

# Problem Set 1

## Applied Statistical Analysis 1

Maiia Skrypnyk 23371609

### Question 1 (40 points): Education

A school counselor was curious about the average of IQ of the students in her school and took a random sample of 25 students' IQ scores. The following is the data set:

```
1 y <- c(105, 69, 86, 100, 82, 111, 104, 110, 87, 108, 87, 90, 94, 113, 112, 98,
      80, 97, 95, 111, 114, 89, 95, 126, 98)
```

1. Find a 90% confidence interval for the average student IQ in the school.

```
1 #1) Calculating the mean of y — our point estimate.
2 mean(y)
3 mean <- mean(y)
4
5 #2) Calculating standard deviation with the help of R function (the
   simple way)...
6 sd(y)
7 s <- sd(y)
8 #... or calculating standard deviation by ourselves as s=sqrt(sum((y -
   mean(y))^2)/(n - 1))
9 n <- length(y)
10 s1 <- sqrt(sum((y-mean(y))^2)/(n-1))
11 #We can see that the results are identical:
12 s
13 s1
14
15 #3) Calculating the standard error using the formula se=s/sqrt(n)
16 sd(y)/sqrt(n) #n being length(y)
17 se <- sd(y)/sqrt(n)
18
19 #4) Calculating how much area do we need under the curve (if CI=90%)...
20 #4.1) ... to the right? (1-Confidence Coefficient)/2
21 (1-.90)/2
22 right <- (1-.90)/2
23 right
24 #4.2) ... to the left? (1+Confidence Coefficient)/2
```

```

25 (1+.90)/2
26 left <- (1+.90)/2
27 left
28
29 #5) Calculating Z-score.
30 #I used the Z-score table, and estimated that Z-score lies close to 1.65.
31 #But I also googled Z-score for 90% CI, and the result was 1.645, which
    is close, but still not the same.
32
33 #I found an explanation in this resource: https://www.vedantu.com/
    question-answer/z-value-for-a-90-95-and-99-percent-confidence-class
    -11-maths-cbse-606c515d034c9021d4c5b4f0
34 #” Looking for this value in the normal distribution table given below,
35 #we can see that this value [area under the curve to the left = 0.95 —
    MS]
36 #lies close to the row containing 1.6 and column containing 0.05.
37 #It also lies close to the row containing 1.6 and column containing 0.04.
38 #So, we take a mean of these values to obtain the z value at this point.
39
40 (1.64+1.65)/2
41 z <- (1.64+1.65)/2
42 z
43
44 #6) Calculating Confidence Interval via formulas 1.1)  $\text{mean}(y)+z(y)+\text{se}(=s/\sqrt{n})$ , 1.2)  $\text{mean}(y)-z(y)+\text{se}(=s/\sqrt{n})$ 
45
46 mean+(z*s/sqrt(n))
47 mean-(z*s/sqrt(n))
48
49 upper_90 <- mean+(z*s/sqrt(n))
50 lower_90 <- mean-(z*s/sqrt(n))
51
52 ConfidenceInterval90 <- c(lower_90, upper_90)
53 ConfidenceInterval90
54
55 #RESULTS: Confidence Interval [94.13244, 102.74756] =
56 #With repeated sampling, 90% of CIs will fall between these bounds:
    [94.13244, 102.74756]

```

2. Next, the school counselor was curious whether the average student IQ in her school is higher than the average IQ score (100) among all the schools in the country.

Using the same sample, conduct the appropriate hypothesis test with  $\alpha = 0.05$ .

```

1 length(y)
2 n <- length(y)
3
4 #1) Assumptions: continuous data; random sample; small sample (as n<30).
5 #As n<30, we should use a t-statistic.
6

```

```

7 #2) Null hypothesis: the average IQ score of the School A is lower than
8 # the average IQ score (100) among all the schools in the country.
9 #H0: mean(y) < 100
10 #Ha: mean(y) > 100
11
12 #3) Test statistic:  $t = (\text{mean}(y) - 100) / \text{se}$ ;  $\text{df} = (n - 1)$ ,  $\text{se} = \text{sd}(y) / \sqrt{n}$ 
13 t <- (mean(y)-100)/se
14 df <- (n-1) #24
15 df
16
17 #4) P-value: as Ha: mean(y) > 100,
18 #P = probability to the right of observed t-value (Agresti 2018, 155)
19 #As I am performing a one-tailed test (right-tailed), I am not
    multiplying pt by 2 (2*pt),
20 #but use the next formula:
21
22 p <- pt((t), df = n-1, lower.tail=FALSE) #if left-tailed, then lower.tail
    =TRUE
23
24 #I have been checking my results in multiple places, so
25 #https://www.statology.org/p-value-of-t-score-r/ suggested the next
    formula (quite the same in meaning)
26
27 p1 <- pt(q=t, df=24, lower.tail=FALSE)
28
29 #(q represents the quantile (or value) for which we want to calculate the
    cumulative probability)
30 #and I also checked the calculator here (got the same result): https://
    www.omnicalculator.com/statistics/p-value
31 #Also, not using pnorm() because we are using a t-statistic, not a z-
    statistic.
32
33 #Let's also run a t-test using t-test() formula to check our results once
    again:
34
35 t_test <- t.test(y, mu=100, alternative="greater")
36 t_test
37
38 p>.05
39 p1>.05
40
41 #As our P-value is larger than Alpha level (.05), our result is NOT
    statistically significant.
42 #Therefore, we CANNOT REJECT the null hypothesis/
43 #there is NOT enough evidence to reject the null hypothesis.

```

## Question 2 (40 points): Political Economy

Researchers are curious about what affects the amount of money communities spend on addressing homelessness. The following variables constitute our data set about social welfare expenditures in the USA.

|        |  |
|--------|--|
| State  | 50 states in US  |
| Y      | per capita expenditure on shelters/housing assistance in state           |
| X1     | per capita personal income in state                                      |
| X2     | Number of residents per 100,000 that are "financially insecure" in state |
| X3     | Number of people per thousand residing in urban areas in state           |
| Region | 1=Northeast, 2= North Central, 3= South, 4=West                          |

Explore the `expenditure` data set and import data into R.

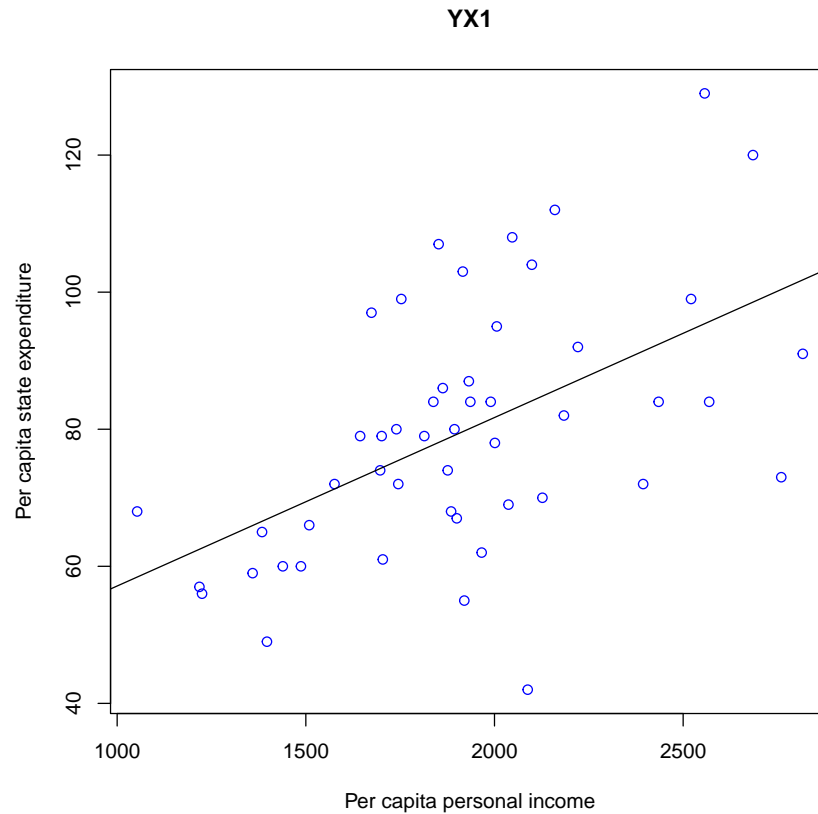
```
1 expenditure <- read.table("https://raw.githubusercontent.com/ASDS-TCD/StatsI_Fall2023/main/datasets/expenditure.txt", header=T)
```

- Please plot the relationships among `Y`, `X1`, `X2`, and `X3`? What are the correlations among them (you just need to describe the graph and the relationships among them)?

### 2.1.1: `YX1`

*(Please see the plot on the next page)*

Figure 1: Relationship between Per capita state expenditure and Per capita personal income

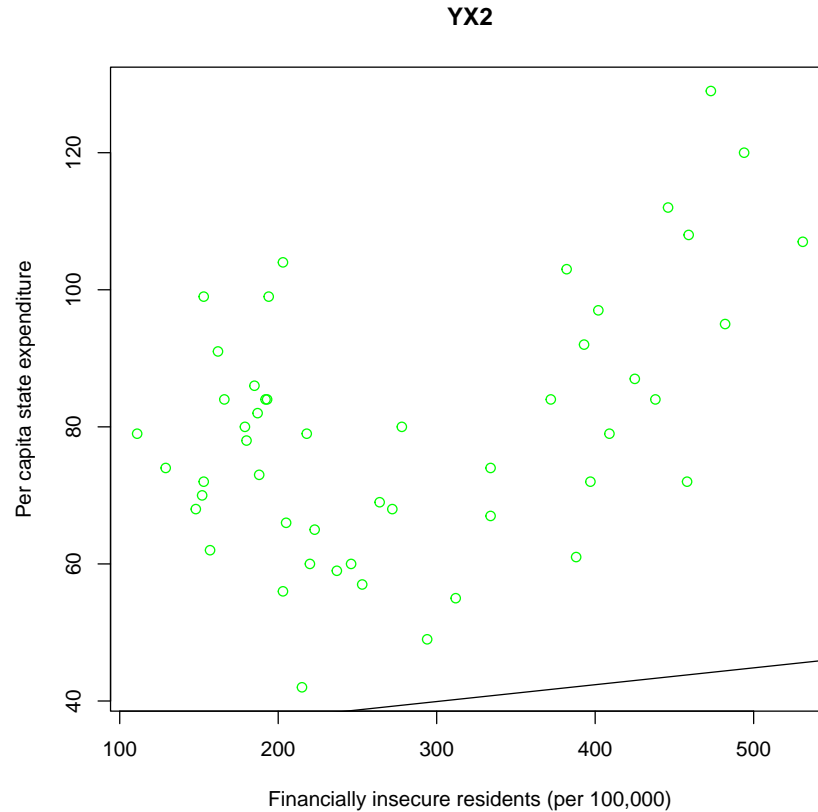


$cor = 0.5$ , positive moderate

- The correlation between Per capita expenditure on shelters/housing assistance in state and Per capita personal income in state is LINEAR, POSITIVE, MODERATE (with a few outliers).
- As Personal income increases, so does Expenditure on shelters/HA, though this relationship is not strong (+there are exceptions).

### 2.1.2: YX2

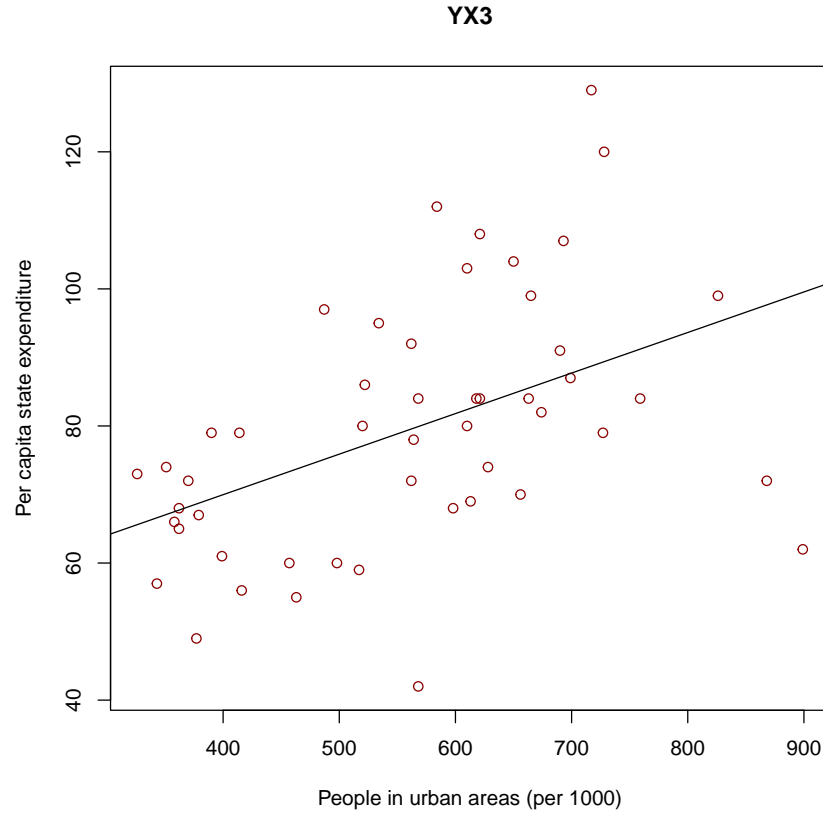
Figure 2: Relationship between Per capita state expenditure and № of Financially insecure residents in state (per 100,000)



- The correlation between Per capita expenditure on shelters/housing assistance in state and Number of residents per 100,000 that are "financially insecure" in state follows a (NON-LINEAR) U-SHAPED pattern.
- *Being honest, this is my first time hearing about a U-shaped relationship, so I searched for the answers on the Internet to understand it: "U-shaped relationship... usually means that the relationship is first decreasing and then increasing, or vice versa. In other words, it means that the relationship is not monotonic (non-monotonic), but instead has exactly one extremum (maximum or minimum)."*
- From my understanding, there is an initial range where changes in Number of "financially insecure" residents do not have a significant effect on Per capita state expenditure, but starts to increase its effect after a 'turning point' in the bottom of the U-shape.

### 2.1.3: YX3

Figure 3: Relationship between Per capita state expenditure and № of Urban residents (per 1000)

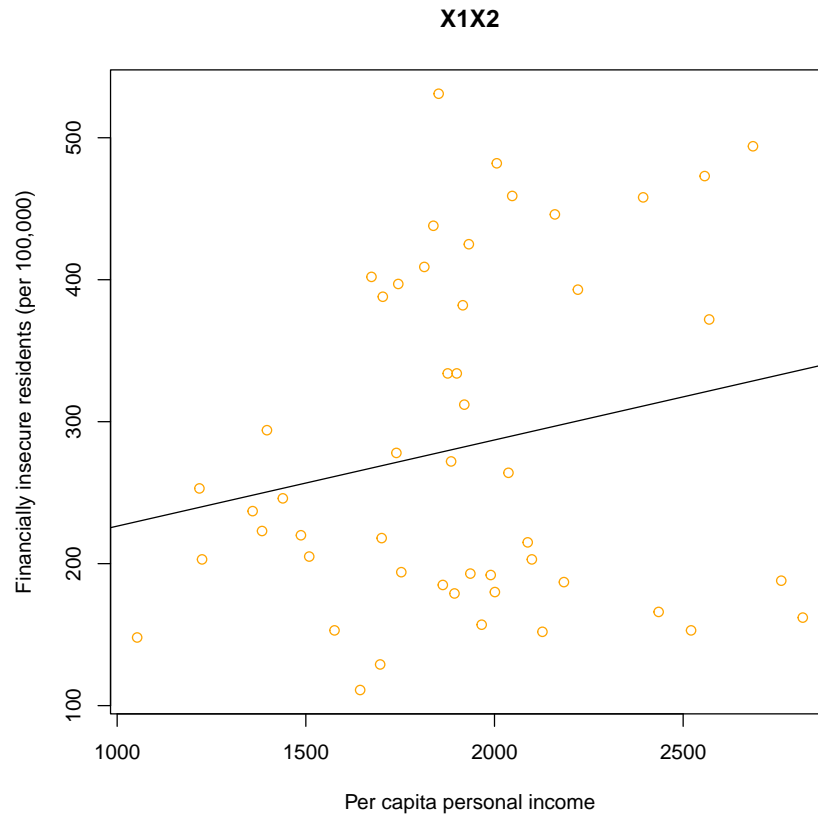


$cor = 0.46$ , positive moderate

- The correlation between Per capita expenditure on shelters/housing assistance in state and Number of people per thousand residing in urban areas in state is LINEAR, POSITIVE, MODERATE (with a few outliers).
- As the Number of urban residents increases, Per capita state expenditure on shelters/HA tends to increase as well, but this correlation is not strong.

### 2.1.4: X1X2

Figure 4: Relationship between № of Financially insecure residents in state (per 100,000) and Per capita personal income in state



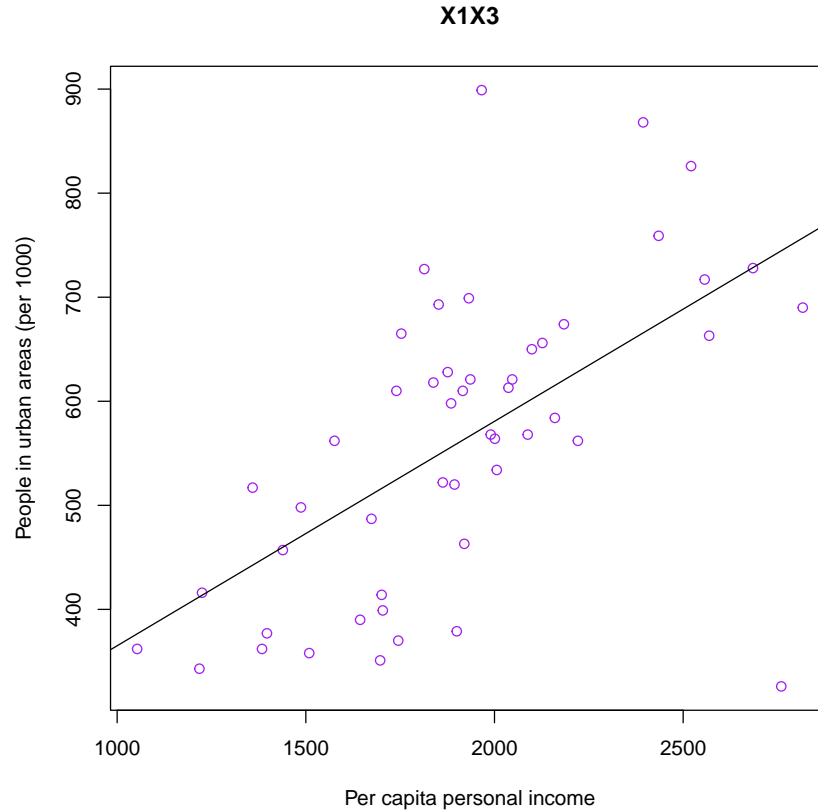
$cor = 0.2$ , *positive weak*

- The correlation between Per capita personal income in state and Number of residents per 100,000 that are "financially insecure" in state is LINEAR, POSITIVE, WEAK.
- As Per capita personal income in state increases, Number of Financially insecure residents in that state tends to decrease, but this tendency is quite weak, and there are probably other factors that influence this variable more strongly.



### 2.1.5:X1X3

Figure 5: Relationship between Per capita personal income in state and № of Urban residents (per 1000).

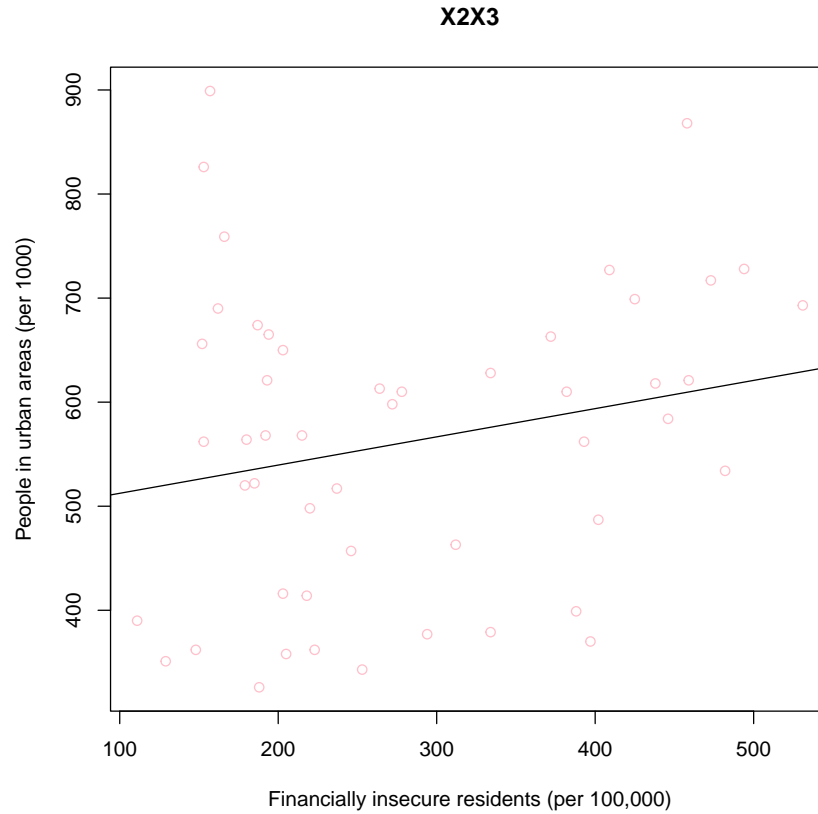


$cor = 0.6$ , positive moderate

- The correlation between Per capita personal income in state and Number of people per thousand residing in urban areas in state is LINEAR, POSITIVE, MODERATE.
- As Per capita personal income in state increases, Number of urban dwellers tends to increase as well. However, the correlation is not very strong(though the strongest of all our cases), and other factors' influence may play its role.

### 2.1.6:X2X3

Figure 6: Relationship between № of Financially insecure residents in state (per 100,000) and № of Urban residents (per 1000).

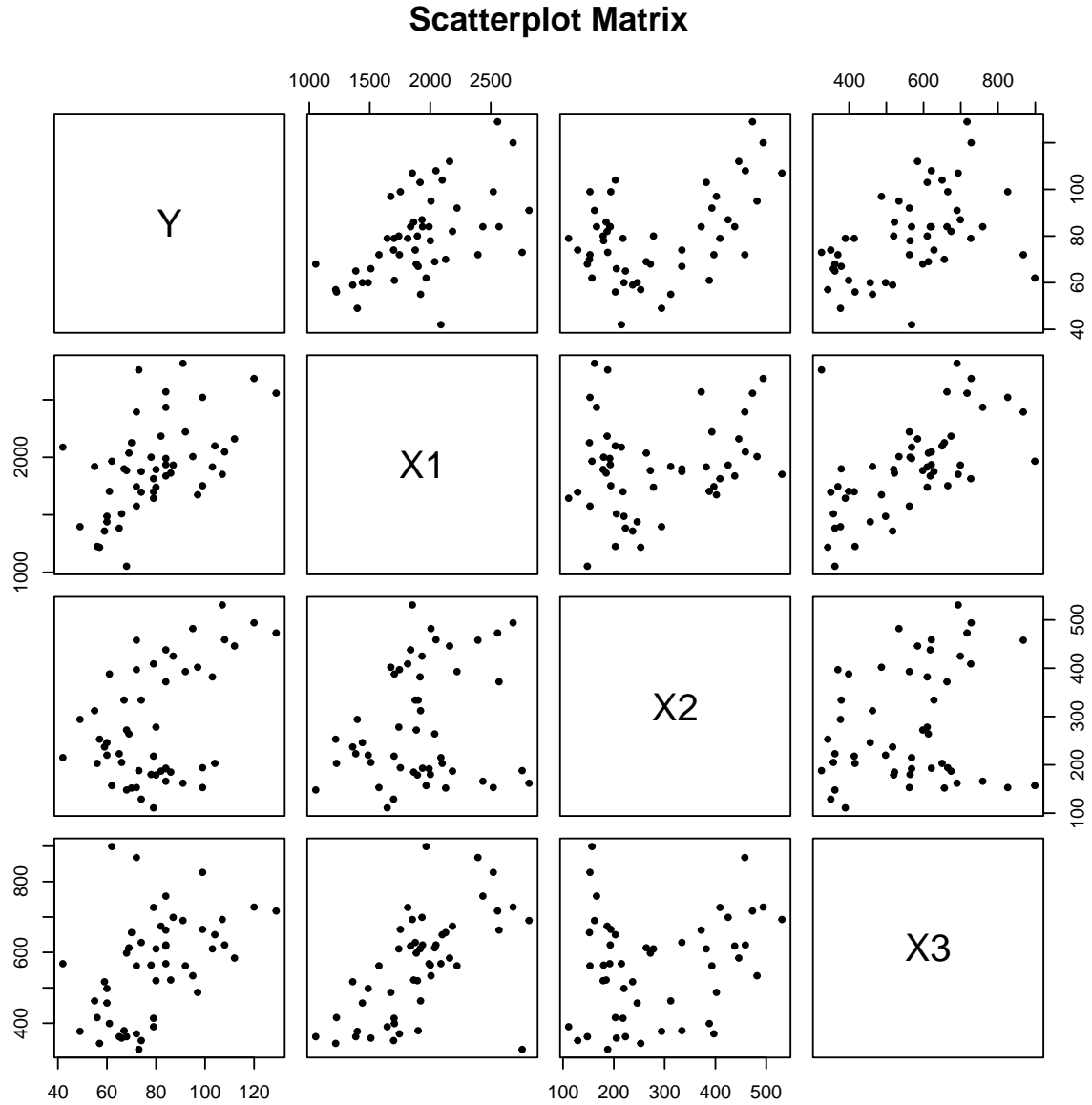


$cor = 0.2$ , positive weak

- The correlation between Number of residents per 100,000 that are "financially insecure" in state and Number of people per thousand residing in urban areas in state is LINEAR, POSITIVE, WEAK.
- As Financial Insecurity increases, Number of Urban dwellers also tends to increase, but this tendency is weak, and there are probably other factors that influence this variable more strongly.

### 2.1.7: Scatterplot Matrix

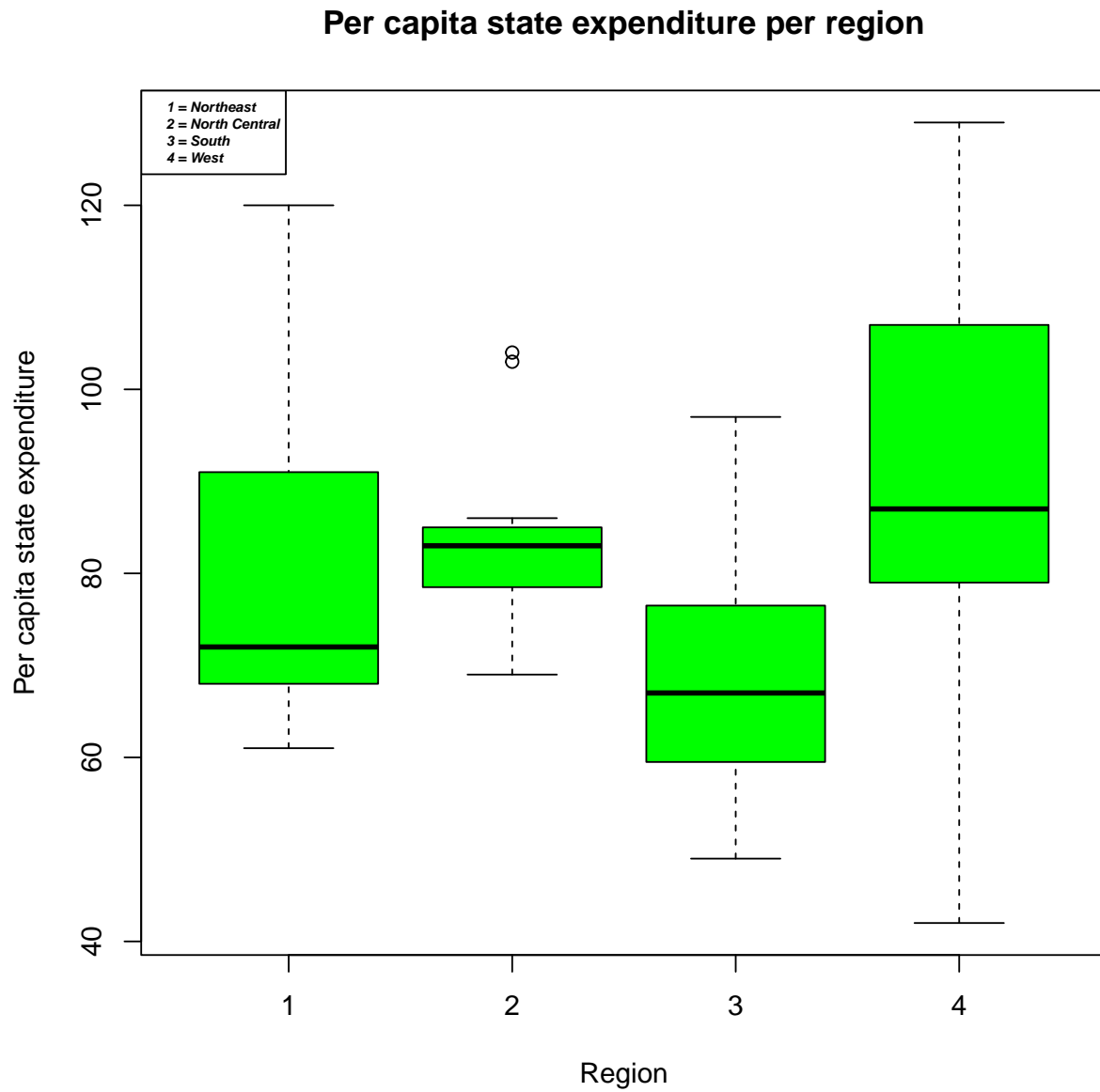
Figure 7: Scatterplot Matrix (Y, X1, X2, X3)



- Please plot the relationship between  $Y$  and  $Region$ ? On average, which region has the highest per capita expenditure on housing assistance?

## 2.2: $YRegion$

Figure 8: Relationship between Per capita expenditure on shelters/housing assistance in state and Region



```

1 #We can calculate our mean per region manually, by subsetting the dataframe
  first...:
2
3 north_east <- expenditure[expenditure$Region == 1,]
4 north_central <- expenditure[expenditure$Region == 2,]
5 south <- expenditure[expenditure$Region == 3,]
6 west <- expenditure[expenditure$Region == 4,]
7
8 #... and using mean() function to calculate average Per capita state
  expenditure
9 #for each of the four regions.
10
11 mean1 <- mean(north_east$Y)
12 mean2 <- mean(north_central$Y)
13 mean3 <- mean(south$Y)
14 mean4 <- mean(west$Y)
15
16 #The Internet suggested I use the aggregate() function
17 #( '...allows you to specify a dataframe, a condition, and a function to apply
18 #to each group of rows that meet the condition' — https://saturncloud.io/blog/aggregate-dataframe-by-condition-in-r-a-comprehensive-guide)
19 #which is a much simpler way and returns the same results.
20
21 mean_y_by_region <- aggregate(Y~Region, data = expenditure, FUN = base::mean)
22 mean_y_by_region
23
24 max(mean_y_by_region)

```

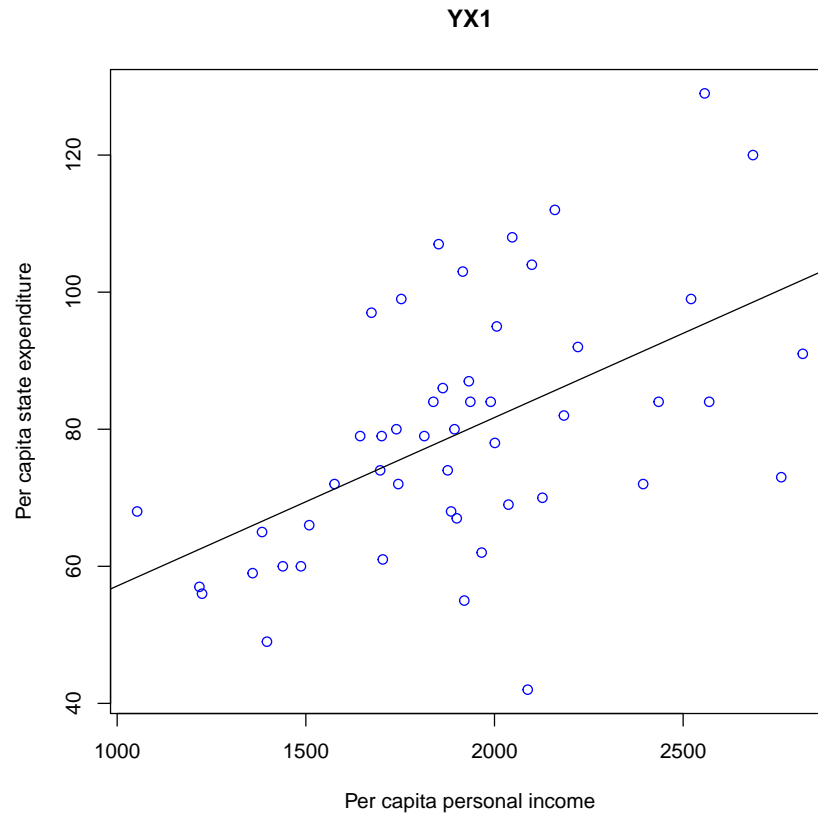
Therefore, REGION 4=WEST, on average, has the highest per capita expenditure on housing assistance (88.30769).

- Please plot the relationship between  $Y$  and  $X1$ ? Describe this graph and the relationship. Reproduce the above graph including one more variable *Region* and display different regions with different types of symbols and colors.

### 2.3.1: $YX1$

*(Please see the plot on the next page)*

Figure 9: Relationship between Per capita state expenditure and Per capita personal income



$cor = 0.5$ , positive moderate

- The correlation between Per capita expenditure on shelters/housing assistance in state and Per capita personal income in state is LINEAR, POSITIVE, MODERATE (with a few outliers).
- As Personal income increases, so does Expenditure on shelters/HA, though this relationship is not strong (+there are exceptions).

### 2.3.2: YX1 by Region

Figure 10: Relationship between Per capita state expenditure and Per capita personal income by Region

