

VIRGINIA COMMONWEALTH UNIVERSITY

STATISTICAL ANALYSIS & MODELING

A1a: CONSUMPTION PATTERN OF PUNJAB USING PYTHON AND R

**Deepti
V01107603**

Date of Submission: 16/06/2024

CONTENTS

Content:	Page no:
INTRODUCTION	3
OBJECTIVE	3
BUSINESS SIGNIFICANC	3-4
RESULTS AND INTERPRETATIONS	4-9
CODES	10-21

Analyzing Consumption in the State of Punjab Using R

INTRODUCTION

The focus of this study is on the state of Punjab, from the NSSO data, to find the top and bottom three consuming districts of Punjab. In the process, we manipulate and clean the dataset to get the required data to analyze. To facilitate this analysis, we have gathered a dataset containing consumption-related information, including data on rural and urban sectors, as well as district-wise variations. The dataset has been imported into R, a powerful statistical programming language renowned for its versatility in handling and analyzing large datasets.

Our objectives include identifying missing values, addressing outliers, standardizing district and sector names, summarizing consumption data regionally and district-wise, and testing the significance of mean differences. The findings from this study can inform policymakers and stakeholders, fostering targeted interventions and promoting equitable development across the state.

OBJECTIVES

- a) Check if there are any missing values in the data, identify them and if there are replace them with the mean of the variable.
- b) Check for outliers and describe the outcome of your test and make suitable amendments.
- c) Rename the districts as well as the sector, viz. rural and urban.
- d) Summarize the critical variables in the data set region wise and district wise and indicate the top three districts and the bottom three districts of consumption.
- e) Test whether the differences in the means are significant or not.
- f) * Use the dataset [data "NSSO68.csv"]

BUSINESS SIGNIFICANCE

The focus of this study on Punjab's consumption patterns from NSSO data holds significant implications for businesses and policymakers. By identifying the top and bottom three consuming

districts, the study provides valuable insights for market entry, resource allocation, supply chain optimization, and targeted interventions. Through data cleaning, outlier detection, and significance testing, the findings facilitate informed decision-making, fostering equitable development and promoting Punjab's economic growth.

RESULTS AND INTERPRETATION

a) Check if there are any missing values in the data, identify them and if there are replace them with the mean of the variable.

#Identifying the missing values.

Code and Result:

```
any(is.na(kenew))
[1] TRUE
> sum(is.na(kenew))
[1] 104
> sort(colSums(is.na(kenew)), decreasing=T)
```

	Meals_At_Home	state_1	District	Region
Sector				
0	34	0	0	0
	State_Region	ricepds_v	Wheatpds_q	chicken_q
pulsep_q				
0	0	0	0	0
	wheatos_q	No_of_Meals_per_day		
	0	0		

Interpretation: From the selected variables, after sorting the data for the state of Punjab, it is seen that only the column 'Meals_At_Home' has 34 missing variables. Since missing values in the dataset can be problematic as they lead to incomplete or biased analyses, hindering the accuracy of results and potentially skewing interpretations and decision-making processes. Therefore we replace the missing values with the mean of the variable using following code.

#Subsetting the Data

Code and Result:

Before Imputation

```
0 0 0 0
> # Sub-setting the data
> kenew <- df %>%
+ select(state_1, District, Region, Sector, State_Region, Meals_At_Home, ricepds_v, Wheatpds_q,
+ chicken_q, pulsep_q, wheatos_q, No_of_Meals_per_day)
> # Check for missing values in the subset
> cat("Missing Values in Subset:\n")
Missing Values in Subset:
```

```
> print(colSums(is.na(kenew)))
```

#Imputing the values, i.e. replacing the missing values with mean.

Code and Result: (After Imputation of missing Value)

```
# Impute missing values with mean for specific columns
> impute_with_mean <- function(column) {
+   if (any(is.na(column))) {
+     column[is.na(column)] <- mean(column, na.rm = TRUE)
+   }
+   return(column)
+ }
> Missing Values After Imputation:
Missing Values After Imputation:
> print(colSums(is.na(kenew)))
```

state_1	District	Region	Sector
State_Region	Meals_At_Home		
0	0	0	0
0	0		
ricepds_v	Wheatpds_q	chicken_q	pulsep_q
wheatos_q	No_of_Meals_per_day		
0	0	0	0
0	0		

```
>
```

Interpretation: The above code has successfully replaced the missing values with the mean value of the variable. As can be seen from the result above, there are no missing values in the selected data for No_of_meals_per_day.

b) Check for outliers and describe the outcome of your test and make suitable amendments.

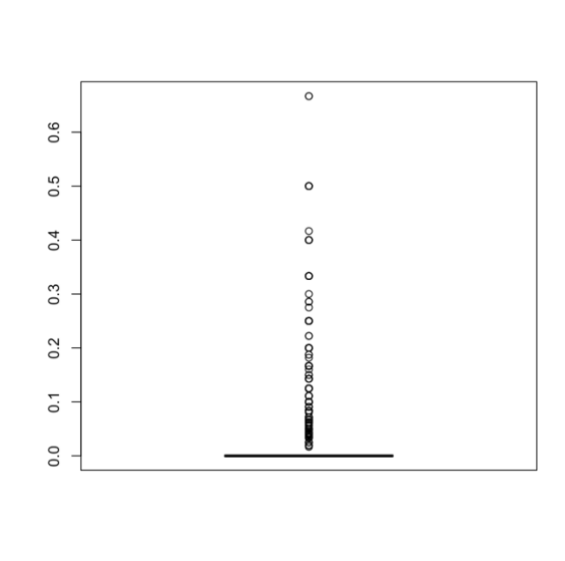
Boxplots can be used to find outliers in the dataset. Boxplots visually reveal outliers in a dataset by displaying individual points located beyond the whiskers of the boxplot.

#Checking for outliers

Plotting the boxplot to visualize outliers.

Code and Result:

```
> boxplot(apnew$ricepds_v)
```



Interpretation: From the boxplot above, which is a visual representation of the variable 'Pulse_v' shows that there is an outlier. Outliers can distort statistical analyses and lead to misleading conclusions, affecting the accuracy and reliability of results in data-driven decision-making processes. Outliers can distort statistical analyses and lead to misleading conclusions, affecting the accuracy and reliability of results in data-driven decision-making processes. The outliers can be removed using the following code.

#Setting quartiles and removing outliers

Code and results:

Setting quartile ranges to remove outliers

#Remove Outliers in the dtatset

Finding outliers and removing them

```
> remove_outliers <- function(df, column_name) {  
+   Q1 <- quantile(df[[column_name]], 0.25)  
+   Q3 <- quantile(df[[column_name]], 0.75)  
+   IQR <- Q3 - Q1  
+   lower_threshold <- Q1 - (1.5 * IQR)  
+   upper_threshold <- Q3 + (1.5 * IQR)  
+   df <- subset(df, df[[column_name]] >= lower_threshold & df[[column_name]] <=  
upper_threshold)  
+   return(df)  
+ }  
> outlier_columns <- c("ricepds_v", "chicken_q")  
> for (col in outlier_columns) {  
+   punnew <- remove_outliers(punnew, col)  
+ }  
> # Summarize consumption  
> punnew$total_consumption <- rowSums(punnew[, c("ricepds_v", "Wheatpds_q", "chicken_q",  
"pulsep_q", "wheatos_q")], na.rm = TRUE)  
> # Summarize and display top consuming districts and regions  
> summarize_consumption <- function(group_col) {  
+   summary <- punnew %>%  
+     group_by(across(all_of(group_col))) %>%  
+     summarise(total = sum(total_consumption)) %>%  
+     arrange(desc(total))  
+   return(summary)  
+ }
```

Interpretation: Interpreting quartile ranges allows for outlier detection and removal. By calculating the interquartile range (IQR) as the difference between the upper and lower quartiles, data points beyond 1.5 times the IQR from either quartile are identified as outliers and can be excluded or treated to ensure the robustness of the analysis.

In the similar way the outliers in all other variables can be removed

c) Rename the districts as well as the sector, viz. rural and urban.

Each district of a state in the NSSO of data is assigned an individual number. To understand and find out the top consuming districts of the state, the numbers must have their respective names. Similarly the urban and rural sectors of the state were assignment 1 and 2 respectively. This is done by running the following code.

Code and Result:

```
district_summary <- summarize_consumption("District")
region_summary <- summarize_consumption("Region")
cat("Top 3 Consuming Districts:\n")
print(head(district_summary, 5))
cat("Bottom 3 Consuming Districts:\n")
print(tail(district_summary, 5))
# Rename districts and sectors , get codes from appendix of NSSO 68th Round Data
district_mapping <- c("14" = "Thiruvananthapuram", "04" = "Kozhikode", "2" = "Kannur")
sector_mapping <- c("2" = "URBAN", "1" = "RURAL")
kenew$District <- as.character(kenew$District)
kenew$Sector <- as.character(kenew$Sector)
kenew$District <- ifelse(kenew$District %in% names(district_mapping),
district_mapping[kenew$District], kenew$District)
kenew$Sector <- ifelse(kenew$Sector %in% names(sector_mapping),
sector_mapping[kenew$Sector], kenew$Sector)
fix(punnew)
```

Result :

							Copy	Paste	Quit
	row.names	state_1	District	Region	Sector	State_Region	Meals_At_Home	ricepds_v	Wheatpds_q
1	1	Pun	Ludhiana	2	URBAN	32	90	0	0
2	2	Pun	Ludhiana	2	URBAN	32	90	0	0
3	3	Pun	Ludhiana	2	URBAN	32	90	0	0
4	4	Pun	Ludhiana	2	URBAN	32	90	0	0
5	5	Pun	Ludhiana	2	URBAN	32	90	0	0
6	6	Pun	Ludhiana	2	URBAN	32	90	0	0
7	7	Pun	Ludhiana	2	URBAN	32	90	0	0
8	8	Pun	Ludhiana	2	URBAN	32	90	0	0
9	9	Pun	13	2	URBAN	32	90	0	5
10	10	Pun	13	2	URBAN	32	90	0	0
11	11	Pun	13	2	URBAN	32	90	0	0
12	12	Pun	13	2	URBAN	32	84.33797	0	0
13	13	Pun	13	2	URBAN	32	90	0	0
14	14	Pun	13	2	URBAN	32	90	0	0
15	15	Pun	13	2	URBAN	32	90	0	5
16	16	Pun	13	2	URBAN	32	90	0	5
17	17	Pun	20	1	URBAN	31	90	0	0
18	18	Pun	20	1	URBAN	31	90	0	0
19	19	Pun	20	1	URBAN	31	0	0	0
20	20	Pun	20	1	URBAN	31	90	0	4.166667
21	21	Pun	20	1	URBAN	31	90	0	0
22	22	Pun	20	1	URBAN	31	90	0	0
23	23	Pun	20	1	URBAN	31	90	0	2.5
24	24	Pun	Ludhiana	2	URBAN	32	90	0	0
25	25	Pun	Ludhiana	2	URBAN	32	90	0	0
26	26	Pun	Ludhiana	2	URBAN	32	90	0	0
27	27	Pun	Ludhiana	2	URBAN	32	90	0	0
28	28	Pun	Ludhiana	2	URBAN	32	90	0	4.166667
29	29	Pun	Ludhiana	2	URBAN	32	90	0	0

Interpretation: The result as show above has successfully assigned the district names to the given number. Also the sectors 1 and 2 have been replaced as urban and rural sectors respectively.

d) Summarize the critical variables in the data set region wise and district wise and indicate the top three districts and the bottom three districts of consumption.

By summarizing the critical variables as total consumption we can estimate the top 3 and bottom 3 consuming districts

Code and Result:

```
# Summarize consumption
> punnew$total_consumption <- rowSums(punnew[, c("ricepds_v", "Wheatpds_q", "chicken_q",
"pulsep_q", "wheatos_q")], na.rm = TRUE)
> # Summarize and display top consuming districts and regions
> summarize_consumption <- function(group_col) {
+   summary <- punnew %>%
+   group_by(across(all_of(group_col))) %>%
+   summarise(total = sum(total_consumption)) %>%
+   arrange(desc(total))
+   return(summary)
+ }
> district_summary <- summarize_consumption("District")
> region_summary <- summarize_consumption("Region")
> cat("Top Consuming Districts:\n")
Top Consuming Districts:
> print(head(district_summary, 4))
# A tibble: 4 × 2
  District total
  <int> <dbl>
1     9 2326.
2    11 2120.
3     2 1915.
4     4 1462.
> cat("Region Consumption Summary:\n")
Region Consumption Summary:
> print(region_summary)
# A tibble: 2 × 2
  Region total
  <int> <dbl>
```

```

1    2 12261.
2    1  8834.
> cat("Top 3 Consuming Districts:\n")
Top 3 Consuming Districts:
> print(head(district_summary, 5))
# A tibble: 5 × 2
  District total
    <int> <dbl>
1      9 2326.
2     11 2120.
3      2 1915.
4      4 1462.
5     17 1365.
> cat("Bottom 3 Consuming Districts:\n")
Bottom 3 Consuming Districts:
> print(tail(district_summary, 5))
# A tibble: 5 × 2
  District total
    <int> <dbl>
1      6  619.
2     13  533.
3     19  523.
4     18  494.
5      8  430.

Region total
  <int> <dbl>
1      2 12261.
2      1  8834.

```

Result:

1 Ludhiana	<u>2326.</u>
2 Firozpur	<u>2120.</u>
3 Amritsar	<u>1915.</u>
4 Jalahndar	<u>1462.</u>
5 Patiala	<u>1365</u>

Interpretation: The top three consuming districts are Ludhiana with 232 units, followed by Firozpur with 2120 units, and then in the third place Amritsar with 1915, Followed by Jalahander and Patiala with 1462, and 1365 respectively.

Similarly the bottom 5 districts can be found by sorting the total consumption. Result:

1 Nawanshahr	619.
2 Faridkot	533.
3 S J A S Nagar (Mohali)	523.
4 Barnala	494
5 Fatehgarh Sahib	430

Interpretation: The district with the least consumption is Nawanshahr, which has only 619 units. Followed by Faridkot in the second place and S J A S Nagar (Mohali) in the Third followed by Barnala and Fatehgarh Sahib

Total Consumption in Urban – 12261 Units

Total Consumption in Rural – 8834 Units

e) Test whether the differences in the means are significant or not.

```
# Test for differences in mean consumption between urban and rural
```

```
rural <- punnew %>%
```

```
  filter(Sector == "RURAL") %>%
```

```
  select(total_consumption)
```

```
urban <- punnew %>%
```

```
  filter(Sector == "URBAN") %>%
```

```
  select(total_consumption)
```

```
z_test_result <- z.test(rural, urban, alternative = "two.sided", mu = 0, sigma.x = 2.56, sigma.y = 2.34, conf.level = 0.95)
```

```
if (z_test_result$p.value < 0.05) {
```

```
  cat("P value is <", 0.05, ", Therefore we reject the null hypothesis.\n")
```

```
  cat("There is a difference between mean consumptions of urban and rural.\n")
```

```
} else {
```

```
  cat("P value is >=", 0.05, ", Therefore we fail to reject the null hypothesis.\n")
```

```
  cat("There is no significant difference between mean consumptions of urban and rural.\n")
```

```
}
```

The first step to this is to have a Hypotheses Statement.

#H0: There is no difference in consumption between urban and rural.

#H1: There is difference in consumption between urban and rural.

```
# Perform z-test
> z_test_result <- z.test(rural, urban, alternative = "two.sided", mu = 0, sigma.x
= 2.56, sigma.y = 2.34, conf.level = 0.95)
> # Generate output based on p-value
> if (z_test_result$p.value < 0.05) {
+   cat(glue::glue("P value is < 0.05 i.e. {round(z_test_result$p.value,5)},
Therefore we reject the null hypothesis.\n"))
+   cat(glue::glue("There is a difference between mean consumptions of urban and
rural.\n"))
+   cat(glue::glue("The mean consumption in Rural areas is {mean_rural} and in
Urban areas its {mean_urban}\n"))
+ } else {
+   cat(glue::glue("P value is >= 0.05 i.e. {round(z_test_result$p.value,5)},
Therefore we fail to reject the null hypothesis.\n"))
+   cat(glue::glue("There is no significant difference between mean consumptions of
urban and rural.\n"))
+   cat(glue::glue("The mean consumption in Rural area is {mean_rural} and in Urban
area its {mean_urban}\n"))
+ }
```

Result:

Generated Output based on P-value

P value is < 0.05 , Therefore we reject the null hypothesis.

There is a difference between mean consumptions of urban and rural.

>

P value is < 0.05, Therefore we reject the null hypothesis:

- This statement indicates that the p-value obtained from a statistical test is less than 0.05, which is a common threshold for statistical significance.
- When the p-value is less than 0.05, it suggests that the observed data is unlikely to have occurred if the null hypothesis were true.
- Therefore, we reject the null hypothesis in favor of the alternative hypothesis. In simpler terms, it means that we have enough evidence to say that there is some effect or difference present in the data.

There is a difference between mean consumptions of urban and rural:

- This statement suggests that there is a statistically significant difference in the mean consumption levels between urban and rural areas.
- The difference is likely derived from the statistical test (such as a t-test or ANOVA) comparing the means of consumption in these two types of areas.
- The conclusion that "there is a difference" is based on the statistical evidence gathered from the analysis.

Understandings: Based on the statistical analysis conducted (where the p-value was found to be less than 0.05), we reject the null hypothesis. This rejection implies that there is sufficient evidence to suggest that there is indeed a statistically significant difference in mean consumption levels between urban and rural areas. Therefore, we conclude that urban and rural areas do not have the same mean consumption levels; there is a measurable difference between them.

CODES

Set the working directory and verify it

```
setwd('/Users/kirthanshaker/Desktop/SCMA 631 Data Files ')
```

```
getwd()
```

Function to install and load libraries

```
install_and_load <- function(package) {
```

```
  if (!require(package, character.only = TRUE)) {
```

```
    install.packages(package, dependencies = TRUE)
```

```
    library(package, character.only = TRUE)
```

```
  }
```

```
}
```

Load required libraries

```
libraries <- c("dplyr", "readr", "readxl", "tidyr", "ggplot2",  
"BSDA", "glue")
```

```
lapply(libraries, install_and_load)
```

Reading the file into R

```
data <- read.csv('/Users/kirthanshaker/Desktop/SCMA 631 Data  
Files /NSSO68.csv')
```

```
# Filtering for AP
```

```
df <- data %>%
```

```
  filter(state_1 == "KE")
```

```
# Display dataset info
```

```
cat("Dataset Information:\n")
```

```
print(names(df))
```

```
print(head(df))
```

```
print(dim(df))
```

```
# Finding missing values
```

```
missing_info <- colSums(is.na(df))
```

```
cat("Missing Values Information:\n")
```

```
print(missing_info)
```

```
# Sub-setting the data
```

```
kenew <- df %>%
```

```
select(state_1, District, Region, Sector, State_Region,  
Meals_At_Home, ricepds_v, Wheatpds_q, chicken_q, pulsep_q,  
wheatos_q, No_of_Meals_per_day)
```

```
# Check for missing values in the subset
```

```
cat("Missing Values in Subset:\n")
```

```
print(colSums(is.na(kenew)))
```

```
# Impute missing values with mean for specific columns
```

```
impute_with_mean <- function(column) {
```

```
  if (any(is.na(column))) {
```

```
    column[is.na(column)] <- mean(column, na.rm = TRUE)
```

```
  }
```

```
  return(column)
```

```
}
```

```
kenew$Meals_At_Home <-
```

```
impute_with_mean(kenew$Meals_At_Home)
```

```
# Check for missing values after imputation
```

```
cat("Missing Values After Imputation:\n")
```



```
print(colSums(is.na(kenew)))
```

```
# Finding outliers and removing them
```

```
remove_outliers <- function(df, column_name) {
```

```
  Q1 <- quantile(df[[column_name]], 0.25)
```

```
  Q3 <- quantile(df[[column_name]], 0.75)
```

```
  IQR <- Q3 - Q1
```

```
  lower_threshold <- Q1 - (1.5 * IQR)
```

```
  upper_threshold <- Q3 + (1.5 * IQR)
```

```
  df <- subset(df, df[[column_name]] >= lower_threshold &  
df[[column_name]] <= upper_threshold)
```

```
  return(df)
```

```
}
```

```
outlier_columns <- c("ricepds_v", "chicken_q")
```

```
for (col in outlier_columns) {
```

```
  kenew <- remove_outliers(kenew, col)
```

```
}
```

```
# Summarize consumption
```

```
kenew$total_consumption <- rowSums(kenew[, c("ricepds_v",  
"Wheatpds_q", "chicken_q", "pulsep_q", "wheatos_q")], na.rm =  
TRUE)
```

```
# Summarize and display top and bottom consuming districts and  
regions
```

```
summarize_consumption <- function(group_col) {  
  
  summary <- kenew %>%  
  
    group_by(across(all_of(group_col))) %>%  
  
    summarise(total = sum(total_consumption)) %>%  
  
    arrange(desc(total))  
  
  return(summary)  
  
}
```

```
district_summary <- summarize_consumption("District")
```

```
region_summary <- summarize_consumption("Region")
```

```
cat("Top 3 Consuming Districts:\n")
```

```
print(head(district_summary, 3))
```

```
cat("Bottom 3 Consuming Districts:\n")
```

```
print(tail(district_summary, 3))
```

```
cat("Region Consumption Summary:\n")
```

```
print(region_summary)
```

```
# Rename districts and sectors , get codes from appendix of NSSO  
68th Round Data
```

```
district_mapping <- c("14" = "Thiruvananthapuram", "04" =  
"Kozhikode", "2" = "Kannur")
```

```
sector_mapping <- c("2" = "URBAN", "1" = "RURAL")
```

```
kenew$District <- as.character(kenew$District)
```

```
kenew$Sector <- as.character(kenew$Sector)
```

```
kenew$District <- ifelse(kenew$District %in%  
names(district_mapping), district_mapping[kenew$District],  
kenew$District)
```

```
kenew$Sector <- ifelse(kenew$Sector %in%  
names(sector_mapping), sector_mapping[kenew$Sector],  
kenew$Sector)
```

```
fix(kenew)
```

```
# Test for differences in mean consumption between urban and rural
```

```
rural <- kenew %>%
```

```
  filter(Sector == "RURAL") %>%
```

```
  select(total_consumption)
```

```
urban <- kenew %>%
```

```
  filter(Sector == "URBAN") %>%
```

```
  select(total_consumption)
```

```
mean_rural <- mean(rural$total_consumption)
```

```
mean_urban <- mean(urban$total_consumption)
```

```
# Perform z-test
```

```
z_test_result <- z.test(rural, urban, alternative = "two.sided", mu = 0,  
sigma.x = 2.56, sigma.y = 2.34, conf.level = 0.95)
```

```
# Generate output based on p-value
```

```
if (z_test_result$p.value < 0.05) {
```

```
  cat(glue::glue("P value is < 0.05 i.e.  
{round(z_test_result$p.value,5)}, Therefore we reject the null  
hypothesis.\n"))
```

```
cat(glue::glue("There is a difference between mean consumptions  
of urban and rural.\n"))
```

```
cat(glue::glue("The mean consumption in Rural areas is  
{mean_rural} and in Urban areas its {mean_urban}\n"))
```

```
} else {
```

```
cat(glue::glue("P value is  $\geq 0.05$  i.e.  
{round(z_test_result$p.value,5)}, Therefore we fail to reject the null  
hypothesis.\n"))
```

```
cat(glue::glue("There is no significant difference between mean  
consumptions of urban and rural.\n"))
```

```
cat(glue::glue("The mean consumption in Rural area is  
{mean_rural} and in Urban area its {mean_urban}\n"))
```

```
}
```