

# Mini Project

Members: Elise Cheng, Stephen Dong, Selena Arias, Shashwat Gupta

## 1. What data do you have?

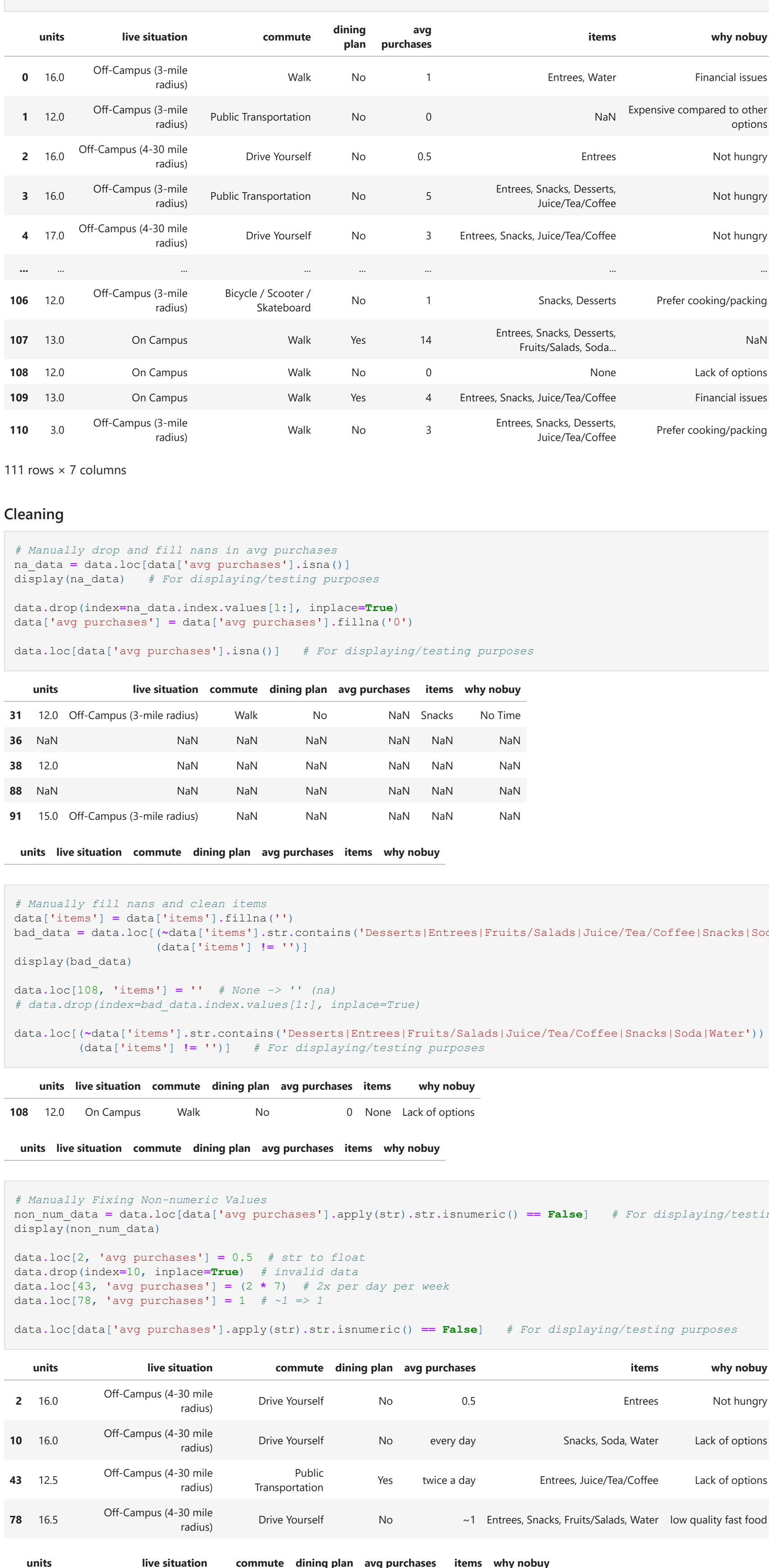
We have data on the following questions from the survey:

- How many units are you taking currently?
- What is your current living situation?
- How do you commute to school?
- Do you have a UCR dining plan?
- How many times a week do you purchase food or drinks from somewhere on campus?
- What items do you purchase?
- What is the biggest reason you do not purchase more food and drinks on campus?

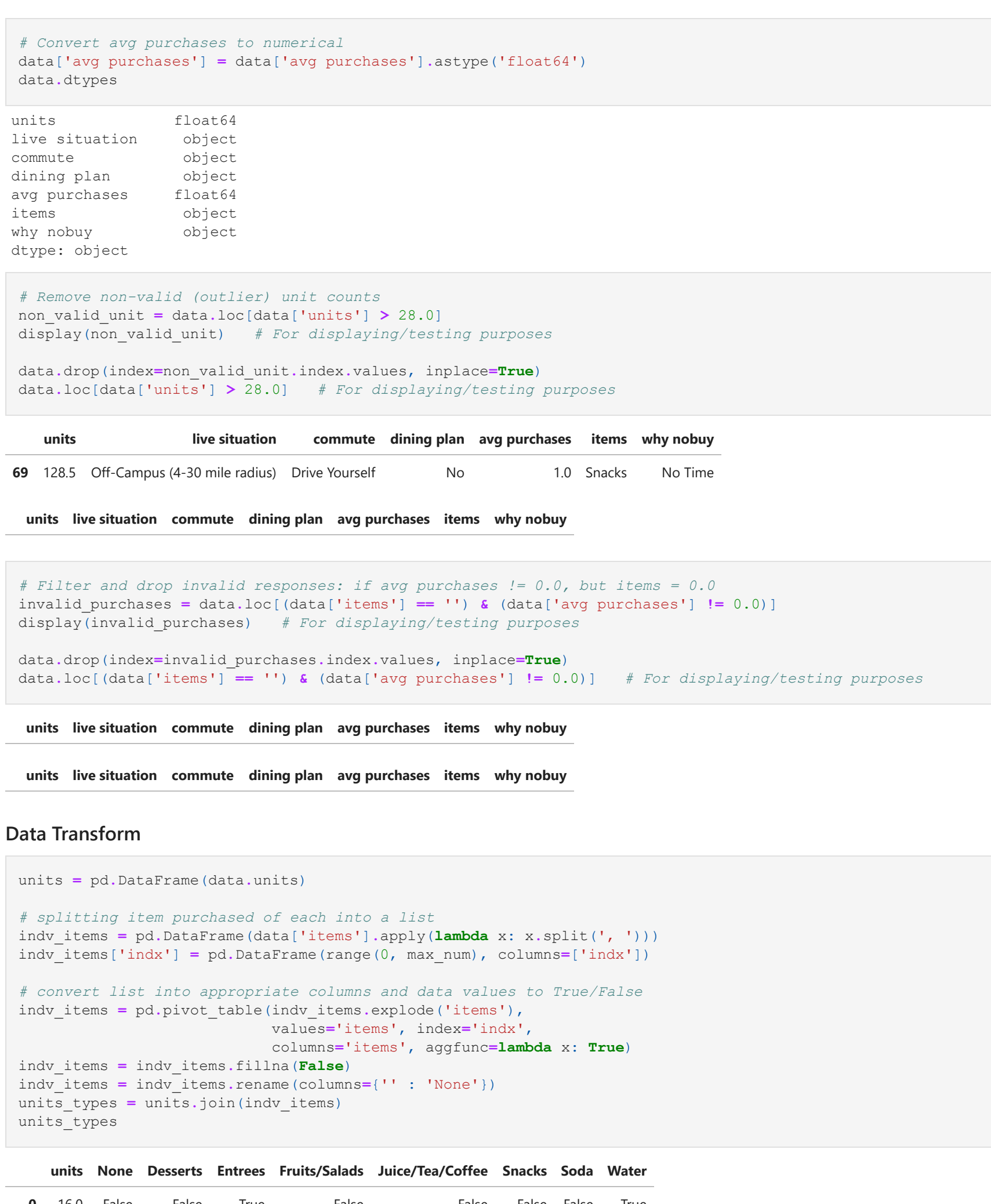
Each column consists of one of these questions, while each row consists of an individual's response to each question. Essentially, the data we have relates to the student's amount of units they are taking, their living situation and if they commute, whether or not they purchase food on campus, and if they don't what is the reason why.

## 2. What would you like to know?

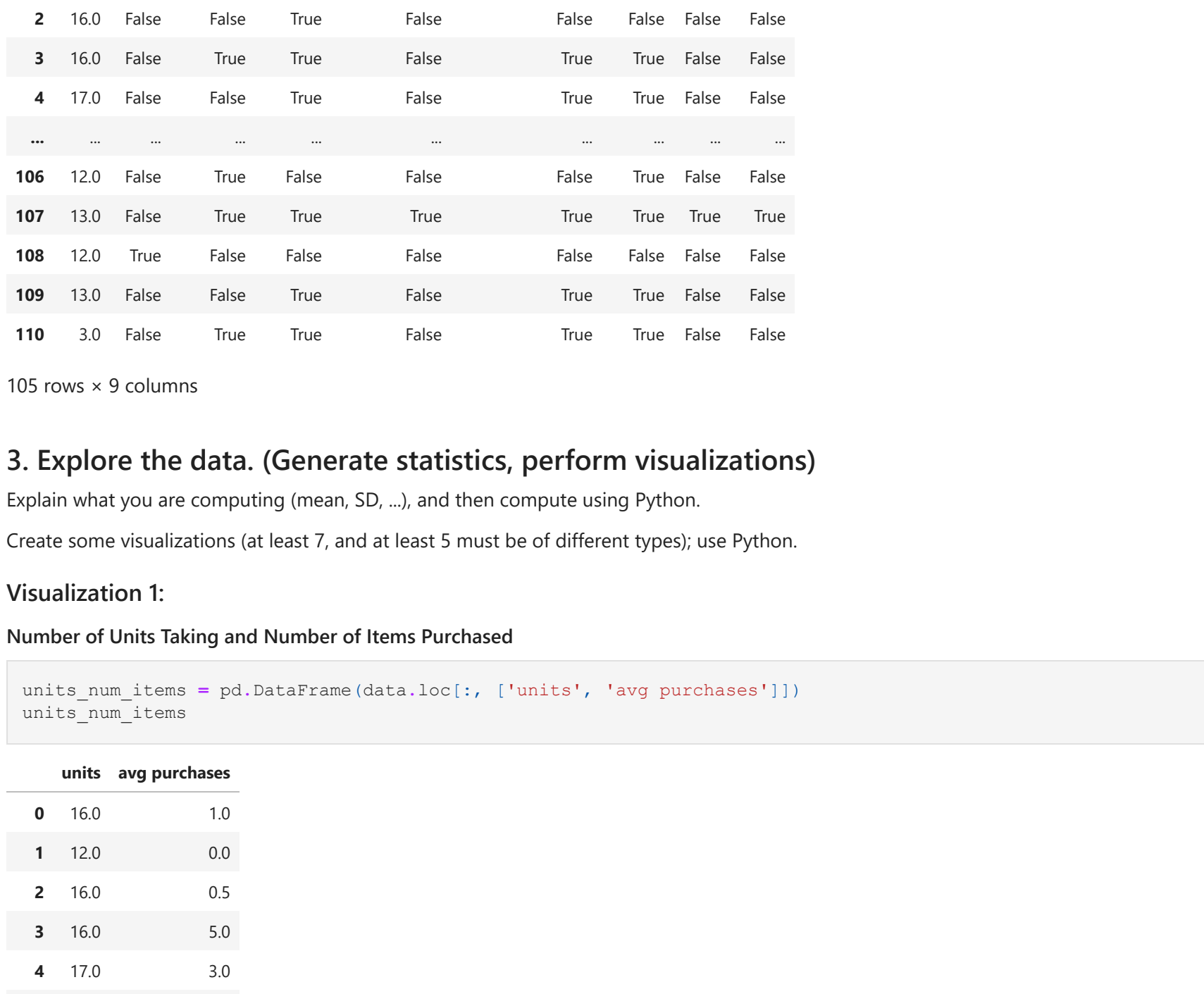
We would like to discover if the amount of enrolled units and living situations of students affects how they obtain food on campus. We are attempting to determine if a student's total enrolled units affect if they buy food more often. Particularly, we are curious if a student who has at least 12 enrolled units buys food on campus more than a student enrolled in less than 12 units. We are also trying to determine if their living situation affects whether or not students purchase food on campus.



## Cleaning



## Data Transform



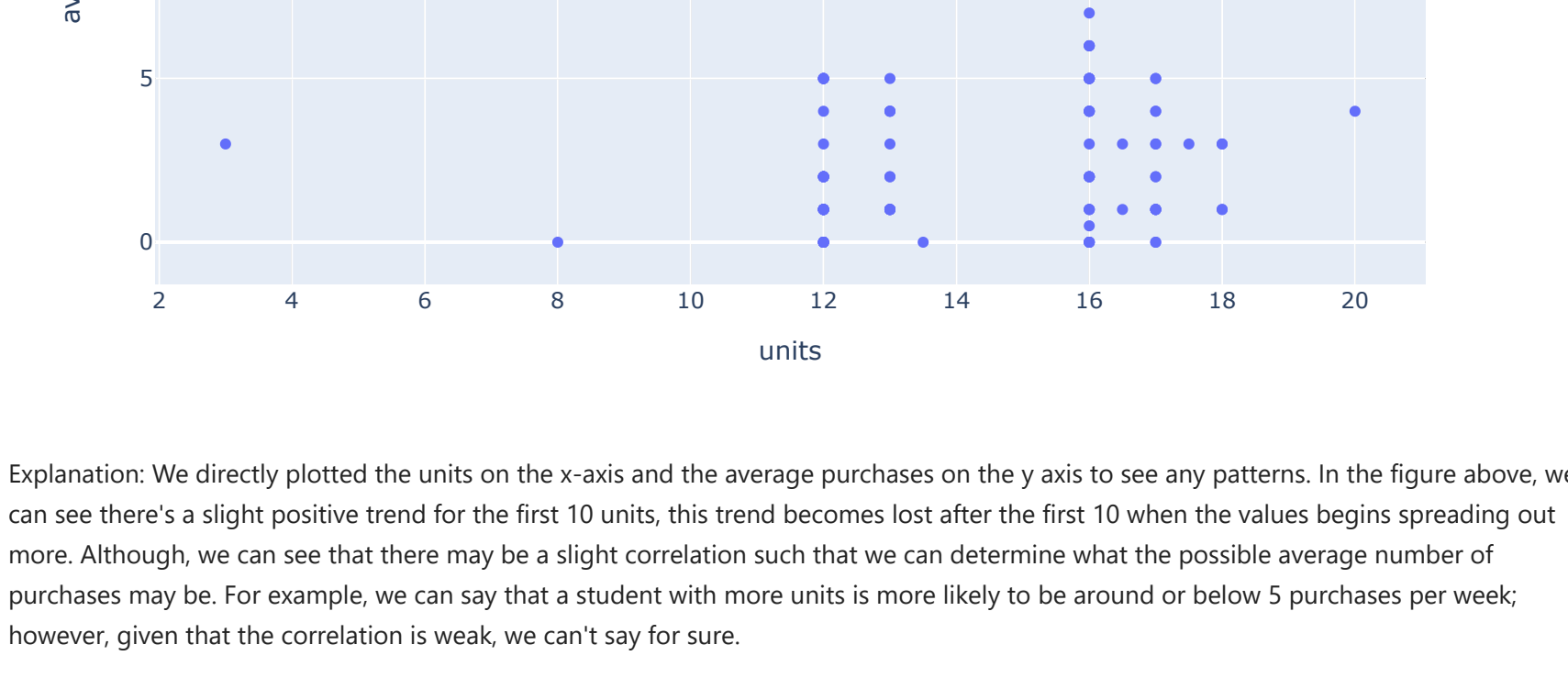
## 3. Explore the data. (Generate statistics, perform visualizations)

Explain what you are computing (mean, SD...), and then compute using Python.

Create some visualizations (at least 7, and at least 5 must be of different types); use Python.

### Visualization 1:

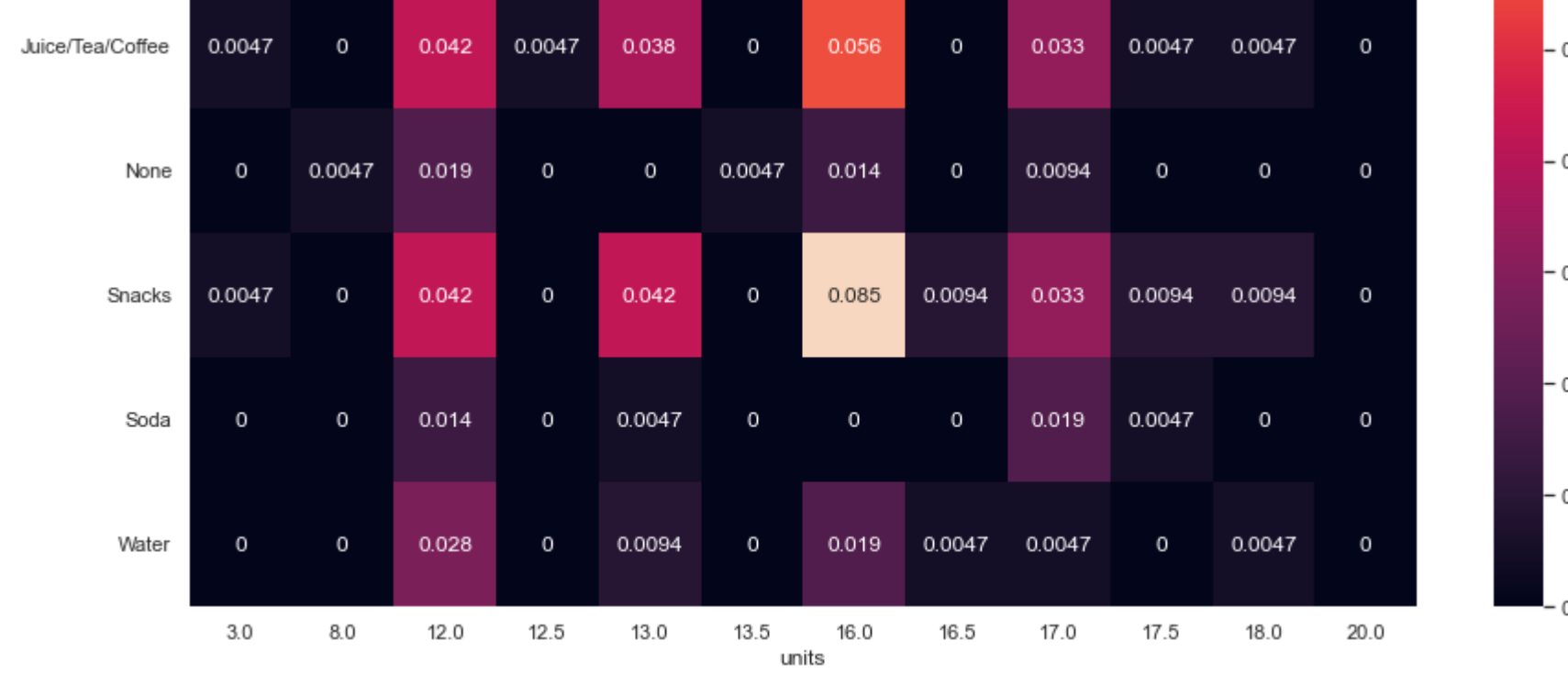
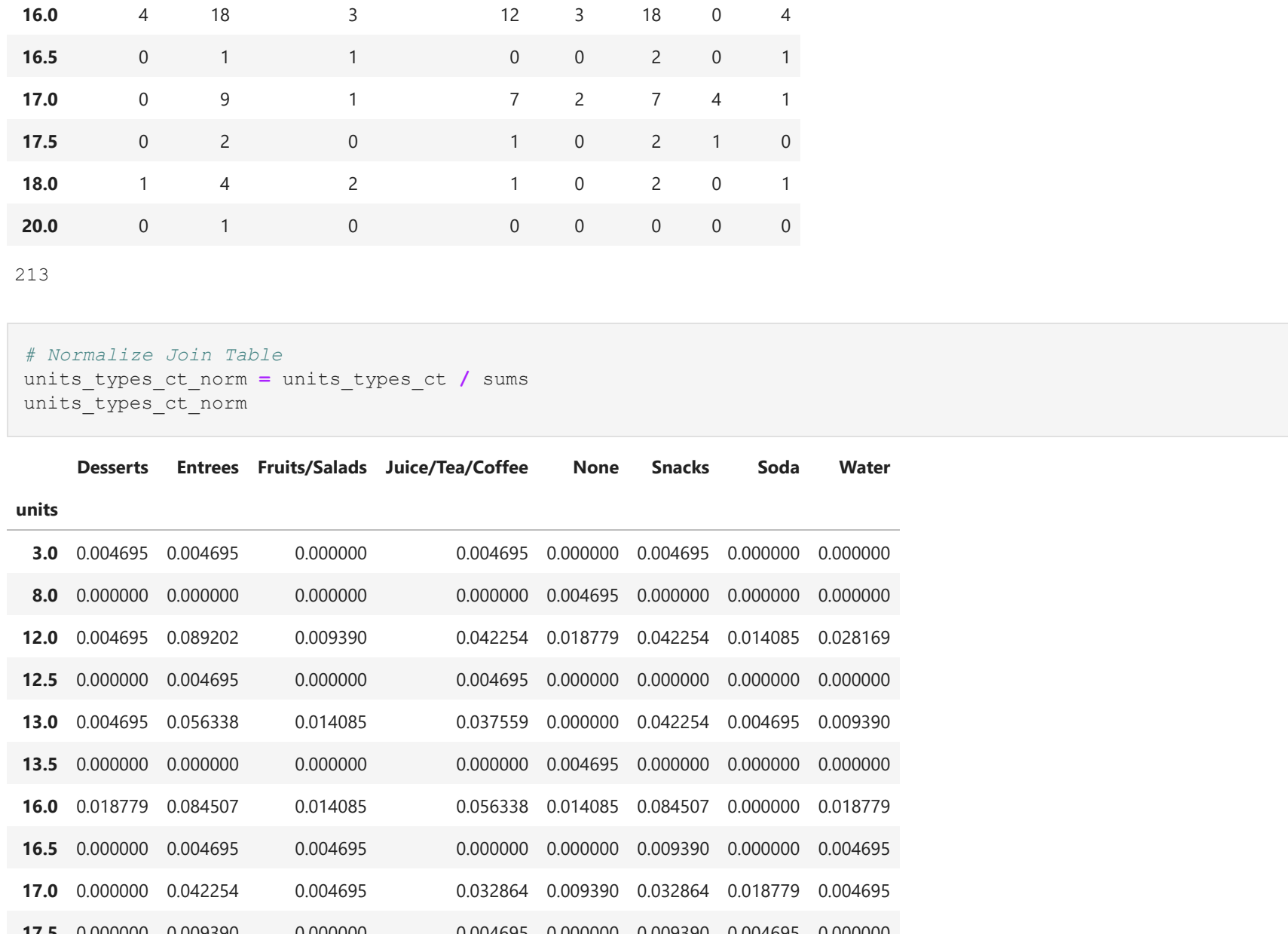
Number of Units Taking and Number of Items Purchased



Explanation: We directly plotted the units on the x-axis and the average purchases on the y-axis to see any patterns. In the figure above, we can see there's a slight positive trend for the first 10 units, this trend becomes lost after the 10th when the values begin spreading out more. Although, we can see that there may be a slight correlation such that we can determine what the possible average number of purchases may be. For example, we can say that a student with more units is more likely to be around or below 5 purchases per week; however, given that the correlation is weak, we can't say for sure.

### Visualization 2:

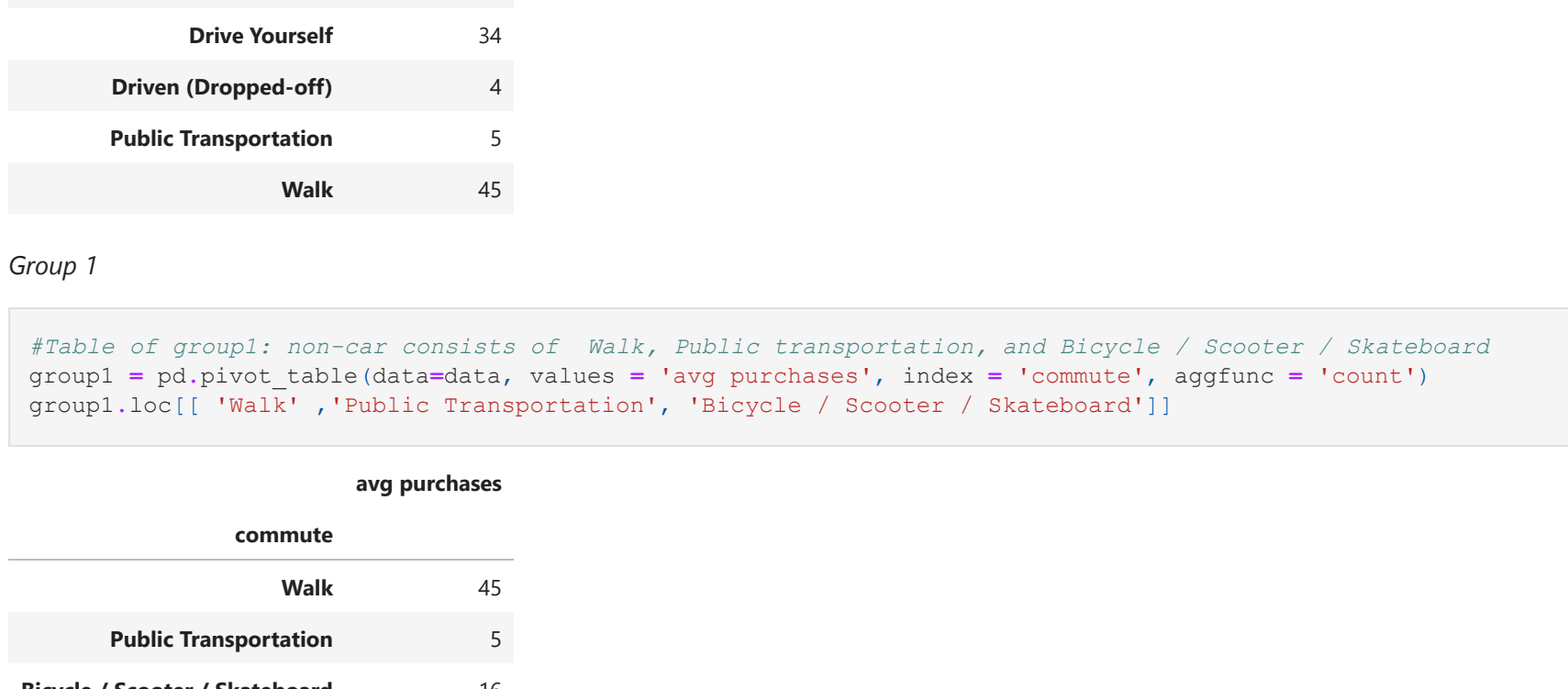
Number of Units Taking and Types of Items Purchased



Explanation: This visualization was done as we are comparing 2 dimensions, but are interested in the values the two variables give. Since the values give a count, we used a heatmap to display when our value is significant versus when it is not. Units is labeled on the x-axis, with item types labeled on the y-axis. For this visualization, the values were normalized as it gave a slightly better view of how much the units affect their choices. We can determine from this visualization that units don't really have an effect on what they tend to buy as it is evenly distributed amongst 3 main types across the board. However, we can determine that most students tend to purchase entrees, and students with higher course loads given the number of units will have a higher chance of purchasing an item than others; although, they are all relatively likely to.

### Visualization 3

How a student's method of commuting affect how many items they purchase





```
In [30]: # Plot of students who do have a DCR dining plan,
# and how many times a week they purchase food/drink on campus
bins = np.arange(26) + 0.5
sns.displot(data=group56_57_yes, bins=bins, kde=True, height=8, legend=False)
plt.xticks(range(33))
plt.xlim((-1,20))
plt.title('Group A: Students who do have dining plan')
plt.xlabel('Food/drinks bought')
plt.ylabel('Students')
plt.show()
```



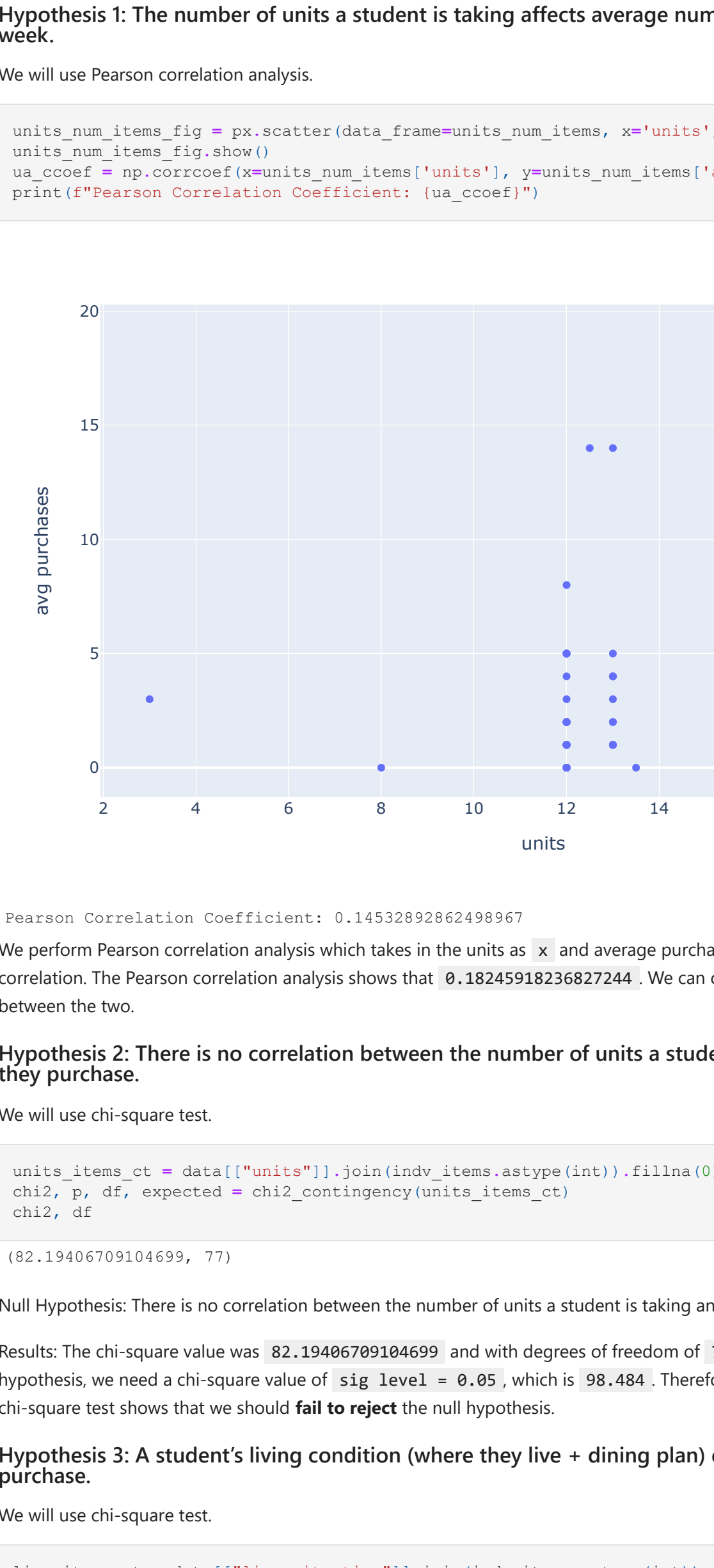
Explanation:  
For Group A: Students who do have dining plan, there are a total of 30 students with a total of 153 food/drinks. By taking the mean of the total amount of food/drinks bought and # of students who do have dining plan we got the average of 5.1 items per week.

```
In [31]: #group_b: students who do not have a dining plan
group56_57 = data[['dining plan', 'avg purchases']]
group56_57_no = group56_57.loc[group56_57['dining plan'] == "No"]
group56_57_no = group56_57_no.drop('dining plan', axis=1)
group56_57_no = group56_57_no.reset_index(drop=True)
```

```
In [32]: #The mean for group_b: students who do not have a dining plan
students = group56_57_no.value_counts().sum()
total_items = group56_57_no['avg purchases'].sum()

total_items / students
#average amount of items bought for group_b
2.2333333333333334
```

```
In [33]: # Plot of students who do NOT have a DCR dining plan,
# and how many times a week they purchase food/drink on campus
bins = np.arange(57) + 0.5
sns.displot(data=group56_57_no, bins=bins, kde=True, height=8, legend=False)
plt.xticks(range(12))
plt.yticks(range(27))
plt.xlim((-1,12))
plt.title('Group B: Students who do not have dining plan')
plt.xlabel('Food/drinks bought')
plt.ylabel('Students')
plt.show()
```



Explanation:  
For Group B: Students who do not have dining plan, there are a total of 75 students with a total of 168 food/drinks. By taking the mean of the total amount of food/drinks bought and # of students who do not have dining plan we got the average of 2.2333333333333334 items.

Summary:  
Here we defined two groups which consists of "Group A: Students who have dining plan" and "Group B: Students who do not have dining plan". We can see that Group A has a bigger average of items than Group B which is surprising due to the fact that we believed that it would be the other way around. We can conclude that those who have a dining plan end up buying more food/drinks than those who do not have a dining plan.

```
In [34]: # Visualization 7
```

#### 4. Can you state any hypotheses or make some predictions? Which tests can you apply to verify your hypothesis?

State clearly each of your hypotheses (at least 3).

**Hypothesis 1: The number of units a student is taking affects average number of items they will purchase per week.**

- We can use Pearson correlation analysis to verify hypothesis 1.

**Hypothesis 2: There is no correlation between the number of units a student is taking and the types of items they purchase.**

- We can use chi-square test to verify hypothesis 2.

**Hypothesis 3: A student's living condition (where they live + dining plan) does not affect the type of items they purchase.**

- We can use chi-square test to verify hypothesis 3.

#### 5. Test your hypotheses.

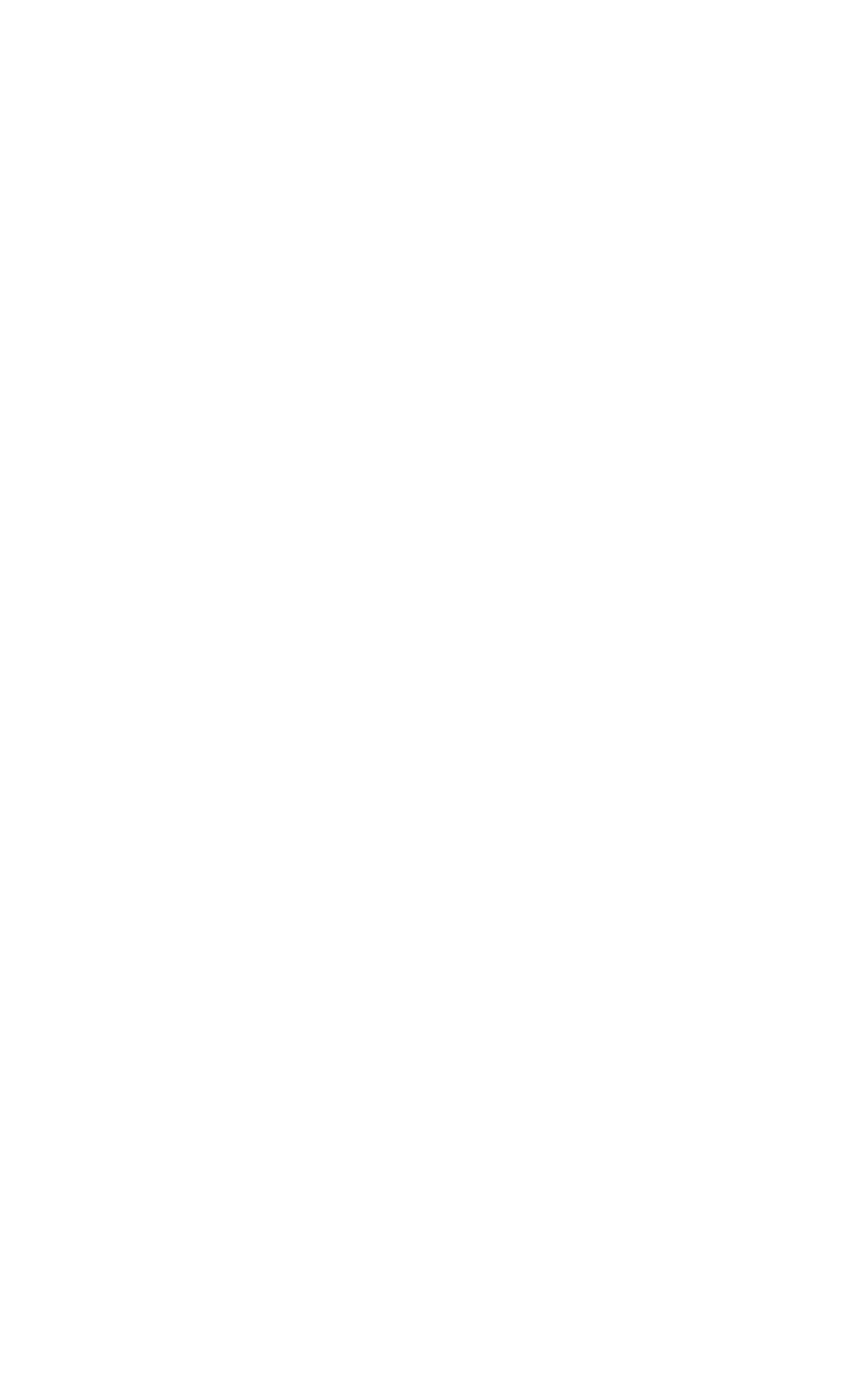
Test your hypotheses and predictions (use at least 2 different tests). For each: describe the test you are using; perform it; analyze the results and draw the conclusion.

**Hypothesis 1: The number of units a student is taking affects average number of items they will purchase per week.**

We will use Pearson correlation analysis.

```
In [35]: units_num_items_fig = px.scatter(data_frame=units_num_items, x='units', y='avg purchases')
units_num_items_fig.show()
```

```
ua_coef = np.corrcoef(x=units_num_items['units'], y=units_num_items['avg purchases'])[1,0]
print(f"Pearson Correlation Coefficient: {ua_coef}")
```



Pearson Correlation Coefficient: 0.14532892862408967

We perform Pearson correlation analysis which takes in the units as  $x$  and average purchases as  $y$  and in order to test whether there is a correlation. The Pearson correlation analysis shows that 0.18245918236827244. We can conclude that there is almost no correlation between the two.

**Hypothesis 2: There is no correlation between the number of units a student is taking and the types of items they purchase.**

We will use chi-square test.

```
In [36]: units_items_ct = data[['units']].join(indv_items.astype(int)).fillna(0).groupby("units").sum()
units_num_items_fig.show()
```

```
chi2, p, df, expected = chi2_contingency(units_items_ct)
chi2, df
```

```
Out[36]: (82.19406709104699, 77)
```

Null Hypothesis: There is no correlation between the number of units a student is taking and the types of items they purchase.

Results: The chi-square value was 82.19406709104699 and with degrees of freedom of 77 the table shows that in order to reject null hypothesis, we need a chi-square value of sig level = 0.05, which is 98.484. Therefore, our chi-square is less than 98.484 so our chi-square test shows that we should fail to reject the null hypothesis.

**Hypothesis 3: A student's living condition (where they live + dining plan) does not affect the type of items they purchase.**

We will use chi-square test.

```
In [37]: live_items_ct = data[['live situation']].join(indv_items.astype(int)).fillna(0).groupby("live situation").sum()
live_num_items_fig.show()
```

```
chi2, p, df, expected = chi2_contingency(live_items_ct)
chi2, df
```

```
Out[37]: (29.451571862592957, 21)
```

Null Hypothesis: A student's living condition (where they live + dining plan) does not affect the type of items they purchase.

The chi-square value was 29.451571862592957 and with degrees of freedom of 21 the table shows that in order to reject null hypothesis, we need a chi-square value of sig level = 0.05, which is 32.671. Therefore, our chi-square is less than 32.671 so our chi-square test shows that we should fail to reject the null hypothesis.

Pearson Correlation Coef:

$x = \text{units}$      $y = \text{avg purchases}$

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

$$\bar{x} = 14.667$$

$$\bar{y} = .5$$

$$16 - 14.667 = 1.333$$

$$\frac{((16 - 14.667)(1 - .5)) + ((12 - 14.667)(0 - .5)) + ((16 - 14.667)(.5 - .5))}{\sqrt{((16 - 14.667)^2 + (12 - 14.667)^2 + (16 - 14.667)^2) \cdot ((1 - .5)^2 + (0 - .5)^2 + (.5 - .5)^2)}}$$

$$\frac{((1.333)(.5)) + ((-2.667)(-.5)) + ((1.333)(0))}{\sqrt{((1.333)^2 + (-2.667)^2 + (1.333)^2) \cdot ((.5)^2 + (-.5)^2 + (0)^2)}}$$

$$\frac{.6665 + 1.3335}{\sqrt{(10.667)(.5)}}$$

$$\frac{2}{2.309} = .866025$$

```
units_num_items.head(3)
```

	units	avg purchases
0	16.0	1.0
1	12.0	0.0
2	16.0	0.5

```
In [30]: ua_ccoef = np.corrcoef(x=units_num_items['units'].head(3), y=units_num_items['avg purchases'].head(3))[1,0]  
print(f"Pearson Correlation Coefficient: {ua_ccoef}")
```

Pearson Correlation Coefficient: 0.8660254037844385

	None	Desserts	Entrees	Fruits/Salads	Juice/Tea/Coffee	Snacks	Soda	Water	Total
units									
3.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0
8.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1
12.0	2.0	0.0	19.0	2.0	9.0	8.0	3.0	6.0	49
Total:	3	0	19	2	9	8	3	6	

	none	Dessert	Entrees	Fruit	Juice	Snacks	Soda	Water
3.0	<del>0-0</del> 0	0	0	0	0	0	0	0
8.0	$\frac{(1 - 6.125)^2}{6.125}$	$\frac{(0 - 6.125)^2}{6.125}$	$\frac{(19 - 6.125)^2}{6.125}$	$\frac{(2 - 6.125)^2}{6.125}$	$\frac{(9 - 6.125)^2}{6.125}$	$\frac{(8 - 6.125)^2}{6.125}$	$\frac{(3 - 6.125)^2}{6.125}$	$\frac{(6 - 6.125)^2}{6.125}$
12.0	$\frac{(2 - 6.125)^2}{6.125}$	$\frac{(0 - 6.125)^2}{6.125}$	$\frac{(19 - 6.125)^2}{6.125}$	$\frac{(2 - 6.125)^2}{6.125}$	$\frac{(9 - 6.125)^2}{6.125}$	$\frac{(8 - 6.125)^2}{6.125}$	$\frac{(3 - 6.125)^2}{6.125}$	$\frac{(6 - 6.125)^2}{6.125}$

$$6.125 + 8(.125) + 2.778 + .625 + 27.063 + 6.125 + 1.349 + 1.573 + 1.544 + .002$$

$$\chi^2 = 47.24$$

	None	Desserts	Entrees	Fruits/Salads	Juice/Tea/Coffee	Snacks	Soda	Water	Total
live situation									
Off-Campus (3-mile radius)	4	3	23	6	11	13	3	8	71
Off-Campus (31+ -mile radius)	2	0	3	0	0	0	0	0	5
Total:	6	3	26	6	11	13	3	8	

	None	Dessert	Entrees	Fruit	Juice	Snacks	Soda	Water
3 mile	$\frac{(4-8.875)^2}{8.875}$	$\frac{(3-8.875)^2}{8.875}$	$\frac{(23-8.875)^2}{8.875}$	$\frac{(6-8.875)^2}{8.875}$	$\frac{(11-8.875)^2}{8.875}$	$\frac{(13-8.875)^2}{8.875}$	$\frac{(3-8.875)^2}{8.875}$	$\frac{(8-8.875)^2}{8.875}$
31+ mile	$\frac{(2-.625)^2}{.625}$	$\frac{(0-.625)^2}{.625}$	$\frac{(3-.625)^2}{.625}$	$\frac{(0-.625)^2}{.625}$	$\frac{(0-.625)^2}{.625}$	$\frac{(0-.625)^2}{.625}$	$\frac{(0-.625)^2}{.625}$	$\frac{(0-.625)^2}{.625}$

$$6(.625) + .3025 + 2.677 + 3.889 + 22.48 + .931 + .508 + 1.917 + 3.889 + .086$$

$$\chi^2 = 40.4297$$