# LAB: INTRO TO STAT ANALYSIS

Zeyi Qian

zeqian@clarku.edu
Office Hours: JC 201, Tuesday 4–5 PM & Thursday 3-4 PM

September 20, 2024

# Question 1

- Fast food is often considered unhealthy because much of it is high in fat and calories. Are **fat** and **calories** related? Analyze the association between fat content and calories by calculating the correlation coefficient and drawing a scatterplot (Use a computer)

| Fat (g) | 19 | 31 | 34 | 35 | 39 | 39 | 41 |
|---|---|---|---|---|---|---|---|
| Calories | 410 | 560 | 585 | 570 | 640 | 680 | 660 |

# Question 1

| Q1 | Fat (g) | 19 | 31 | 34 | 35 | 39 | 39 | 41 |
|---|---|---|---|---|---|---|---|---|
| | Calories | 410 | 560 | 585 | 570 | 640 | 680 | 660 |

=CORREL(C1:I1,C2:I2)

**Function Arguments**                                                        ?    ✕

CORREL

    **Array1**    `C1:I1`     ⬆   =   {19,31,34,35,39,39,41}

    **Array2**    `C2:I2`     ⬆   =   {410,560,585,570,640,680,660}

                                                   =   0.977514742

Returns the correlation coefficient between two data sets.

                        **Array1**    is a cell range of values. The values should be numbers, names, arrays, or references that contain numbers.

Formula result =  0.977514742

Help on this function                                                        OK         Cancel
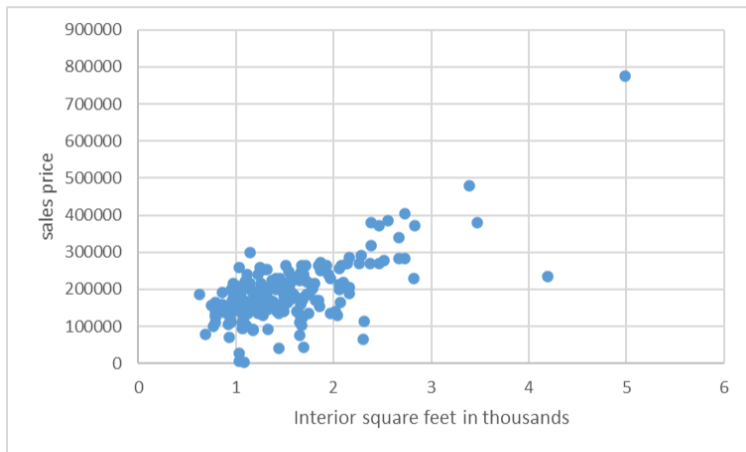
# Question 1

## Question 2

This question asks you to work with housing sales data from the 01602 zip code of Worcester, which is across Park Ave from the Clark neighborhood. All of the houses are single-family houses. Please use the Excel Spreadsheet that is on the page for this assignment.

- a. Please use the "Single Family" Tab first. Please find the scatterplot for **floor area of the house (interior square feet in thousnds)** and **sales price**. Why would it seem to make the most sense to put the floor area on the "x-axis and the sales price on the "y-axis."?

It is logical to treat sales price as something that responds to changes in floor area

# Question 2: a



- It is logical to treat sales price as something that responds to changes in floor area

# Question 2: b

- b. Please describe the scatterplot you have created according to these criteria:
    - i. direction of association (positive or negative?) **positive**
    - ii. is the form curved, straight, or "exotic?" **straight (linear)**
    - iii. is the strength of association apparently strong or weak? **strong**

# Question 2: c

- c. Using the data from the "Single Family" tab, calculate the correlation coefficient by hand and then check your work using the excel command. **(0.688743)**

$$Z_{x_i} = \frac{x_i - \bar{x}}{Stdev_x}, \ Z_{y_i} = \frac{y_i - \bar{y}}{Stdev_y}$$

$$r = \sum_{i=1}^{n} \frac{Z_{x_i} Z_{y_i}}{n-1}$$

- Order of calculation by hand
  - Mean of x & y
  - Deviation from x & deviation from y
  - Stdev of x & y
  - Z score of x & y
  - Correlation

## Question 2: d

- d. Using the data on the "Three Family" tab, you will find three-family houses that were sold in the adjacent 01603 zip code that is in South Worcester further down Main Street away from downtown. Please find the same kind of scatter plot that you did in part a., but now it will be for three-family houses and describe it using the criteria in part b.:
  - i. direction of association (positive or negative?)
  - ii. is the form curved, straight, or "exotic?"
  - iii. is the strength of association apparently strong or weak?

# Question 2: d

- Positive
- Straight, though there is a considerable spread in the data points
- (Moderate to) Strong
- Other factors may also play a significant role, leading to the observed variability

# Question 2: e & f

- e. Using the data on the "Three Family" tab, test the regression assumptions
  - Linearity Assumption
  - Independence Assumption
  - Equal Variance Assumption
  - Normal Population Assumption
- f. Estimate a regression model using Excel. What is the interpretation of the coefficients (i.e., put the coefficients in a sentence)?

# Question 2: e & f

# Question 2: e & f

# Question 2: e & f

- Independence Assumption



interior square feet in thousands Residual Plot

# Question 2: e & f

- Linearity Assumption **Put y into new sheet**

# Question 2: e & f

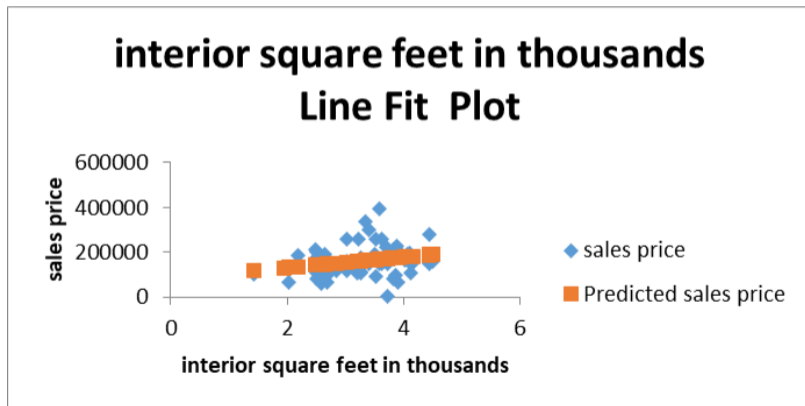- Linearity Assumption

# Question 2: e & f

- Equal Variance Assumption

# Question 2: e & f

- Normal Population Assumption

# Question 2: e & f

# Question 2: e & f

# Question 2: e & f

- All else being equal, for each additional 1,000 square feet of house area, the expected increase in the sales price is $23,318.94

# Question 3

For this question, please use the data from the student survey on social media use ("socialmedia&sleep" Tab).

- a. Find the scatter diagram for social media posts per day and hours of sleep per week.
- b. Find the correlation coefficient using excel.
- c. Suppose that we multiplied the "posts per day" variable by 7 so that it was measured in terms of weeks (rather than days). Would that influence the correlation coefficient (r)? How?
- d. Calculate the mean and standard deviation of the two variables. Then calculate the regression slope and regression intercept by hand. Write out the regression model.
- e. Estimate the regression model using Excel. Do the answers in d. and e. match?

# Question 3: a & b & c



- r=CORREL(A2:A40,B2:B40)=-0.477231403
- If we multiplied the "posts per day" variable by 7, r =-0.477231403 still holds

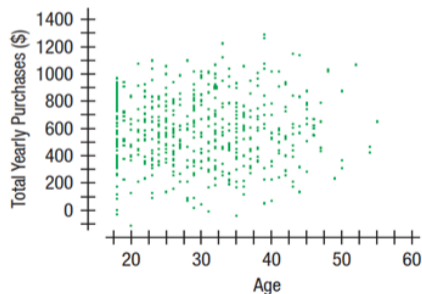## Question 3: d & e

- hours of sleep per week $= b0 + b1$ social media posts per day
- mean of x $=$ AVERAGE(A2:A40)$=1.59$
- mean of y $=$ AVERAGE(B2:B40)$=49.76$
- Stdev of x $=$ STDEV(A2:A40)$=2.59$
- Stdev of y $=$ STDEV(B2:B40)$=9.53$
- b1$=$r Sy/Sx $=$ SLOPE(B2:B40,A2:A40)$=1.75$
- b0$=$mean of y - b1 mean of x $=$ INTERCEPT(B2:B40,A2:A40)$=52.55$
- **Data Analysis** will give you same answwaer

# Question 4

Online clothes  An online clothing retailer keeps track of its
customers' purchases. For those customers who signed up for the
company's credit card, the company also has information on
the customer's *Age* and *Income*. A random sample of 500 of
these customers shows the following scatterplot of *Total Yearly
Purchases* by *Age*:



The correlation between *Total Yearly Purchases* and *Age* is
$r = 0.037$. Summary statistics for the two variables are:

# Question 4

| | Mean | SD |
|---|---|---|
| **Age** | 29.67 yr | 8.51 yr |
| **Total Yearly Purchase** | $572.52 | $253.62 |

a) What is the linear regression equation for predicting *Total Yearly Purchase* from *Age*?

b) Do the assumptions and conditions for regression appear to be met?

c) What is the predicted *Total Yearly Purchase* for an 18-year-old? For a 50-year-old?

d) What percent of the variability in *Total Yearly Purchases* is accounted for by this model?

e) Do you think the regression might be a useful one for the company? Explain.

# Question 4

- Total Yearly Purchases $= b0 + b1$ Age
- $b1 = r\ Sy/Sx = 11.37$
- $b0 =$ mean of y - b1 mean of x $= 268.56$
- Total Yearly Purchases $= 268.56 + 11.37 \times$ Age
- $R2 = 0.001369$
- If age $= 18$, predicted y $= 268.56 + 11.37 \times 18 = 473.22$
- If age $= 50$, predicted y $= 268.56 + 11.37 \times 50 = 837.06$