

## Repo link:

<https://github.com/tang-xin72/DSC520/tree/main/scratch/week4>

DSC520 / scratch / week4 /

tang-xin72 create sub folder	
Name	
..	
figure	
DSC520_week4_housing_TangXin.R	
DSC520_week4_housing_TangXin.log	
DSC520_week4_housing_TangXin.tex	
DSC520_week_4_Testscore_TangXin.R	
DSC520_week_4_Testscore_TangXin.log	
DSC520_week_4_Testscore_TangXin.pdf	
DSC520_week_4_Testscore_TangXin.tex	
DSC520_wk4_Tangxin-assignment(housing+Testscore).pdf	
DSC520_wk4_Tangxin-assignment.docx	

The housing\_tangXin.R can run but failed to create pdf due to memory overflow. The output is inside the DSC520\_wk4\_Tangxin\_assignment(housing\_testscore).pdf

## Exercise 1 test score:

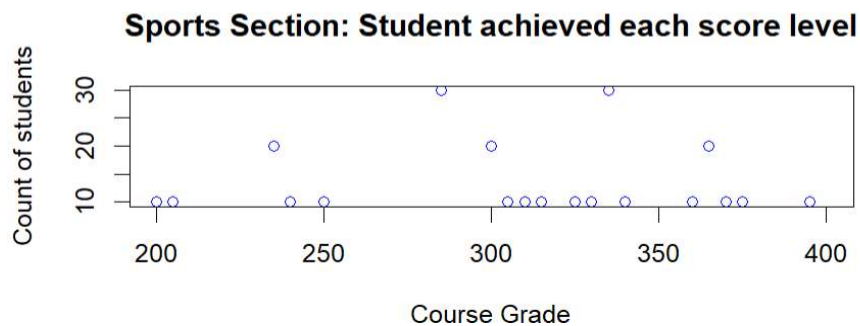
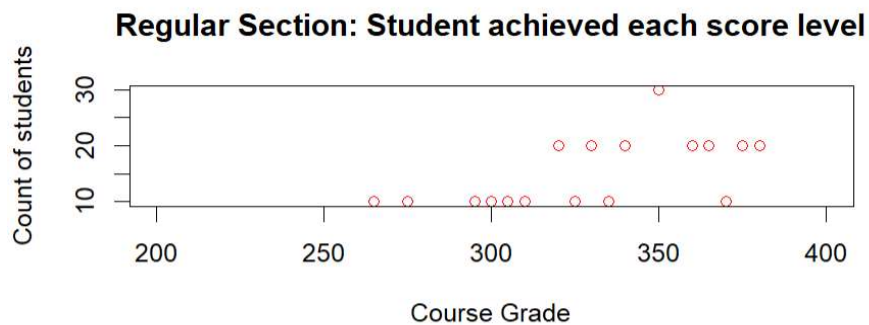
```
> data <- read.csv(file = theUrl, header = TRUE, sep = ',')
> head(data)
  Count Score Section
1     10    200 Sports
2     10    205 Sports
3     20    235 Sports
4     10    240 Sports
5     10    250 Sports
6     10    265 Regular
```

1. What are the observational units in this study?  
A: There are total 38 units of observations of the scores in the data, grouped by score. The score and count are numeric, and section is a category variable.
2. Identify the variables mentioned in the narrative paragraph and determine which are categorical and quantitative?  
A: There are 3 variables in the data. The score and count are quantitative variables, and section is a category variable.
3. Create one variable to hold a subset of your data set that contains only the Regular Section and one variable for the Sports Section.  
A: see R code in program or the screen shot here.

```
regular_data <- subset(data, Section=='Regular')
```

```
sports_data <- subset(data, Section=='sports')
```

4. Use the Plot function to plot each Sections scores and the number of students achieving that score.



4.1 Can you say that one section tended to score more points than the other? Justify and explain your answer.

**A:** look at the plot, for regular section, more data points of count of 20 and 30 exist in score >300 side. For sports section, more data points of count 10 exist in score > 300 side. Also, there is no data points for score <250 in regular section, while there are some in sports section.

So, more students from regular section scored high score, especially above 350, also there are no student from regular section scored less than 250. In summary, students from regular section performance better in test.

4.2 Did every student in one section score more points than every student in the other section? If not, explain what a statistical tendency means in this context.

**A:** look at the plot, it is not true that every student from regular section outperformed all students from sport section. Some students from sport section also get high score, but it is rare. The lowest score the student from regular section gets is 265. While the lowest score in sport section is just 200.

So the mean score for regular section is higher than sport section. Statistically the students in regular section get better test score than those in sport section.

4.3 What could be one additional variable that was not mentioned in the narrative that could be influencing the point distributions between the two sections?

**A:** the data did not mention if the students are taking the same test (quiz, mid term or final test).

## Exercise 2 housing data:

See week4\_housing\_TangXin.R for coding

### New variables:

Total\_sale\_price: created by 'apply' function

Zipcode\_SalePrice\_mean: created by 'aggregate' function

The following variables created using 'plyr' package

Count\_by\_sales\_room\_avg, Count\_by\_sales\_room, Count\_by\_Sales and Housing\_by\_zip\_bedroom.

### Data analysis:

1. Houses with bedroom of 2-4 have most sales.
2. Sales price is not a normal distribution, it is positively skewed, which makes sense since most people cannot afford multimillion mansions.
3. There are a few outliers. Like one house with 7 bedroom sold for 4 million, well above its mean of 2 million. Also the sales of the houses with more than 9 rooms are also outlier in this dataset.

