
0.1 Question 2e

If we were to drop businesses with MISSING postal code values, what specific types of businesses would we be excluding? In other words, is there a commonality among businesses with missing postal codes?

Hint: You may want to look at the names of the businesses with missing postal codes. Feel free to reuse parts of your code from 2d, but we will not be grading your code.

SOLUTION:

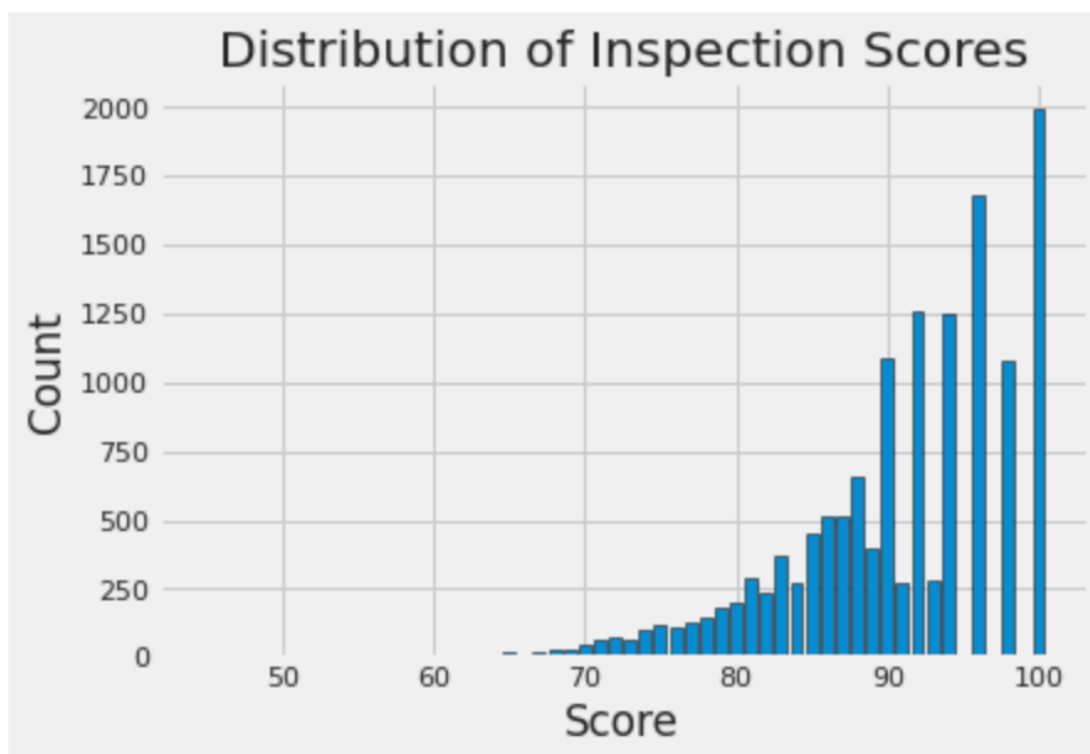
The businesses with the -9999 ZIP code appear to be primarily food trucks and concession establishments.

0.2 Question 5a

Let's look at the distribution of inspection scores. As we saw before when we called `head` on this data frame, inspection scores appear to be integer values. The discreteness of this variable means that we can use a bar plot to visualize the distribution of the inspection score. Make a bar plot of the counts of the number of inspections receiving each score.

It should look like the image below. It does not need to look exactly the same (e.g., no grid), but **make sure that all labels and axes are correct**.

You should use the `ins` dataframe, and should ignore any score that is less than 0.

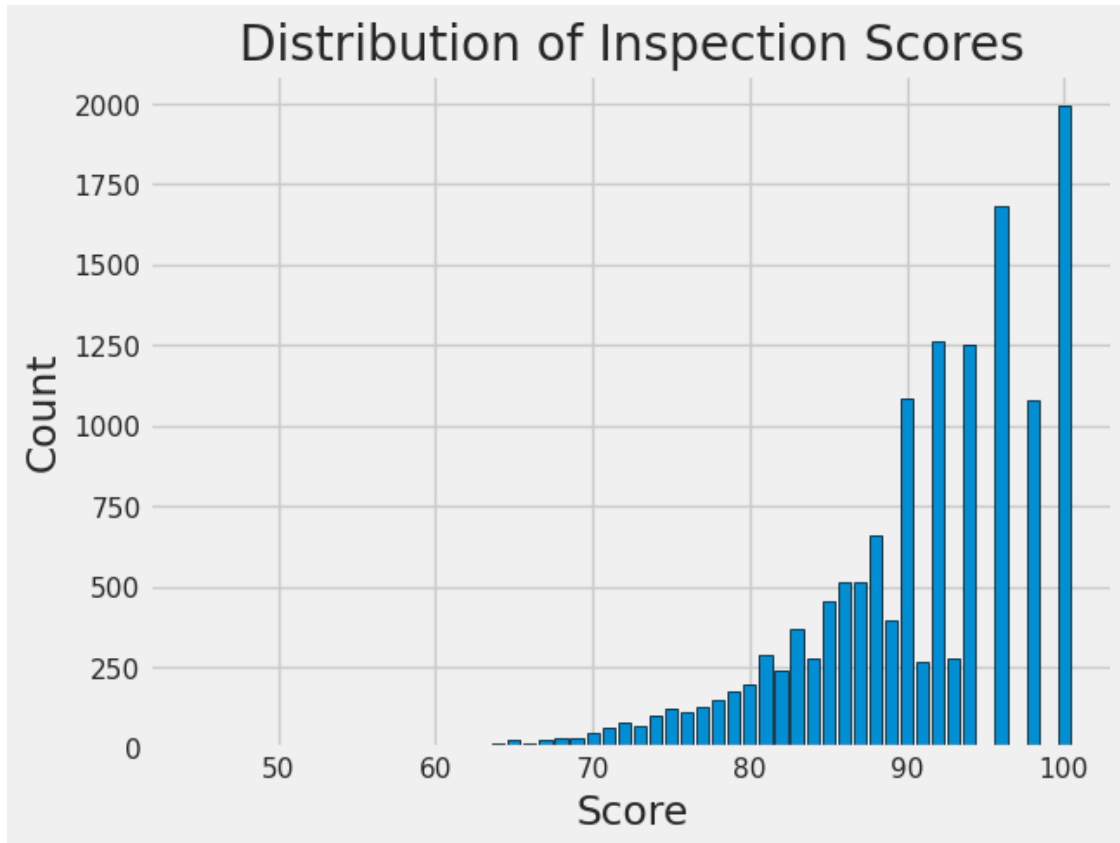


You might find this [matplotlib.pyplot tutorial](#) useful. Key syntax that you'll need:

```
plt.bar
plt.xlabel
plt.ylabel
plt.title
```

To set the color of the edges for your bars, include 'edgecolor = 'black'.

```
In [65]: score_counts = ins.loc[ins["score"] > 0, 'score'].value_counts()
plt.bar(score_counts.index, score_counts, edgecolor = 'black')
plt.xlabel("Score")
plt.ylabel("Count")
plt.title("Distribution of Inspection Scores");
```



0.2.1 Question 5b

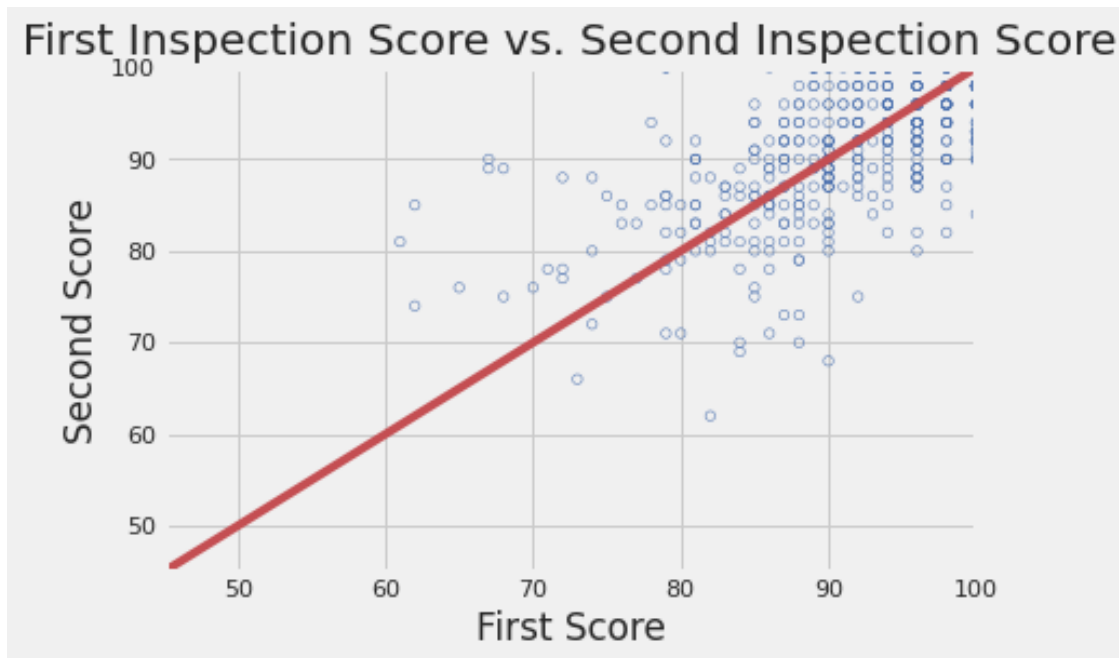
Now let's actually reflect on the histogram that we generated before with a bin size of 1.

Describe the qualities of the distribution of the inspections scores based on your bar plot. Consider the mode(s), symmetry, tails, gaps, and anomalous values. Are there any unusual features of this distribution? What do your observations imply about the scores?

The distribution is unimodal with a peak at 100. It is skewed left (as expected with a variable bounded on the right). The distribution has a long left tail with some restaurants receiving scores that are in the 50s, 60s, and 70s. One unusual feature of the distribution is the bumpiness with even numbers having higher counts than odd. This may be because the violations result in penalties of 2, 4, 10, etc. points.

Now let's make a scatter plot to display these pairs of scores. Include on the plot a reference line with slope 1 and y-intercept 0. Since restaurant scores bottom out at 45 points, we'll only focus on ratings between 45 and 100. Thus your reference line should start at [45, 45] and go up to [100, 100].

Create your scatter plot in the cell below. It does not need to look exactly the same (e.g., no grid) as the sample below, but make sure that all labels, axes and data itself are correct.



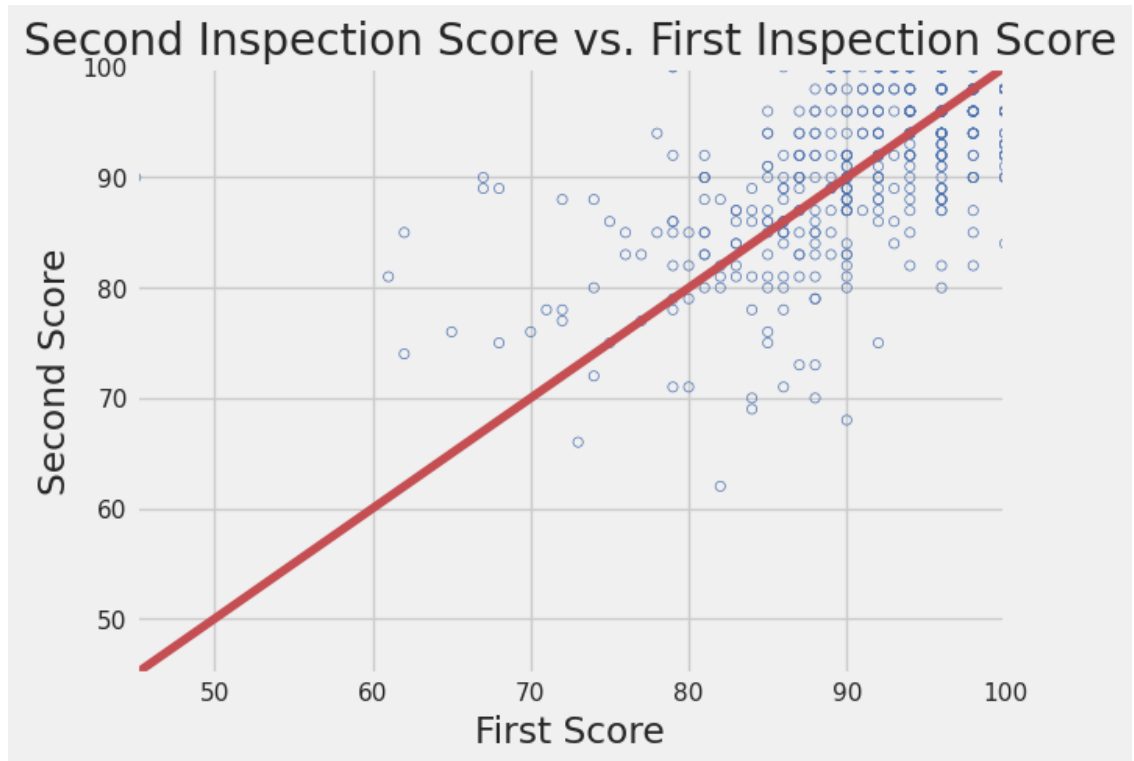
Key pieces of syntax you'll need:

`plt.scatter` plots a set of points. Use `facecolors='none'` and `edgecolors='b'` to make circle markers with blue borders.

`plt.plot` for the reference line. Using the argument `r` will make the line red.

`plt.xlabel`, `plt.ylabel`, `plt.axis`, and `plt.title`.

```
In [69]: # BEGIN SOLUTION
x = scores["first score"]
y = scores["second score"]
plt.scatter(x,y,s = 20, facecolors='none',edgecolors='b')
plt.plot([45,100],[45,100], 'r')
plt.xlabel('First Score')
plt.ylabel('Second Score')
plt.axis([45,100,45,100])
plt.title("Second Inspection Score vs. First Inspection Score");
# END SOLUTION
```



0.2.2 Question 6c

If restaurants' scores tend to improve from the first to the second inspection, what do you expect to see in the scatter plot that you made in question 6b? What do you observe from the plot? Are your observations consistent with your expectations?

Hint: What does the slope represent?

bhbhbh

