

```

---
title: "R Notebook"
output: html_notebook
---

# DSC520
# Assignment Week 5
# Zemelak Goraga
# 2024-1-13"

```

Set the working directory to the correct path
setwd("C:\\Users\\MariaStella\\Downloads\\wk5_house")
```

```R
Install and load required packages
install.packages("dplyr")
install.packages("tidyr")
install.packages("magrittr")
install.packages("purrr")
```

```R
Load the installed packages
library(dplyr)
library(tidyr)
library(magrittr)
library(purrr)
```

```R
Import the dataset scores.csv as df
df <- read.csv("C:\\Users\\MariaStella\\Downloads\\wk5_house\\house2.csv")
```

```R
Display the first few rows of the dataset
head(df)
```


	sale_date	sale_price	sale_reason	sale_instrument	sale_warning	sitetype
1	1/3/2006	698000	1	3		R1
2	1/3/2006	649990	1	3		R1
3	1/3/2006	572500	1	3		R1
4	1/3/2006	420000	1	3		R1
5	1/3/2006	369900	1	3	15	R1
6	1/3/2006	184667	1	15	18 51	R1

  


	addr_full	zip5	ctyname	postalctyn	lon	lat	building_grade
1	17021 NE 113TH CT	98052	REDMOND	REDMOND	-122.1124	47.70139	9
2	11927 178TH PL NE	98052	REDMOND	REDMOND	-122.1022	47.70731	9
3	13315 174TH AVE NE	98052		REDMOND	-122.1085	47.71986	8
4	3303 178TH AVE NE	98052	REDMOND	REDMOND	-122.1037	47.63914	8
5	16126 NE 108TH CT	98052	REDMOND	REDMOND	-122.1242	47.69748	7
6	8101 229TH DR NE	98053		REDMOND	-122.0341	47.67545	7

  


	square_feet_total_living	bedrooms	bath_full_count	bath_half_count
1	2810	4	2	1
2	2880	4	2	0
3	2770	4	1	1
4	1620	3	1	0
5	1440	3	1	0
6	4160	4	2	1

  


	bath_3qtr_count	year_built	year_renovated	current_zoning	sq_ft_lot	prop_type
--	-----------------	------------	----------------	----------------	-----------	-----------


```

| | | | | | | |
|---|---|------|---|-------|------|---|
| 1 | 0 | 2003 | 0 | R4 | 6635 | R |
| 2 | 1 | 2006 | 0 | R4 | 5570 | R |
| 3 | 1 | 1987 | 0 | R6 | 8444 | R |
| 4 | 1 | 1968 | 0 | R4 | 9600 | R |
| 5 | 1 | 1980 | 0 | R6 | 7526 | R |
| 6 | 1 | 2005 | 0 | URPSO | 7280 | R |

```
present_use
```

| | |
|---|---|
| 1 | 2 |
| 2 | 2 |
| 3 | 2 |
| 4 | 2 |
| 5 | 2 |
| 6 | 2 |

```
```R
Display column names using the names() function
column_names <- names(df)
print(column_names)
```
```

| | | |
|------|-------------------|----------------------------|
| [1] | "sale_date" | "sale_price" |
| [3] | "sale_reason" | "sale_instrument" |
| [5] | "sale_warning" | "sitetype" |
| [7] | "addr_full" | "zip5" |
| [9] | "ctyname" | "postalctyn" |
| [11] | "lon" | "lat" |
| [13] | "building_grade" | "square_feet_total_living" |
| [15] | "bedrooms" | "bath_full_count" |
| [17] | "bath_half_count" | "bath_3qtr_count" |
| [19] | "year_built" | "year_renovated" |
| [21] | "current_zoning" | "sq_ft_lot" |
| [23] | "prop_type" | "present_use" |

```
```R
The 'sale_price' variable was used to perform the under mentioned operations (Task 1
to Task4).
null values can be removed from the 'sale_price' data as follows.
sale_price_without_nulls <- na.omit(df$sale_price)
print(head(sale_price_without_nulls))
```
```

```
[1] 698000 649990 572500 420000 369900 184667
```

```
# Task 1
```

a. Using the dplyr package, use the 6 different operations to analyze/transform the data - GroupBy, Summarize, Mutate, Filter, Select, and Arrange.

```
```R
dplyr operations using the 'sale_price' variable.

GroupBy and Summarize
grouped_summary <- df %>%
 dplyr::group_by(ctyname) %>%
 dplyr::summarize(
 Total_Sale_Price = sum(`sale_price`),
 Average_Square_Feet = mean(square_feet_total_living)
)
```

```

Mutate
mutated_df <- df %>%
 dplyr::mutate(Sale_Price_Double = `sale_price` * 2)

Filter
filtered_df <- df %>%
 dplyr::filter(bedrooms == 4 & bath_full_count == 2)

Select
selected_columns <- df %>%
 dplyr::select(ctyname, `sale_price`, square_feet_total_living)

Arrange
arranged_df <- df %>%
 dplyr::arrange(desc(year_built))

Other similar operations:
Count distinct values in ctyname
distinct_count <- df %>%
 dplyr::count(ctyname, name = "City_Count")

Calculate cumulative sum of `sale_price` within each ctyname
cumulative_sum <- df %>%
 dplyr::arrange(ctyname, `sale_price`) %>%
 dplyr::group_by(ctyname) %>%
 dplyr::mutate(Cumulative_Sale_Price = cumsum(`sale_price`))

Print the results
cat("GroupBy and Summarize:\n")
print(grouped_summary)

cat("Mutate:\n")
head(mutated_df)

cat("Filter:\n")
head(filtered_df)

cat("Select:\n")
head(selected_columns)

cat("Arrange:\n")
head(arranged_df)

cat("Count Distinct:\n")
print(distinct_count)

cat("Cumulative Sum:\n")
head(cumulative_sum)
` ``

```

Output:

```

GroupBy and Summarize:
> print(grouped_summary)

```

```

A tibble: 3 × 3
 ctyname Total_Sale_Price Average_Square_Feet
 <chr> <dbl> <dbl>
1 "" 4102485158 2612.
2 "REDMOND" 4333722291 2461.
3 "SAMMAMISH" 64183700 3788.
>

```

Mutate:

```
> head(mutated_df)
```

	sale_date	sale_price	sale_reason	sale_instrument	sale_warning	sitetype
1	1/3/2006	698000	1	3		R1
2	1/3/2006	649990	1	3		R1
3	1/3/2006	572500	1	3		R1
4	1/3/2006	420000	1	3		R1
5	1/3/2006	369900	1	3	15	R1
6	1/3/2006	184667	1	15	18 51	R1

  

	addr_full	zip5	ctyname	postalctyn	lon	lat	building_grade
1	17021 NE 113TH CT	98052	REDMOND	REDMOND	-122.1124	47.70139	9
2	11927 178TH PL NE	98052	REDMOND	REDMOND	-122.1022	47.70731	9
3	13315 174TH AVE NE	98052		REDMOND	-122.1085	47.71986	8
4	3303 178TH AVE NE	98052	REDMOND	REDMOND	-122.1037	47.63914	8
5	16126 NE 108TH CT	98052	REDMOND	REDMOND	-122.1242	47.69748	7
6	8101 229TH DR NE	98053		REDMOND	-122.0341	47.67545	7

  

	square_feet_total_living	bedrooms	bath_full_count	bath_half_count
1	2810	4	2	1
2	2880	4	2	0
3	2770	4	1	1
4	1620	3	1	0
5	1440	3	1	0
6	4160	4	2	1

  

	bath_3qtr_count	year_built	year_renovated	current_zoning	sq_ft_lot	prop_type
1	0	2003	0	R4	6635	R
2	1	2006	0	R4	5570	R
3	1	1987	0	R6	8444	R
4	1	1968	0	R4	9600	R
5	1	1980	0	R6	7526	R
6	1	2005	0	URPSO	7280	R

  

	present_use	Sale_Price_Double
1	2	1396000
2	2	1299980
3	2	1145000
4	2	840000
5	2	739800
6	2	369334

>

Filter:

```
> head(filtered_df)
```

	sale_date	sale_price	sale_reason	sale_instrument	sale_warning	sitetype
1	1/3/2006	698000	1	3		R1
2	1/3/2006	649990	1	3		R1
3	1/3/2006	184667	1	15	18 51	R1
4	1/4/2006	875000	1	3		R1
5	1/4/2006	660000	1	3		R1
6	1/6/2006	765000	1	3		R1

  

	addr_full	zip5	ctyname	postalctyn	lon	lat	building_grade
1	17021 NE 113TH CT	98052	REDMOND	REDMOND	-122.1124	47.70139	9
2	11927 178TH PL NE	98052	REDMOND	REDMOND	-122.1022	47.70731	9
3	8101 229TH DR NE	98053		REDMOND	-122.0341	47.67545	7
4	21404 NE 67TH ST	98053		REDMOND	-122.0555	47.66510	10
5	7525 238TH AVE NE	98053		REDMOND	-122.0227	47.67208	9
6	8944 237TH PL NE	98053		REDMOND	-122.0230	47.68150	9

  

	square_feet_total_living	bedrooms	bath_full_count	bath_half_count
1	2810	4	2	1
2	2880	4	2	0
3	4160	4	2	1
4	3720	4	2	1
5	4160	4	2	1

```

6 4000 4 2 1
bath_3qtr_count year_built year_renovated current_zoning sq_ft_lot prop_type
1 0 2003 0 R4 6635 R
2 1 2006 0 R4 5570 R
3 1 2005 0 URPSO 7280 R
4 0 1988 0 RA5 30649 R
5 1 1978 0 RA5 42688 R
6 1 2005 0 URPSO 7611 R
present_use
1 2
2 2
3 2
4 2
5 2
6 2
>

```

Select:

```
> head(selected_columns)
```

```

ctyname sale_price square_feet_total_living
1 REDMOND 698000 2810
2 REDMOND 649990 2880
3 572500 2770
4 REDMOND 420000 1620
5 REDMOND 369900 1440
6 184667 4160
>

```

```
> cat("Arrange:\n")
```

Arrange:

```
> head(arranged_df)
```

```

sale_date sale_price sale_reason sale_instrument sale_warning sitetype
1 3/28/2006 270000 1 3 R1
2 12/4/2006 562000 1 3 R1
3 12/17/2007 2300000 1 3 45 R1
4 7/21/2011 120527 1 15 18 22 R1
5 12/20/2012 1536000 1 22 45 R1
6 6/3/2013 302500 1 3 R1
addr_full zip5 ctyname postalctyn lon lat building_grade
1 5806 249TH CT NE 98053 REDMOND -122.0053 47.65706 11
2 16326 NE 43RD CT 98052 REDMOND -122.1219 47.64860 11
3 9113 258TH AVE NE 98053 REDMOND -121.9957 47.68123 10
4 17132 NE 80TH ST 98052 REDMOND -122.1103 47.67558 7
5 9113 258TH AVE NE 98053 REDMOND -121.9957 47.68123 10
6 11320 244TH AVE NE 98053 REDMOND -122.0129 47.69933 10
square_feet_total_living bedrooms bath_full_count bath_half_count
1 5060 4 23 1
2 5040 5 2 1
3 5100 4 3 1
4 940 2 1 1
5 5100 4 3 1
6 5820 3 2 1
bath_3qtr_count year_built year_renovated current_zoning sq_ft_lot prop_type
1 0 2016 0 RA5 89734 R
2 3 2016 0 R4 11761 R
3 0 2016 0 RA5 131301 R
4 0 2016 0 R5 12230 R
5 0 2016 0 RA5 131301 R
6 1 2016 0 RA10P 436507 R
present_use
1 0
2 2
3 2
4 2
5 2

```

```
6 300
>
```

```
Count Distinct:
> print(distinct_count)
```

```
 ctyname City_Count
1 6078
2 REDMOND 6721
3 SAMMAMISH 66
>
```

```
Cumulative Sum:
> head(cumulative_sum)
```

```
A tibble: 6 × 25
Groups: ctyname [1]
 sale_date sale_price sale_reason sale_instrument sale_warning sitetype
 <chr> <int> <int> <int> <chr> <chr>
1 7/6/2010 698 1 26 24 R1
2 7/6/2010 698 1 26 24 R1
3 12/29/2009 873 1 26 24 R1
4 1/28/2010 873 1 26 24 32 R1
5 12/22/2009 998 1 26 24 R1
6 3/20/2007 1000 1 15 18 51 R1
```

```
Task 2
```

Using the purrr package - perform 2 functions on your dataset. You could use zip\_n, keep, discard, compact, etc.

```
```R
# Use purrr functions on the df dataset
# Function 1: Use zip_n to combine `sale_price` and square_feet_total_living columns
into a list
combined_cols <- df %>%
  dplyr::mutate(combined = purrr::pmap(list(`sale_price`, square_feet_total_living),
c))

# Function 2: Use filter to retain rows where bedrooms is greater than 3
filtered_df <- df %>%
  dplyr::filter(bedrooms > 3)

# Function 3: Use discard to remove rows where bath_full_count is 0
filtered_df_no_zero_bath <- df %>%
  dplyr::filter(bath_full_count != 0)

# Function 4: Use compact to remove NULL elements from a list
compact_list <- list(a = 1, b = NULL, c = 3) %>%
  purrr::compact()
```

```
# Display the results, the first few rows of each
cat("Combined Columns:\n")
print(head(combined_cols$combined))

cat("Filtered Rows (bedrooms > 3):\n")
print(head(filtered_df))

cat("Filtered Rows (bath_full_count not 0):\n")
print(head(filtered_df_no_zero_bath))

cat("Compact List:\n")
print(head(compact_list))
```

```

Output:

```
Combined Columns:
> print(head(combined_cols$combined))
[[1]]
[1] 698000 2810

[[2]]
[1] 649990 2880

[[3]]
[1] 572500 2770

[[4]]
[1] 420000 1620

[[5]]
[1] 369900 1440

[[6]]
[1] 184667 4160

>
```

```
Filtered Rows (bedrooms > 3):
> print(head(filtered_df))
```

|   | sale_date | sale_price | sale_reason | sale_instrument | sale_warning | sitetype |
|---|-----------|------------|-------------|-----------------|--------------|----------|
| 1 | 1/3/2006  | 698000     | 1           | 3               |              | R1       |
| 2 | 1/3/2006  | 649990     | 1           | 3               |              | R1       |
| 3 | 1/3/2006  | 572500     | 1           | 3               |              | R1       |
| 4 | 1/3/2006  | 184667     | 1           | 15              | 18 51        | R1       |
| 5 | 1/4/2006  | 1050000    | 1           | 3               |              | R1       |
| 6 | 1/4/2006  | 875000     | 1           | 3               |              | R1       |

  

|   | addr_full          | zip5  | ctyname | postalctyn | lon       | lat      | building_grade |
|---|--------------------|-------|---------|------------|-----------|----------|----------------|
| 1 | 17021 NE 113TH CT  | 98052 | REDMOND | REDMOND    | -122.1124 | 47.70139 | 9              |
| 2 | 11927 178TH PL NE  | 98052 | REDMOND | REDMOND    | -122.1022 | 47.70731 | 9              |
| 3 | 13315 174TH AVE NE | 98052 |         | REDMOND    | -122.1085 | 47.71986 | 8              |
| 4 | 8101 229TH DR NE   | 98053 |         | REDMOND    | -122.0341 | 47.67545 | 7              |
| 5 | 21634 NE 87TH PL   | 98053 |         | REDMOND    | -122.0507 | 47.68053 | 10             |
| 6 | 21404 NE 67TH ST   | 98053 |         | REDMOND    | -122.0555 | 47.66510 | 10             |

  

|   | square_foot_total | living | bedrooms | bath_full_count | bath_half_count |
|---|-------------------|--------|----------|-----------------|-----------------|
| 1 |                   | 2810   | 4        | 2               | 1               |
| 2 |                   | 2880   | 4        | 2               | 0               |
| 3 |                   | 2770   | 4        | 1               | 1               |
| 4 |                   | 4160   | 4        | 2               | 1               |

```

5 3960 5 3 0
6 3720 4 2 1
 bath_3qtr_count year_built year_renovated current_zoning sq_ft_lot prop_type
1 0 2003 0 R4 6635 R
2 1 2006 0 R4 5570 R
3 1 1987 0 R6 8444 R
4 1 2005 0 URPSO 7280 R
5 1 1993 0 RA5 97574 R
6 0 1988 0 RA5 30649 R
 present_use
1 2
2 2
3 2
4 2
5 2
6 2
>

```

```

Filtered Rows (bath_full_count not 0):
> print(head(filtered_df_no_zero_bath))

```

```

 sale_date sale_price sale_reason sale_instrument sale_warning sitetype
1 1/3/2006 698000 1 3 R1
2 1/3/2006 649990 1 3 R1
3 1/3/2006 572500 1 3 R1
4 1/3/2006 420000 1 3 R1
5 1/3/2006 369900 1 3 15 R1
6 1/3/2006 184667 1 15 18 51 R1
 addr_full zip5 ctynome postalctyn lon lat building_grade
1 17021 NE 113TH CT 98052 REDMOND REDMOND -122.1124 47.70139 9
2 11927 178TH PL NE 98052 REDMOND REDMOND -122.1022 47.70731 9
3 13315 174TH AVE NE 98052 REDMOND REDMOND -122.1085 47.71986 8
4 3303 178TH AVE NE 98052 REDMOND REDMOND -122.1037 47.63914 8
5 16126 NE 108TH CT 98052 REDMOND REDMOND -122.1242 47.69748 7
6 8101 229TH DR NE 98053 REDMOND REDMOND -122.0341 47.67545 7
 square_feet_total_living bedrooms bath_full_count bath_half_count
1 2810 4 2 1
2 2880 4 2 0
3 2770 4 1 1
4 1620 3 1 0
5 1440 3 1 0
6 4160 4 2 1
 bath_3qtr_count year_built year_renovated current_zoning sq_ft_lot prop_type
1 0 2003 0 R4 6635 R
2 1 2006 0 R4 5570 R
3 1 1987 0 R6 8444 R
4 1 1968 0 R4 9600 R
5 1 1980 0 R6 7526 R
6 1 2005 0 URPSO 7280 R
 present_use
1 2
2 2
3 2
4 2
5 2
6 2
>

```

```

Task 3

```



Use the cbind and rbind function on your dataset

```
```R
# Use cbind and rbind functions on the df dataset
# Create a new data frame to demonstrate cbind
df2 <- data.frame(`sale_price` = rep(c(500000, 600000, 700000), length.out = nrow(df)),
                  square_feet_total_living = rep(c(2500, 3000, 3500), length.out =
nrow(df)),
                  ctynome = rep(c("REDMOND", "BELLEVUE", "SEATTLE"), length.out =
nrow(df)))

# Use cbind to combine df and df2 by columns
combined_by_columns <- cbind(df, df2)
cat("Combined by Columns:\n")
print(combined_by_columns)

# Use rbind to combine df and df2 by rows
combined_by_rows <- rbind(df, df2)
cat("Combined by Rows:\n")
print(head(combined_by_rows))
```
```

Output

Combined by Columns:

> print(combined\_by\_columns)

|    | sale_date | sale_price | sale_reason | sale_instrument | sale_warning | sitetype |
|----|-----------|------------|-------------|-----------------|--------------|----------|
| 1  | 1/3/2006  | 698000     | 1           | 3               |              | R1       |
| 2  | 1/3/2006  | 649990     | 1           | 3               |              | R1       |
| 3  | 1/3/2006  | 572500     | 1           | 3               |              | R1       |
| 4  | 1/3/2006  | 420000     | 1           | 3               |              | R1       |
| 5  | 1/3/2006  | 369900     | 1           | 3               | 15           | R1       |
| 6  | 1/3/2006  | 184667     | 1           | 15              | 18 51        | R1       |
| 7  | 1/4/2006  | 1050000    | 1           | 3               |              | R1       |
| 8  | 1/4/2006  | 875000     | 1           | 3               |              | R1       |
| 9  | 1/4/2006  | 660000     | 1           | 3               |              | R1       |
| 10 | 1/4/2006  | 650000     | 1           | 3               |              | R1       |
| 11 | 1/4/2006  | 599950     | 1           | 3               |              | R1       |
| 12 | 1/4/2006  | 526787     | 1           | 3               |              | R1       |
| 13 | 1/4/2006  | 470000     | 1           | 3               |              | R1       |
| 14 | 1/4/2006  | 165000     | 1           | 3               |              | R1       |
| 15 | 1/5/2006  | 803000     | 1           | 3               |              | R1       |
| 16 | 1/5/2006  | 507950     | 1           | 3               |              | R1       |
| 17 | 1/6/2006  | 765000     | 1           | 3               |              | R1       |
| 18 | 1/6/2006  | 589950     | 1           | 3               |              | R1       |
| 19 | 1/9/2006  | 501000     | 1           | 3               |              | R1       |
| 20 | 1/9/2006  | 372500     | 1           | 3               |              | R1       |
| 21 | 1/10/2006 | 513262     | 1           | 3               |              | R1       |
| 22 | 1/10/2006 | 482000     | 1           | 3               |              | R1       |
| 23 | 1/11/2006 | 765000     | 1           | 3               |              | R1       |
| 24 | 1/11/2006 | 372500     | 1           | 3               |              | R2       |
| 25 | 1/11/2006 | 265000     | 1           | 3               |              | R1       |
| 26 | 1/12/2006 | 1392000    | 1           | 3               |              | R1       |
| 27 | 1/12/2006 | 717390     | 1           | 3               |              | R1       |
| 28 | 1/12/2006 | 552000     | 1           | 3               |              | R1       |
| 29 | 1/12/2006 | 470000     | 1           | 3               |              | R1       |
| 30 | 1/13/2006 | 523935     | 1           | 3               |              | R1       |
| 31 | 1/13/2006 | 399900     | 1           | 3               |              | R1       |
| 32 | 1/13/2006 | 335105     | 1           | 3               |              | R1       |
| 33 | 1/16/2006 | 572950     | 1           | 3               |              | R1       |
| 34 | 1/17/2006 | 949950     | 1           | 3               |              | R1       |
| 35 | 1/17/2006 | 905000     | 1           | 3               | 41           | R1       |
| 36 | 1/17/2006 | 750073     | 1           | 3               |              | R1       |

```

37 1/17/2006 526718 1 3 R1
 addr_full zip5 ctynome postalctyn lon lat
1 17021 NE 113TH CT 98052 REDMOND REDMOND -122.1124 47.70139
2 11927 178TH PL NE 98052 REDMOND REDMOND -122.1022 47.70731
3 13315 174TH AVE NE 98052 REDMOND REDMOND -122.1085 47.71986
4 3303 178TH AVE NE 98052 REDMOND REDMOND -122.1037 47.63914
5 16126 NE 108TH CT 98052 REDMOND REDMOND -122.1242 47.69748
6 8101 229TH DR NE 98053 REDMOND REDMOND -122.0341 47.67545
7 21634 NE 87TH PL 98053 REDMOND REDMOND -122.0507 47.68053
8 21404 NE 67TH ST 98053 REDMOND REDMOND -122.0555 47.66510
9 7525 238TH AVE NE 98053 REDMOND REDMOND -122.0227 47.67208
10 17703 NE 26TH ST 98052 REDMOND REDMOND -122.1039 47.63341
11 14924 NE 74TH CT 98052 REDMOND REDMOND -122.1411 47.67142
12 7858 148TH CT NE 98052 REDMOND REDMOND -122.1425 47.67407
13 17905 NE 26TH ST 98052 REDMOND REDMOND -122.1010 47.63319
14 2921 288TH AVE NE 98053 REDMOND REDMOND -121.9577 47.63382
15 3624 264TH AVE NE 98053 REDMOND REDMOND -121.9857 47.64184
16 7850 148TH CT NE 98052 REDMOND REDMOND -122.1425 47.67390
17 8944 237TH PL NE 98053 REDMOND REDMOND -122.0230 47.68150

```

```

[reached 'max' / getOption("max.print") -- omitted 12828 rows]
>
> # Use rbind to combine df and df2 by rows
> combined_by_rows <- rbind(df, df2)
Error in rbind(deparse.level, ...) :
 numbers of columns of arguments do not match

```

# Task 4

Split a string, then concatenate the results back together

```

```R
# Using one variable, split a string, then concatenate
# the results back together.

# Do it in more detail by using another variable.
# Example using one variable (ctynome)
df$City_Split <- strsplit(as.character(df$ctynome), "")
df <- cbind(df, t(sapply(df$City_Split, c)))

# Display the updated dataset
head(df)

# Concatenate the results for 'ctynome'
df$City_Joined <- sapply(df$City_Split, function(x) paste(x, collapse = ""))
df <- df[, -which(colnames(df) %in% c("City_Split"))] # Remove intermediate columns

# Display the final updated dataset
head(df)
```

```

Output:

```

sale_date sale_price sale_reason sale_instrument sale_warning
sitetype addr_full zip5 ctynome postalctyn ... 12856 12857
12858 12859 12860 12861 12862 12863 12864 12865

```

```

<fct> <int> <int> <int> <fct> <fct> <fct> <int> <fct> <fct> ...
<list> <list> <list> <list> <list> <list> <list> <list> <list> <list>
1 1/3/2006 698000 1 3 R1 17021 NE 113TH CT
98052 REDMOND REDMOND ... R, E, D, M, O, N, D R, E, D, M, O,
N, D R, E, D, M, O, N, D R, E, D, M, O, N, D R, E, D, M, O,
N, D R, E, D, M, O, N, D
2 1/3/2006 649990 1 3 R1 11927 178TH PL NE
98052 REDMOND REDMOND ... R, E, D, M, O, N, D R, E, D, M, O,
N, D R, E, D, M, O, N, D R, E, D, M, O, N, D R, E, D, M, O,
N, D R, E, D, M, O, N, D
3 1/3/2006 572500 1 3 R1 13315 174TH AVE NE
98052 REDMOND ... R, E, D, M, O, N, D R, E, D, M, O,
N, D R, E, D, M, O, N, D R, E, D, M, O, N, D R, E, D, M, O,
N, D R, E, D, M, O, N, D
4 1/3/2006 420000 1 3 R1 3303 178TH AVE NE
98052 REDMOND REDMOND ... R, E, D, M, O, N, D R, E, D, M, O,
N, D R, E, D, M, O, N, D R, E, D, M, O, N, D R, E, D, M, O,
N, D R, E, D, M, O, N, D
5 1/3/2006 369900 1 3 15 R1 16126 NE 108TH CT
98052 REDMOND REDMOND ... R, E, D, M, O, N, D R, E, D, M, O,
N, D R, E, D, M, O, N, D R, E, D, M, O, N, D R, E, D, M, O,
N, D R, E, D, M, O, N, D
6 1/3/2006 184667 1 15 18 51 R1 8101 229TH DR NE
98053 REDMOND ... R, E, D, M, O, N, D R, E, D, M, O,
N, D R, E, D, M, O, N, D R, E, D, M, O, N, D R, E, D, M, O,
N, D R, E, D, M, O, N, D

```

A data.frame: 6 × 12890

```

sale_date sale_price sale_reason sale_instrument sale_warning
sitetype addr_full zip5 ctynome postalctyn ... 12857 12858
12859 12860 12861 12862 12863 12864 12865 City_Joined
<fct> <int> <int> <int> <fct> <fct> <fct> <fct> ...
<list> <list> <list> <list> <list> <list> <list> <list> <chr>
1 1/3/2006 698000 1 3 R1 17021 NE 113TH CT
98052 REDMOND REDMOND ... R, E, D, M, O, N, D R, E, D, M, O,
N, D R, E, D, M, O, N, D R, E, D, M, O, N, D R, E, D, M, O,
N, D REDMOND
2 1/3/2006 649990 1 3 R1 11927 178TH PL NE
98052 REDMOND REDMOND ... R, E, D, M, O, N, D R, E, D, M, O,
N, D R, E, D, M, O, N, D R, E, D, M, O, N, D R, E, D, M, O,
N, D REDMOND
3 1/3/2006 572500 1 3 R1 13315 174TH AVE NE
98052 REDMOND ... R, E, D, M, O, N, D R, E, D, M, O,
N, D R, E, D, M, O, N, D R, E, D, M, O, N, D R, E, D, M, O,
N, D REDMOND
4 1/3/2006 420000 1 3 R1 3303 178TH AVE NE
98052 REDMOND REDMOND ... R, E, D, M, O, N, D R, E, D, M, O,
N, D R, E, D, M, O, N, D R, E, D, M, O, N, D R, E, D, M, O,
N, D REDMOND
5 1/3/2006 369900 1 3 15 R1 16126 NE 108TH CT
98052 REDMOND REDMOND ... R, E, D, M, O, N, D R, E, D, M, O,
N, D R, E, D, M, O, N, D R, E, D, M, O, N, D R, E, D, M, O,
N, D REDMOND
6 1/3/2006 184667 1 15 18 51 R1 8101 229TH DR NE
98053 REDMOND ... R, E, D, M, O, N, D R, E, D, M, O,
N, D R, E, D, M, O, N, D R, E, D, M, O, N, D R, E, D, M, O,
N, D REDMOND

```

# Discussion

This R code focuses on data manipulation tasks using the 'housing.xlsx' dataset, which was first saved as 'house2.csv' in excel, and then imported and saved as df dataset in

R. The exercise was divided into four tasks, each employing different packages and functions to analyze, transform, and combine the data.

#### Task 1: Using dplyr operations

**GroupBy and Summarize:** The dataset is grouped by city name (ctyname), and summary statistics like total sale price and average square feet are calculated for each group.

**Mutate:** A new column, Sale\_Price\_Double, is created by doubling the 'sale\_price' variable.

**Filter:** Rows are filtered based on specific conditions, such as having 4 bedrooms and 2 full bathrooms.

**Select:** Specific columns (ctyname, sale\_price, and square\_feet\_total\_living) are selected for further analysis.

**Arrange:** The dataset is arranged in descending order based on the 'year\_built' variable.

**Other Operations:** Additional operations include counting distinct values in ctyname and calculating cumulative sums within each city.

#### Task 2: Using purrr functions on the df dataset

**Function 1 (zip\_n):** Utilizes pmap to combine 'sale\_price' and 'square\_feet\_total\_living' into a list.

**Function 2 (filter):** Retains rows where the number of bedrooms is greater than 3.

**Function 3 (discard):** Removes rows where the count of full bathrooms is 0.

**Function 4 (compact):** Eliminates NULL elements from a list.

#### Task 3: Using cbind and rbind functions on the df dataset

Two new data frames (df2) are created to demonstrate cbind and rbind.

**cbind:** Combines the original dataset (df) and df2 by columns.

**rbind:** Combines the original dataset (df) and df2 by rows.

#### Task 4: String Manipulation

The 'ctyname' variable is split into individual characters and added as new columns to the dataset (City\_Split).

The split elements are then concatenated back into a single string (City\_Joined).

Intermediate columns are removed, resulting in a final updated dataset.

In summary, this R code showcases a comprehensive set of data manipulation tasks, including filtering, grouping, summarizing, mutating, and combining datasets, using the 'house' dataset. The explanations provided aim to make the code accessible to users with varying levels of R proficiency.