

Discussion 4 Solutions

Note: Starting this week, I'm going to release solutions as an answer document instead of the filled worksheet to have space to explain everything, but just FYI: all of the following is from practice.dsc80.com — the purpose of this is just so you don't have to cross-reference anything yourself!

FA23 Midterm Problem 3

Problem

- A. The researchers chose the 30 donkeys with the largest `'weight'` values to reweigh.
- B. The researchers drew 30 donkeys uniformly at random without replacement from the donkeys with `BCS` score of 4 or greater.
- C. The researchers set `i` as a number drawn uniformly at random between 0 and 514, then reweighed the donkeys in `donkeys.iloc[i:i+30]`.
- D. The researchers reweighed all the donkeys, but deleted all the values in `'weightAlt'` except for the 30 lowest values.
- E. The researchers split up the donkeys into the 6 different age groups, then sampled 5 donkeys uniformly at random without replacement within each age group.

Solution

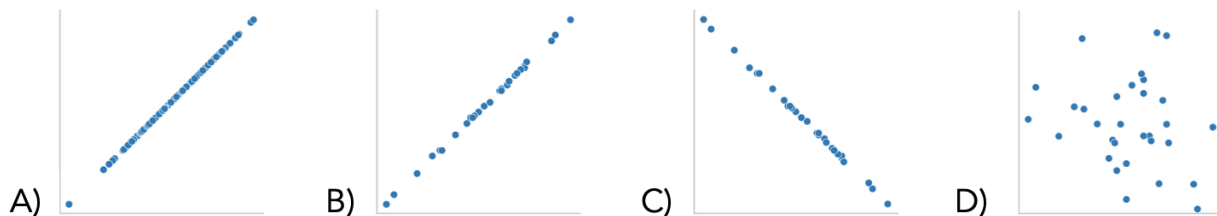
- A. **Missing at random.** This means missing values depend on another column in the DataFrame. In this case, the missing values of `'weightAlt'` depend on the `'weight'` column since we select the 30 largest.
- B. **Missing at random.** This means missing values depend on another column in the DataFrame. In this case, the missing values of `'weightAlt'` depend on the `'BCS'` column since we choose from those with a score of 4 or greater.
- C. **Missing completely at random or, possibly, Missing at Random.** The argument for MAR is as follows: this means missing values depend on another column in the DataFrame. The missing values depend on the index since index 0 can only be selected if `i = 0`, but index 29 could be chosen if `i` is any value between 0 and 29, so it has a higher probability of being chosen. The original solution was MCAR as we did not account for edge case of `i` being small, but it is technically MAR. Credit was given for either answer.

- **D. Not missing at random.** This means missing values depend on the column they're missing from. The missing values here are all values that are not the 30 lowest in 'Weight', and so they depend on the column itself.
- **E. Missing completely at random or Missing at random.** If the data was assumed to be evenly distributed, then the data is missing completely at random since the six age groups would all be chosen from uniformly. However, if the data was assumed to possibly have skewed age data, then samples from small sample size age groups had a higher probability of being chosen than those of large sample size age group. Credit was given for either answer.

NOTE: Despite the fact that we accepted multiple answers for a couple of these, you should make *as few assumptions about the data as possible* to get your solutions — but if you're unsure, feel free to ask!

Problem

For this next question, assume that the researchers chose the 30 donkeys to reweigh by drawing a simple random sample of 30 underweight donkeys: donkeys with BCS values of 1, 1.5, or 2. The researchers weighed these 30 donkeys one day later and stored the results in 'WeightAlt'. Which of the following shows the scatter plot of 'WeightAlt' - 'Weight' on the y-axis and 'Weight' on the x-axis? Assume that missing values are not plotted.

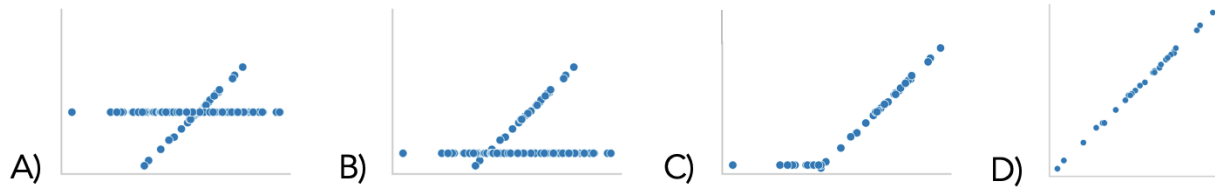


Solution

We are measuring the difference in weight from just one day on the y-axis, which means we can't expect any noticeable pattern of weight gain or loss no matter the original weight of the donkey. Therefore, a random scatterplot makes sense. Options A through C all suggest that the single-day weight change correlates with the starting weight, which is not a good assumption.

Problem

Suppose we use mean imputation to fill in the missing values in `'weightAlt'`. Select the scatter plot `'weightAlt'` on `'weight'` after imputation.



Solution

Note we are now plotting `'weight'` on the y-axis, not the difference of `'weightAlt' - 'weight'`. Therefore, it makes sense that we would have 30 data points with a positive slope as the initial weight and re-weight are likely very similar.

Then, mean imputation is the process of filling in missing values with the average of the non-missing values. Therefore, all missing values will be the same, and should be at the center of the sloped line since the line is roughly evenly distributed.

FA23 Final Problem 3

The `bus` table (left) records bus arrivals over 1 day for all the bus stops within a 2 mile radius of UCSD. The data dictionary (right) describes each column.

	time	line	stop	late
0	12pm	201	Gilman Dr & Mandeville Ln	-1.1
1	1:15pm	30	Gilman Dr & Mandeville Ln	2.8
2	11:02am	101	Gilman Dr & Myers Dr	-0.8
3	8:04am	202	Gilman Dr & Myers Dr	NaN
4	9am	30	Gilman Dr & Myers Dr	-3.0

`time` Time of arrival (`str`). Note that the times are inconsistently entered (e.g. 12pm vs. 1:15pm).

`line` Bus line (`int`). There are multiple buses per bus line each day.

`stop` Bus stop (`str`).

`late` The number of minutes the bus arrived after its scheduled time. Negative numbers mean that the bus arrived early (`float`). Some entries in this column are missing.

For each of the following questions, select the correct procedure to simulate a single sample under the null hypothesis, and the correct test statistic for the hypothesis test. Assume that the `time` column of the bus DataFrame has already been parsed into timestamps.

Problem

Are buses equally likely to be early or late? *Note: while the problem says there is only one solution, post-exam two options for the test statistic were given credit. Pick one of the two.*

Simulation procedure:

- A. `np.random.choice([-1, 1], bus.shape[0])`
- B. `np.random.choice(bus['late'], bus.shape[0], replace = True)`
- C. Randomly permute the `'late'` column

Test statistic:

- A. Number of values below 00
- B. `np.mean`
- C. `np.std`
- D. TVD
- E. K-S statistic

Solution

Simulation procedure: The sample we have here is something like 152 early buses, 125 late buses (these numbers are made up – in practice, these two numbers need to add to `bus.shape[0]`). The question is whether this sample looks like it was drawn from a population that is 50-50 (an equal number of early and late buses), which makes this a hypothesis test. In terms of examples from class, this most closely resembles the very first hypothesis testing example we looked at – the “coin flipping” example.

`np.random.choice([-1, 1], bus.shape[0])` will return an array of length `bus.shape[0]`, where each element is equally likely to be either `-1` (late) or `1` (early). (Note that we could also take `-1` to mean early and `1` to mean late – it doesn't really matter.)

Test statistic: Each time we simulate an arrays of `-1` s and `1` s, we'd like to compute a statistic

that helps us differentiate between the number of late (`-1`) and the number of early (`1`) simulated buses. The number of values below 0 will give us the number of late simulated buses, so we could use that. ~~The mean of the `-1` s and `1` s will give us a value that is negative if there were more late buses and positive if there were more early buses, so we could use that too.~~

NOTE: this problem accepted `np.mean` for the test statistic, but I am pretty confident that this won't work with some pretty simple assumptions about the data, and I'll see about getting that fixed

Problem

Is the `'late'` column MAR dependent on the `'line'` column?

Simulation procedure:

- `np.random.choice([-1, 1], bus.shape[0])`
- `np.random.choice(bus['late'], bus.shape[0], replace = True)`
- Randomly permute the `'late'` column

Test statistic:

- Absolute difference in means
- Absolute difference in proportions
- TVD
- K-S statistic

Solution

Answer: Simulation procedure: Randomly permute the `'late'` column; Test statistic: TVD

Simulation procedure: To determine if `'late'` is missing at random dependent on the `'line'` column, we conduct a permutation test and compare (1) the distribution of the `'line'` column when the `'late'` column is missing to (2) the distribution of the `'line'` column when the `'late'` column is not missing to see whether they're significantly different. If the distributions are indeed significantly different, then it is likely that the `'late'` column is MAR dependent on `'line'`.

Test statistic: Since we are comparing the distributions of *categorical* data (`'line'` is categorical) for our permutation test, Total Variation Distance is the best test statistic to use.