It's much harder to see the outliers in this dataset. If you look closely, you can see two tiny little rectangles at about 1500. However, the real indication that there are outliers is simply the scale of the x -axis: since it ranges from 0 up to 1500, that means that there must be some points around 1500.

Part (c): Assigning NA

For this problem, we're going to assign the value NA to the outliers. The challenge here is to specify the outliers. I'm going to suggest that you select a threshold value, and anything above that threshold value is considered an outlier. You'll have a lot of flexibility as to how you define the threshold, but choose something reasonable.

There are two approaches here:

- Use the `which()` function with a vectorized comparison operation to determine the exact locations of the outliers in `problem.3.data`, and then use positive integer indexing to assign `NA` to these values.
- Use a logical comparison operation to determine the exact locations of the outliers in `problem.3.data`, and then use logical indexing to assign `NA` to these values.

You can do this with just a few lines of code, so if you're writing lots of code, then you're doing too much work.

Solution

I'm going to use a threshold value of 1000:

```
threshold.value <- 1000
```

Now we can use the `which()` function to determine the exact locations of the outliers:

```
outlier.locations <-  
  which( problem.3.data > threshold.value )  
  
outlier.locations
```

```
## [1] 520 521
```

Let's check this:

```
problem.3.data[ outlier.locations ]
```

```
## [1] 1600 1500
```

Those are the outliers!

Now we can assign `NA` to these values:

```
cleaned.problem.3.data <-  
  problem.3.data  
  
cleaned.problem.3.data[ outlier.locations ] <- NA
```

Part (d): Stripchart

Construct a stripchart of the values in `problem.3.data` with the outliers removed. Do you think that there are any outliers in this data? Explain your answer with one or two sentences.

Solution

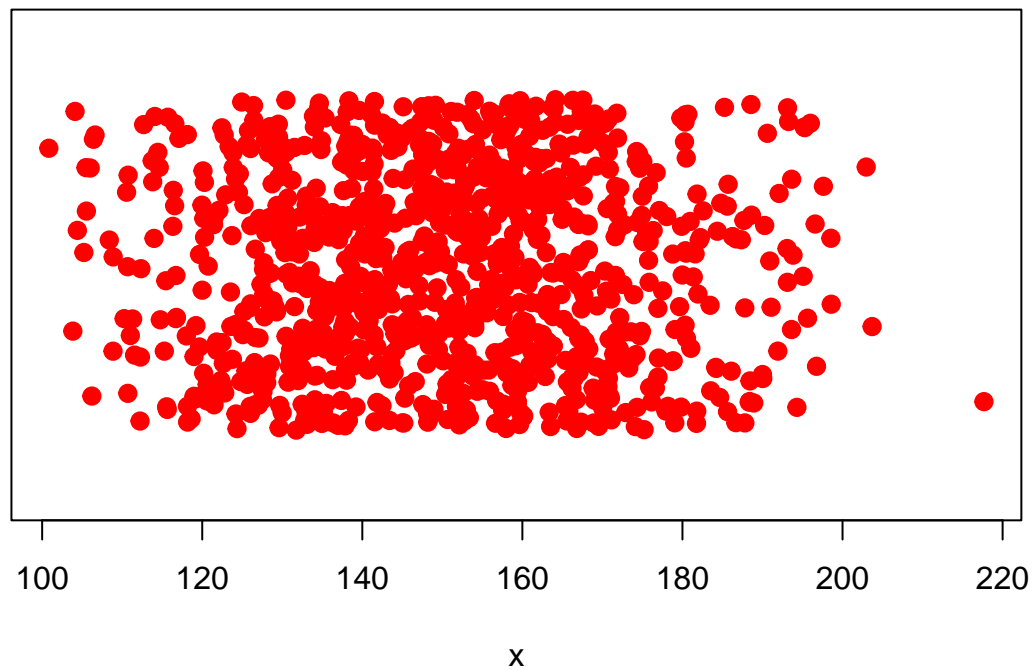
```
stripchart(  
  cleaned.problem.3.data,  
  ylim = c(0, 2),  
  main = "Stripchart of problem.3.data",  
  xlab = "x",
```

```

method = "jitter",
jitter = 0.7,
pch = 19,
cex = 1.2,
col = "red"
)

```

Stripchart of problem.3.data



Very different!

Part (e): Histogram

Construct a histogram of the values in your cleaned version of `problem.3.data`. Then superimpose the best-fitting normal density curve on the data.

There's a subtle point here: when we fixed the outliers, we changed their values to `NA`. But when we fit the normal density curve, we calculate the sample mean and sample standard deviation. What do we have to do to calculate the sample mean and sample standard deviation when there are `NA` values in the data?

Solution

```

hist(
  cleaned.problem.3.data,
  prob = TRUE,
  main = "Histogram of cleaned problem.3.data",
  xlab = "x",
  ylab = "Density",
)

```

```

col = "azure2",
breaks = 50
)

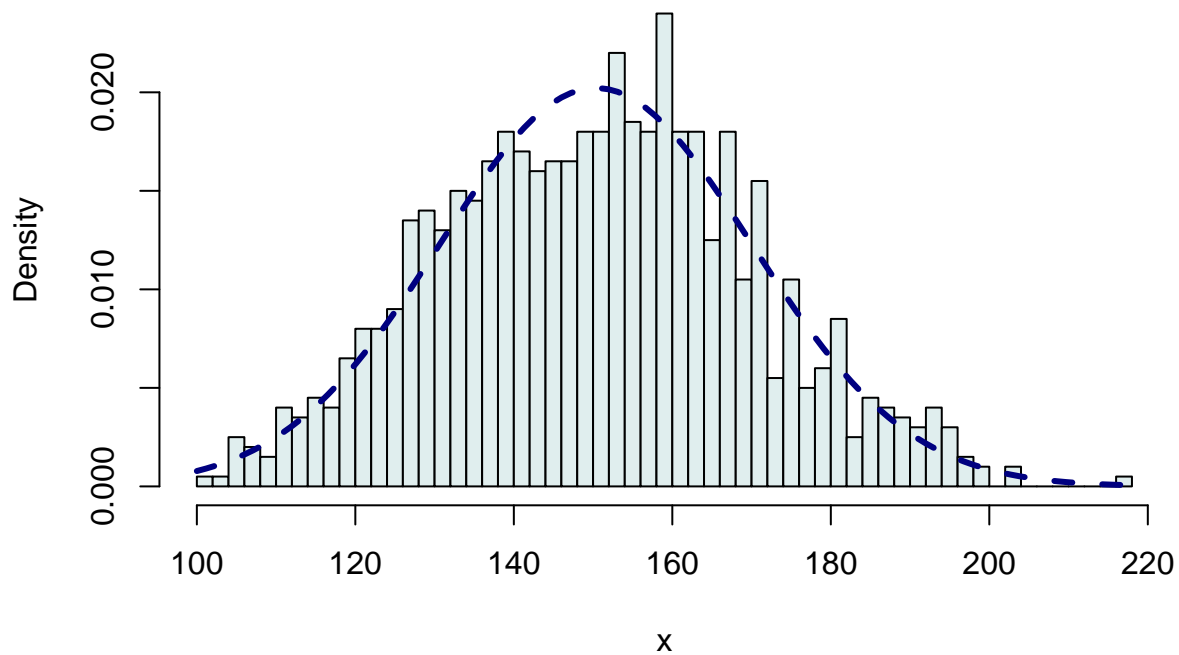
sample.mean <-
  mean( cleaned.problem.3.data, na.rm = TRUE )

sample.standard.deviation <-
  sd( cleaned.problem.3.data, na.rm = TRUE )

curve(
  dnorm(
    x,
    mean = sample.mean,
    sd = sample.standard.deviation
  ),
  lty = "dashed",
  lwd = 3,
  col = "navy",
  add = TRUE
)

```

Histogram of cleaned problem.3.data



End of problem 3

Problem 4: Grocery Store Sales

Let's return to our grocery store example.

This week, the prices per box for each brand of breakfast cereal are:

Brand	Price
SBZ	2.79
KYM	3.99
HKT	8.49

Here we have a sequence of transactions:

Transaction	Brand	Number of Boxes
1	KYM	2
2	SBZ	4
3	SBZ	3
4	HKT	1
5	SBZ	3
6	KYM	2
7	SBZ	5

Our goal in this problem is to calculate the total amount in sales for Sugar Bomz (SBZ).

Part (a): Constructing vectors

Construct three vectors:

- A vector consisting of the brand names for each transaction.
- A vector consisting of the number of boxes sold for each transaction.
- A pricing lookup vector that converts the brand name to the price per box.

There's nothing to report here, but write your code clearly so the TAs can understand what you're doing.

Solution

```
brand.vector <-  
  c( "KYM", "SBZ", "SBZ", "HKT",  
      "SBZ", "KYM", "SBZ")  
  
number.of.bboxes.vector <-  
  c( 2, 4, 3, 1, 3, 2, 5 )  
  
pricing.lookup.vector <-  
  c(  
    "SBZ" = 2.79,  
    "KYM" = 3.99,  
    "HKT" = 8.49  
  )
```

Part (b): Converting the brand name vector

Use the pricing lookup vector from part (a) to convert the brand name vector from part (a) into a vector of prices per box. Display this vector of prices per box directly.

Solution

```
price.per.box.vector <-  
  pricing.lookup.vector[  
    brand.vector  
  ]  
  
price.per.box.vector
```

```
## KYM SBZ SBZ HKT SBZ KYM SBZ  
## 3.99 2.79 2.79 8.49 2.79 3.99 2.79
```

Part (c): Transaction sales amount vector

Use a vectorized operation with the number of boxes sold vector from part (a) and the price per box vector from part (b) to construct a vector consisting of the sales amount for each transaction. Display this vector of sales amounts for each transaction directly.

Solution

```
transaction.sales.amount.vector <-  
  price.per.box.vector *  
  number.of.boxes.vector  
  
transaction.sales.amount.vector
```

```
## KYM SBZ SBZ HKT SBZ KYM SBZ  
## 7.98 11.16 8.37 8.49 8.37 7.98 13.95
```

Part (d): Sugar Bomz transaction sales amount vector

We saw in lecture how to use values in one vector to select values from another vector.

In this part, we're going to use the values in the brand name vector from part (a) to select values from the transaction sales amount vector we constructed in part (c).

In particular, use logical indexing to select the values in the transaction sales amount vector from part (c) corresponding to a Sugar Bomz (SBZ) sale.

Display this vector of total Sugar Bomz sales amounts directly.

Solution

```
sbz.sales.amount.vector <-  
  transaction.sales.amount.vector[  
    brand.vector == "SBZ"  
  ]  
  
sbz.sales.amount.vector
```

```
## SBZ SBZ SBZ SBZ  
## 11.16 8.37 8.37 13.95
```


Part (e): Sum of Sugar Bomz transaction sales

Calculate the sum of the values in the vector of the Sugar Bomz transaction sales amounts that you constructed in part (d). Report your result using a `cat()` statement, displaying this value to 2 decimal places.

Solution

```
cat(  
  "Sugar Bomz total sales amount:",  
  formatC(  
    sum( sbz.sales.amount.vector ),  
    format = "f",  
    digits = 2  
  )  
)
```

```
## Sugar Bomz total sales amount: 41.85
```

End of problem 4

Problem 5: Two-Tone Stripchart

We can use logical indexing techniques to create an informative data visualization that I call a “two-tone” stripchart.

The idea of a two-tone stripchart is that we have a threshold value, and we want to emphasize or highlight the values that are greater than this threshold value.

For this problem, we’ll use a threshold value of 200:

```
threshold <- 200
```

The data for this problem is contained in the vector `problem.5.data`.

Part (a): Low-pass filtering

Use logical indexing to select all the values in `problem.5.data` that are less than or equal to the threshold value. Store these values in a variable, and report the first 5 elements using a `cat()` statment, displaying these values using 2 decimal places.

When we select values that are less than or equal to a cutoff value, this is called “low-pass filtering”, because we are removing high values from the data but allowing the low values to pass through.

Solution

```
low.values.vector <-  
  problem.5.data[ problem.5.data <= threshold ]  
  
cat(  
  "Low values:",  
  formatC(  
    head( low.values.vector, n = 5 ),  
    format = "f",  
    digits = 2  
  )  
)
```

```
## Low values: 0.16 73.70 142.26 41.21 5.13
```

Part (b): High-pass filtering

Use logical indexing to select all the values in `problem.5.data` that are strictly greater than the threshold value. Store these values in a variable, and report the first 4 elements using a `cat()` statment, displaying these values using 2 decimal places.

When we select values that are strictly greater than a cutoff value, this is called “high-pass filtering”, because we are removing low values from the data but allowing the high values to pass through.

Solution

```
high.values.vector <-  
  problem.5.data[ problem.5.data > threshold ]  
  
cat(  
  "High values:",  
  formatC(  
    head( high.values.vector, n = 4 ),  
    format = "f",  
    digits = 2  
  )  
)
```

```

    "High values:",
    formatC(
      head( high.values.vector, n = 5 ),
      format = "f",
      digits = 2
    )
  )
)

```

```
## High values: 343.92 214.80 231.94 206.76 211.29
```

Part (c): Two-tone stripchart

Now we're going to construct a stripchart for this data, where the data below the threshold is displayed using one color and the data above the threshold above the threshold is displayed using a different color.

- First, use the low values from part (a) to construct a stripchart, including titles and using jitter. You should explicitly select the point shape, size, and color. Finally, specify the range of the x -axis to be from 0 to 400.
- Next, we're going to add another stripchart to this graph. Construct another stripchart, this time using the high values from part (b). You don't have to specify the titles or the range of the x -axis, because those have already been determined. You *do* have to use jitter for these points, and explicitly specify the shape, size, and color of these points. Use a different color from what you used for the low values. Finally, include `add = TRUE`, just like with the `curve()` function.
- Finally, draw a vertical line with an x -value equal to the threshold.

All of these graphing actions have to occur within the same code chunk.

If you do this properly, you should end up with a stripchart where all the points below the threshold value have one color, and all the points above the threshold value have another color, and we can easily see the threshold value because it's represented by the vertical line.

Solution

```

stripchart(
  low.values.vector,
  xlim = c(0, 400),
  ylim = c(0, 2),
  main = "Two-tone stripchart",
  xlab = "x",
  method = "jitter",
  jitter = 0.7,
  pch = 19,
  cex = 1.2,
  col = "darkslategray3"
)

segments(
  x0 = threshold,
  y0 = 0.3,
  x1 = threshold,
  y1 = 1.7,

```

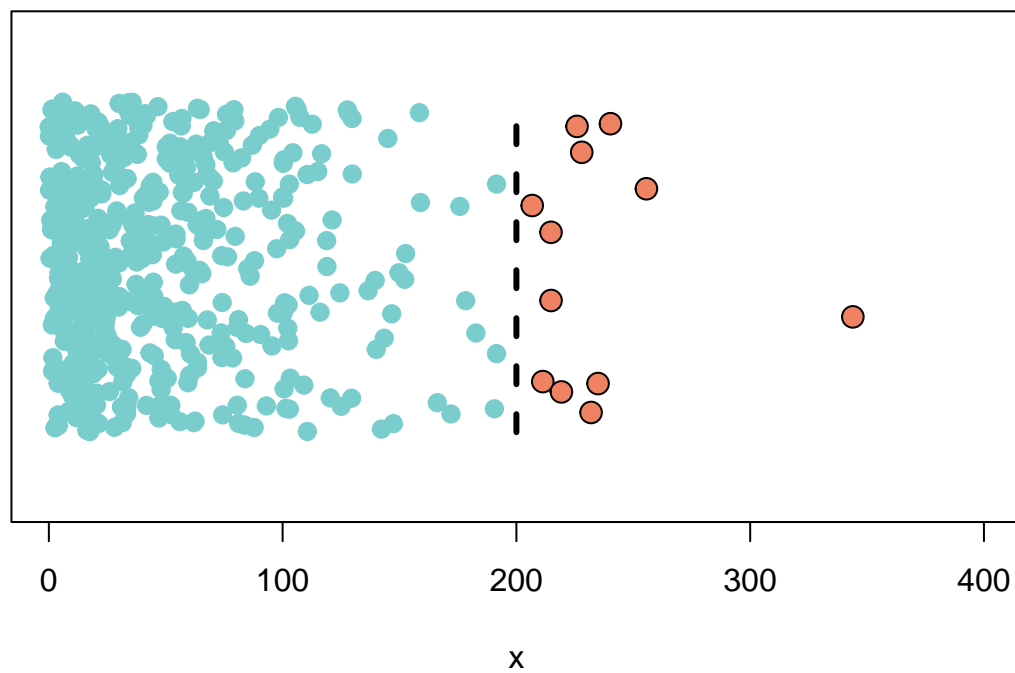
```

    lty = "dashed",
    lwd = 3,
    col = "black"
)

stripchart(
  high.values.vector,
  method = "jitter",
  jitter = 0.7,
  pch = 21,
  cex = 1.5,
  col = "black",
  bg = "salmon2",
  add = TRUE
)

```

Two-tone stripchart



End of problem 5

Problem 6: Final Grades

We finally have the tools to calculate a final course score, given the raw scores. In this problem, we'll go from the initial raw scores all the way to the final course score, using R techniques from every lecture so far.

Five students have these final raw scores:

Status	Problem Sets	Midterm	Comprehensive Assessment
Graduate	58	74	74
Undergraduate	65	65	58
Graduate	64	76	74
Graduate	60	66	72
Undergraduate	65	57	61

For this course:

- A raw score of 68 points on the problem sets is equivalent to a standardized score of 100.
- A raw score of 80 points on the Midterm Assessment is equivalent to a standardized score of 100.
- A raw score of 80 points on the Comprehensive Assessment is equivalent to a standardized score of 100.

Part (a): Preliminary Score 1

Using vectorized operations, calculate the Preliminary Score 1 for all 5 students. Report this vector using a `cat()` statement, displaying these values with 2 decimal places.

Solution

First, let's make some vectors to represent the data:

```
registration.status.vector <-  
  c( "Graduate", "Undergraduate", "Graduate",  
      "Graduate", "Undergraduate" )  
  
problem.set.raw.score.vector <-  
  c( 58, 65, 64, 60, 65 )  
  
midterm.exam.raw.score.vector <-  
  c( 74, 65, 76, 66, 57 )  
  
cca.raw.score.vector <-  
  c( 74, 58, 74, 72, 61 )
```

Now let's standardize these scores:

```
problem.set.standardized.score.vector <-  
  problem.set.raw.score.vector / 68 * 100  
  
midterm.exam.standardized.score.vector <-  
  midterm.exam.raw.score.vector / 80 * 100  
  
cca.standardized.score.vector <-  
  cca.raw.score.vector / 80 * 100
```

Now we can calculate the Preliminary Score 1:

```
preliminary.score.1.vector <-  
  (0.20 * problem.set.standardized.score.vector) +  
  (0.30 * midterm.exam.standardized.score.vector) +  
  (0.50 * cca.standardized.score.vector)  
  
cat(  
  "Preliminary Score 1 vector: ",  
  formatC(  
    preliminary.score.1.vector,  
    format = "f",  
    digits = 2  
  )  
)
```

```
## Preliminary Score 1 vector:  91.06 79.74 93.57 87.40 78.62
```

Part (b): Preliminary Score 2

Using vectorized operations, calculate the Preliminary Score 2 for all 5 students. Report this vector using a `cat()` statement, displaying these values with 2 decimal places.

Solution

```
preliminary.score.2.vector <-  
  (0.35 * midterm.exam.standardized.score.vector) +  
  (0.65 * cca.standardized.score.vector)  
  
cat(  
  "Preliminary Score 2 vector: ",  
  formatC(  
    preliminary.score.2.vector,  
    format = "f",  
    digits = 2  
  )  
)
```

```
## Preliminary Score 2 vector:  92.50 75.56 93.38 87.38 74.50
```

Part (c): Graduate Final Course Score

Recall that the graduate final course score is calculated as the maximum of the Preliminary Scores 1 and 2.

Using the vectors of preliminary scores that you created in parts (a) and (b), construct a vector of the graduate final course scores for the five students. You'll have to figure out how to take the maximum of the two scores for each student, but here are two suggestions:

- You could write a `for` loop.
- You could use the `ifelse()` function.

It's up to you.

Report this vector of graduate course scores using a `cat()` statement, displaying these values with 2 decimal places.

Solution

```
graduate.course.score.vector <-  
  numeric( length( preliminary.score.1.vector ) )  
  
for( index in 1:length(graduate.course.score.vector) ) {  
  
  graduate.course.score.vector[ index ] <-  
    max(  
      preliminary.score.1.vector[ index ],  
      preliminary.score.2.vector[ index ]  
    )  
}  
  
cat(  
  "Graduate course score vector:",  
  formatC(  
    graduate.course.score.vector,  
    format = "f",  
    digits = 2  
  )  
)
```

```
## Graduate course score vector: 92.50 79.74 93.57 87.40 78.62
```

Part (d): Scaling factor

Recall that I multiply the graduate course by a scaling factor, depending on the student's registration status:

Status	Scaling Factor
Graduate	1
Undergraduate	4/3

Create a lookup vector for this scaling factor. Then use this lookup vector to convert the registration status to a numeric vector of scaling factors. Report this vector of scaling factors using a `cat()` statement, displaying the values with 2 decimal places.

Solution

```
undergraduate.scaling.lookup.vector <-  
  c(  
    "Graduate" = 1.0,  
    "Undergraduate" = 4/3  
  )  
  
undergraduate.scaling.lookup.vector[  
  registration.status.vector  
]
```

##	Graduate	Undergraduate	Graduate	Graduate	Undergraduate
##	1.000000	1.333333	1.000000	1.000000	1.333333

Part (e): Final course score

Using a vectorized operation, multiply the values in the vector of graduate course scores from part (c) by the vector of scaling factors from part (d). This will be a vector consisting of the final course scores. Report this vector using a `cat()` statement, displaying the values with 2 decimal places.

Solution

```
final.course.score <-
  graduate.course.score.vector *
  undergraduate.scaling.lookup.vector[
    registration.status.vector
  ]

cat(
  "Final course score:",
  formatC(
    final.course.score,
    format = "f",
    digits = 2
  )
)
```

```
## Final course score: 92.50 106.32 93.57 87.40 104.82
```