

PSTAT 10 (Practice) Final Exam

Name:

TA:

You have 80 minutes to complete 6 problems, some of which contain subproblems.

You are allowed one sheet, two sides of handwritten notes.

If you run out of space for your code, write your code on the backs of the pages and clearly indicate which problem it corresponds to.

Please make sure your handwritten code is legible. I cannot give credit to code that is unreadable.

3. Write code to simulate, using 1,000 replications, the expected number of passengers on the bus after 3 bus stops. (2 pts)

```
# Your code here
```

```
# End
```

Problem 2: Poisson Distribution (6 pts)

How many emails will I receive in an hour? How many earthquakes will occur in a year in some region of the world? A popular probability distribution that models such events is the *Poisson distribution*. We write $X \sim \text{Pois}(\lambda)$ if X is a Poisson random variable with parameter λ , where λ is the *rate* of occurrence of an event given in emails/hour, earthquakes/year, etc.

R has built-in functions for the Poisson distribution. They are:

- `dpois(x, lambda)`
- `ppois(q, lambda)`
- `rpois(n, lambda)`

Let $X \sim \text{Pois}(2)$ be the number of emails I receive in an hour.

1. State the support of X . Is X continuous or discrete? (2 pts)

2. Determine the probability I receive greater than 5 but less than 10 emails in an hour. (2 pts)

```
# Your code here
```

```
# End
```

3. Determine the probability I receive no emails in an hour. (2 pts)

```
# Your code here
```

```
# End
```

Problem 3: Majority (5 pts)

Write a function `majority(v)` that takes a numeric vector and returns the *majority element*; i.e., the element that appears more than $n/2$ times, where n is the length of `v`. The input `v` is guaranteed to be nonempty and contain positive integers. A majority element of `v` is guaranteed to exist.

```
majority(c(3, 2, 3))
```

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

```
majority(c(2, 2, 1, 1, 1, 2, 2))
```

```
## [1] 2
```

```
majority <- function(v) {
# Your code here

# End
}
```

Hint: You may use the `tabulate(v)` function, which returns a vector of the counts of elements of a numeric vector `v` in the position indexed by that element.

```
# Each number from 1 through 3 appears once, 5 appears once, but 4 appears zero times.
tabulate(c(1:3, 5))
```

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

```
# 1 appears zero times, 2 appears two times, 3 appears one time,
# 4 appears zero times, 5 appears two times.
tabulate(c(2, 2, 3, 5, 5))
```

```
## [1] 0 2 1 0 2
```

Problem 4: Second largest (5 pts)

Write a function `second_largest(vec)` that takes a numeric vector `vec` and returns its second largest element. The elements of `vec` are guaranteed to be nonnegative and unique.

```
second_largest(1:10)
```

```
## [1] 9
```

```
second_largest(c(5, 0, 6.6, 3.13, 1.5))
```

```
## [1] 5
```

```
second_largest <- function(vec) {  
  # Your code here
```

```
  # End  
}
```

Problem 5: Some data (6 pts)

The tibble `some_data` contains some randomly generated values.

```
some_data

## # A tibble: 100 x 3
##       a         b c
##   <int> <dbl> <fct>
## 1     7 0.327 J
## 2     7 0.389 J
## 3     8 0.0411 I
## 4     5 0.361 K
## 5     8 0.571 J
## 6     8 0.685 I
## 7    10 0.971 J
## 8     7 0.702 K
## 9     8 0.0115 J
## 10    6 0.536 K
## # ... with 90 more rows
```

1. Create a tibble called `some_data1` of cases with `c` level of `K` and only variables `c` and `d` where `d` is the sum of variables `a` and `b`. Match the provided result exactly. (2 pts)

```
## # A tibble: 35 x 2
##       c         d
##   <fct> <dbl>
## 1 K     5.36
## 2 K     7.70
## 3 K     6.54
## 4 K     7.20
## 5 K     9.58
## 6 K     8.35
## 7 K    12.0
## 8 K     8.55
## 9 K     9.58
## 10 K    9.04
## # ... with 25 more rows
```

```
# Your code here
```

```
# End
```

2. Create a tibble `some_data2` containing, for each level of `c`, the sum of `a` and the mean of `b` for that level of `c`. Match the provided result exactly. (2 pts)

```
## # A tibble: 3 x 3
##   c      sum_a sum_b
##   <fct> <int> <dbl>
## 1 I         224 0.476
## 2 J         310 0.477
## 3 K         285 0.525
```

```
# Your code here
```

```
# End
```

3. Using `some_data2`, create the ggplot that's given on the next page. (2 pts)

```
# Your code here
```

```
# End
```

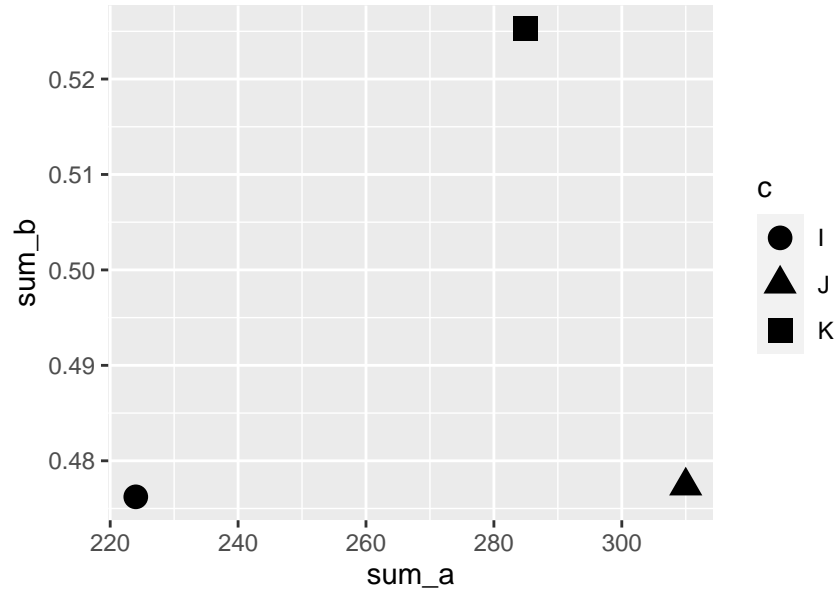



Figure 1: ggplot for problem 5.3

Problem 6: tinyclothes (8 pts)

The `tinyclothes` database stores information about customers buying products at a store. It contains four tables; `Customer` contains customer information, `SalesOrder` contains customer orders, `SalesOrderLine` is an association table between orders and products within an order, and `Product` contains product information.

The data are given below.

Customer			SalesOrder		
CustomerId	Name	Address	CustomerId	OrderId	Date
1	Alex	State	1	1	11/11/19
2	Bob	Hollister	3	2	7/9/19
3	Carol	Ocean	6	9	8/16/19
6	Juan	Phelps	6	10	10/12/19

SalesOrderLine			Product		
OrderId	ProductId	Quantity	ProductId	Name	Color
1	1	10	1	Pants	Blue
2	1	10	2	Pants	Khaki
2	4	20	3	Socks	Green
9	1	5	4	Socks	White
10	1	5	5	Shirts	White

Answer the questions on the next page.

Write SQL queries to answer the following questions. **For full credit, write a single query that matches the provided output.** Also, write only SQL and nothing else.

1. Retrieve the CustomerId, Name, Address, OrderId, and Date for customer Alex. (2 pts)

```
## CustomerId Name Address OrderId Date
## 1          1 Alex   State          1 11/11/19
```

```
# Your SQL query here
```

```
# End
```

2. Retrieve the CustomerId, Name, and total quantity of products ordered by customer Carol. (3 pts)

```
## CustomerId Name TotalQuantity
## 1          3 Carol             30
```

```
# Your SQL query here
```

```
# End
```

3. Retrieve the *unique* CustomerId, customer name, ProductId, product name, and color for all customers who ordered Blue Pants. (3 pts)

##	CustomerId	CustomerName	ProductId	ProductName	Color
## 1	1	Alex	1	Pants	Blue
## 2	3	Carol	1	Pants	Blue
## 3	6	Juan	1	Pants	Blue

Your SQL query here

End