

DSC 80 Discussion 2 Worksheet

1 WI24 Midterm Problem 3

Jasmine is a veterinarian. Below, you'll find information about some of the dogs in her care, separated by district and breed.

| | Beagle | | Cocker Spaniel | |
|------------|-------------|-------|----------------|-------|
| | Mean Weight | Count | Mean Weight | Count |
| District 1 | 25 | 3 | 20 | 2 |
| District 2 | 45 | 1 | x | y |

What is the mean weight of all beagles in the table above, across both districts?

$$((25 * 3) + (45 * 1)) / 4 = 30$$

Notice that the table above has two unknowns, x and y . Find **positive integers** x and y such that the mean weight of all beagles is equal to the mean weight of all cocker spaniels, *where x is as small as possible*.

$x = 31, y = 20$ x must be > 30 to raise Cocker Spaniel mean to 30, and if we try $x = 31$, we can solve for y using same weighted mean equation as above

*for the additional problem about Simpson's paradox happening both ways, I'll make a post on Ed

2 FA23 Midterm Problem 1

| | date | name | food | weight |
|---|------------|--------|-------------|--------|
| 0 | 2023-01-01 | Sam | Ribeye | 0.20 |
| 1 | 2023-01-01 | Sam | Pinto beans | 0.10 |
| 2 | 2023-01-01 | Lauren | Mung beans | 0.25 |
| 3 | 2023-01-02 | Lauren | Lima beans | 0.30 |
| 4 | 2023-01-02 | Sam | Sirloin | 0.30 |

Find the total kg of food eaten for each day and each person in `df` as a Series.

```
df.groupby(['date', 'name'])[['weight']].sum()
```

Find all the unique people who did not eat any food containing the word "beans".

```
def foo(x):  
    return x['food'].str.contains('beans').sum() == 0
```

```
df.groupby('name').filter(lambda x: foo(x))['name'].unique()
```

3 FA23 Final Problem 1

The `bus` table (left) records bus arrivals over 11 day for all the bus stops within a 22 mile radius of UCSD.

| | time | line | stop | late | |
|---|---------|------|---------------------------|------|---|
| 0 | 12pm | 201 | Gilman Dr & Mandeville Ln | -1.1 | time Time of arrival (str). Note that the times are inconsistently entered (e.g. 12pm vs. 1:15pm). |
| 1 | 1:15pm | 30 | Gilman Dr & Mandeville Ln | 2.8 | line Bus line (int). There are multiple buses per bus line each day. |
| 2 | 11:02am | 101 | Gilman Dr & Myers Dr | -0.8 | stop Bus stop (str). |
| 3 | 8:04am | 202 | Gilman Dr & Myers Dr | NaN | 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. |
| 4 | 9am | 30 | Gilman Dr & Myers Dr | -3.0 | |

The `stop` table (left) contains information for all the bus lines in San Diego (not just the ones near UCSD).

| | line | stop | next | |
|---|------|---------------------------|--------------------------------|--|
| 0 | 201 | Gilman Dr & Mandeville Ln | VA Hospital | line Bus line (int). |
| 1 | 201 | VA Hospital | La Jolla Village Dr & Lebon Dr | stop Bus stop (str). |
| 2 | 30 | VA Hospital | Villa La Jolla Dr & Holiday Ct | next The next bus stop for a particular bus line (str). For example, the first row of the table shows that after the 201 stops at Gilman Dr & Mandeville Ln, it will stop at the VA Hospital next. A missing value represents the end of a line. |
| 3 | 30 | UTC | NaN | |

Compute the number of buses in `bus` whose next stop is 'UTC'.

```
x = stop.merge(bus, on = ['line', 'stop'], how = 'inner')
x[x['next'] == 'UTC'].shape[0]
```

Compute the number of unique pairs of bus stops that are exactly two stops away from each other. For example, if you only use the first four rows of the `stop` table, then your code should evaluate to the number 2, since you can go from 'Gilman Dr & Mandeville Ln' to 'La Jolla Village Dr & Lebon Dr' and from 'Gilman Dr & Mandeville Ln' to 'Villa La Jolla Dr & Holiday Ct' in two stops. Hint: The `suffixes = (1, 2)` argument to `merge` appends a 1 to column labels in the left table and a 2 to column labels in the right table whenever the merged tables share column labels.

```
m = stop.merge(stop, left_on = 'next', right_on = 'stop',
               how = 'inner', suffixes=(1, 2))
(m[['stop1', 'next2']].drop_duplicates()).shape[0]
```

4 FA22 Midterm Problem 7

| | category | completed | minutes | urgency | client |
|------|--------------|-----------|---------|---------|------------------------------|
| 0 | work | False | NaN | 2.0 | NaN |
| 1 | work | False | NaN | 1.0 | NaN |
| 2 | work | True | 13.5 | 2.0 | NaN |
| 3 | work | False | NaN | 1.0 | NaN |
| 4 | relationship | True | 5.3 | NaN | NaN |
| ... | ... | ... | ... | ... | ... |
| 9831 | consulting | True | 71.7 | 2.0 | San Diego Financial Analysts |
| 9832 | finance | True | 36.4 | 1.0 | NaN |
| 9833 | work | True | 31.1 | 1.0 | NaN |
| 9834 | work | True | 24.8 | 3.0 | NaN |
| 9835 | work | False | NaN | 2.0 | NaN |

The code below creates a pivot table.

```
pt = tasks.pivot_table(index='urgency', columns='category', values='completed', aggfunc='sum')
```

Which of the below snippets of code will produce the same result as `pt.loc[3.0, 'consulting']`? **Select all that apply.**

Snippet 1:

```
tasks[(tasks['category'] == 'consulting') & (tasks['urgency'] == 3.0)]['completed'].sum()
```

Snippet 2:

```
tasks[tasks['urgency'] == 3].groupby('category')['completed'].sum().loc['consulting']
```

Snippet 3:

```
tasks.groupby('urgency')['completed'].sum().loc[3.0, 'consulting']
```

Snippet 4:

```
tasks.groupby(['urgency', 'category'])['completed'].sum().loc[(3.0, 'consulting')]
```

Snippet 5

```
tasks.groupby('completed').sum().loc[(3.0, 'consulting')]
```

Discussion 3 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!

WI23 Midterm Problem 6

Problem:

Given the above information, what does the following expression evaluate to?

```
tv_excl.groupby(["Age", "Service"]).sum().shape[0]
```

Solution:

Note that the DataFrame `counts` is a pivot table, created using `tv_excl.pivot_table(index="Age", columns="Service", aggfunc="size")`. As we saw in lecture, pivot tables contain the same information as the result of grouping on two columns.

The DataFrame `tv_excl.groupby(["Age", "Service"]).sum()` will have one row for every unique combination of `"Age"` and `"Service"` in `tv_excl`. (The same is true even if we used a different aggregation method, like `.mean()` or `.max()`.) As `counts` shows us, `tv_excl` contains every possible combination of a single element in `{"13+", "16+", "18+", "7+", "all"}` with a single element in `{"Disney+", "Hulu", "Netflix", "Prime Video"}`, except for `("13+", "Disney+")` and `("18+", "Disney+")`, which were not present in `tv_excl`; if they were, they would have non-null values in `counts`.

As such, `tv_excl.groupby(["Age", "Service"]).sum()` will have $20 - 2 = 18$ rows, and

```
tv_excl.groupby(["Age", "Service"]).sum().shape[0]
```

 evaluates to 18.

Problem:

Tiffany would like to compare the distribution of Age for Hulu and Netflix. Specifically, she'd like to test the following hypotheses:

- **Null Hypothesis:** The distributions of Age for Hulu and Netflix are drawn from the same population distribution, and any observed differences are due to random chance.

- **Alternative Hypothesis:** The distributions of Age for Hulu and Netflix are drawn from different population distributions.

Is this a hypothesis test, or a permutation test? Why?

Solution:

Answer: Permutation test

A permutation test is a statistical test in which we aim to determine if two samples look like they were drawn from the same unknown population. Here, our two samples are the distribution of

"Age" s for Hulu and the distribution of "Age" s for Netflix.

Problem:

Consider the DataFrame `distr`, defined below.

```
hn = counts[["Hulu", "Netflix"]]
distr = (hn / hn.sum()).T # Note that distr has 2 rows and 5 columns.
```

To test the hypotheses above, Tiffany decides to use the total variation distance as her test statistic. Which of the following expressions DO NOT correctly compute the observed statistic for her test?

Solution:

Answer:

`distr.diff().sum().sum().abs() / 2` only

First, note that the difference between the TVD calculation here and those in lecture is that our pivot table contains one **row** for each distribution, rather than one **column** for each distribution. This is because of the `.T` in the code snippet above. `distr` may look something like:

| Age | 13+ | 16+ | 18+ | 7+ | all |
|----------------|----------|----------|----------|----------|----------|
| Service | | | | | |
| Hulu | 0.004103 | 0.415385 | 0.228718 | 0.252308 | 0.099487 |
| Netflix | 0.001720 | 0.275150 | 0.382631 | 0.210662 | 0.129837 |

As such, here we need to apply the `.diff()` method to each column first, not each row (meaning we should supply `axis=0` to `diff`, not `axis=1`; `axis=0` is the default, so we don't need to explicitly specify it). `distr.diff()` may look something like:

| Age | 13+ | 16+ | 18+ | 7+ | all |
|----------------|-----------|-----------|----------|-----------|----------|
| Service | | | | | |
| Hulu | NaN | NaN | NaN | NaN | NaN |
| Netflix | -0.002383 | -0.140234 | 0.153913 | -0.041646 | 0.030349 |

With that in mind, let's look at each option, remembering that the TVD is the **sum of the absolute differences in proportions, divided by 2**.

- `distr.diff().iloc[-1].abs().sum() / 2`:
 - `distr.diff().iloc[-1]` contains the differences in proportions.
 - `distr.diff().iloc[-1].abs()` contains the absolute differences in proportions.
 - `distr.diff().iloc[-1].abs().sum() / 2` contains the sum of the absolute differences in proportions, divided by 2. **This is the TVD.**
- `distr.diff().sum().abs().sum() / 2`:
 - `distr.diff().sum()` is a Series containing just the last row in `distr.diff()`; remember, null values are ignored when using methods such as `.mean()` and `.sum()`.
 - `distr.diff().sum().abs()` contains the absolute differences in proportions, and hence `distr.diff().sum().abs().sum() / 2` contains the sum of the absolute differences in proportions, divided by 2. **This is the TVD.**

- `distr.diff().sum().sum().abs() / 2` :
 - `distr.diff().sum()` contains the differences in proportions (explained above).
 - `distr.diff().sum().sum()` contains the sum of the differences in proportions. **This is 0**; remember, the reason we use the absolute value is to prevent the positive and negative differences in proportions from cancelling each other out. As a result, this option **does not compute the TVD**; in fact, it errors, because `distr.diff().sum().sum()` is a single `float`, and `float`s don't have an `.abs()` method.
- `(distr.sum() - 2 * distr.iloc[0]).abs().sum() / 2` :
 - This option seems strange, but does actually compute the TVD. The key idea is the fact that $a - b$ is the same as $(a + b) - (2 \cdot b)$. `distr.sum()` is the same as `distr.iloc[0] + distr.iloc[1]`, so `distr.sum() - 2 * distr.iloc[0]` is `distr.iloc[0] + distr.iloc[1] - 2 * distr.iloc[0]` which is `distr.iloc[1] - distr.iloc[0]`, which is just `distr.diff().iloc[-1]`.
 - Then, this option reduces to `distr.diff().iloc[-1].abs().sum() / 2`, which is the same as Option 1. **This is the TVD.**
- `distr.diff().abs().sum(axis=1).iloc[-1] / 2` :
 - `distr.diff().abs()` is a DataFrame in which the last row contains the absolute differences in proportions.
 - `distr.diff().abs().sum(axis=1)` is a Series in which the first element is null and the second element is the sum of the absolute differences in proportions.
 - As such, `distr.diff().abs().sum(axis=1).iloc[-1] / 2` is the sum of the absolute differences in proportions divided by 2. **This is the TVD.**

WI23 Final Problem 2.5

Problem:

Suppose $\vec{a} = [a_1 \ a_2 \ \dots \ a_n]^T$ and $\vec{b} = [b_1 \ b_2 \ \dots \ b_n]^T$ are both vectors containing proportions that add to 1. As we've seen before, the TVD is defined as follows:

$$\text{TVD}(\vec{a}, \vec{b}) = \frac{1}{2} \sum_{i=1}^n |a_i - b_i|$$

The TVD is not the only metric that can quantify the distance between two categorical distributions. Here are three other possible distance metrics:

$$\text{dis1}(\vec{a}, \vec{b}) = \vec{a} \cdot \vec{b} = a_1 b_1 + a_2 b_2 + \dots + a_n b_n$$

$$\text{dis2}(\vec{a}, \vec{b}) = \frac{\vec{a} \cdot \vec{b}}{|\vec{a}| |\vec{b}|} = \frac{a_1 b_1 + a_2 b_2 + \dots + a_n b_n}{\sqrt{a_1^2 + a_2^2 + \dots + a_n^2} \sqrt{b_1^2 + b_2^2 + \dots + b_n^2}}$$

$$\text{dis3}(\vec{a}, \vec{b}) = 1 - \frac{\vec{a} \cdot \vec{b}}{|\vec{a}| |\vec{b}|}$$

Of the above three possible distance metrics, only one of them has the same range as the TVD (i.e. the same minimum possible value and the same maximum possible value) *and* has the property that smaller values correspond to more similar vectors. Which distance metric is it?

Solution:

Note: the solution here has a lot of inline LaTeX that's hard to copy — you can read the source solution here: <https://practice.dsc80.com/wi23-final/index.html>

FA23 Midterm Problem 4

Problem:

The `donkeys` table contains data from a research study about donkey health. The researchers measured the attributes of 544 donkeys. The next day, they selected 30 donkeys to reweigh. The first few rows of the `donkeys` table are shown below (left), and the table contains the following columns (right):

| | id | BCS | Age | Weight | WeightAlt |
|---|-----|-----|-----|--------|-----------|
| 0 | d01 | 3.0 | <2 | 77 | NaN |
| 1 | d02 | 2.5 | <2 | 100 | NaN |
| 2 | d03 | 1.5 | <2 | 74 | NaN |

| | |
|-----------|--|
| id | A unique identifier for each donkey (d01, d02, etc.). |
| BCS | Body condition score: from 1 (emaciated) to 3 (healthy) to 5 (obese) in increments of 0.5. |
| Age | Age in years: <2, 2–5, 5–10, 10–15, 15–20, and over 20 years. |
| Weight | Weight in kilograms. |
| WeightAlt | Second weight measurement taken for 30 donkeys. NaN if the donkey was not reweighed. |

Alan wants to see whether donkeys with `'BCS' >= 3` have larger `'Weight'` values on average compared to donkeys that have `'BCS' < 3`. Select **all the possible test statistics** that Alan could use to conduct this hypothesis test. Let μ_1 be the mean weight of donkeys with `'BCS' >= 3` and μ_2 be the mean weight of donkeys with `'BCS' < 3`.

- A. μ_1
- B. $\mu_1 - \mu_2$

- C. $2\mu_2 - \mu_1$
- D. $|\mu_1 - \mu_2|$
- E. Total variation distance
- F. Kolmogorov-Smirnov test statistic

Answer: B and C

- A: Incorrect. μ_1 does not tell compare the two groups, and so cannot be used to see which is larger on average.
- B: Correct. $\mu_1 - \mu_2$ tells us the difference between the average weight of both groups as well as which group would be larger (a test statistic greater than zero means μ_1 is larger).
- C: Correct. $2\mu_2 - \mu_1$ tells us the difference between the average weight of both groups as well as which group would be larger (a test statistic greater than μ_2 means μ_2 is larger).
- D: Incorrect. $|\mu_1 - \mu_2|$ tells us the difference between the average weight of both groups, but we cannot tell which group is larger due to the absolute value sign.
- E: Incorrect. Total variation distance is defined as $\frac{1}{2} \sum_{i=1}^k |a_i - b_i|$. This has the same issue as D where we cannot tell which group is larger due to the absolute value sign.
- F: Incorrect. Kolmogorov-Smirnov is a measurement of the maximum absolute difference between two cumulative distribution functions. It does not look at the average, nor does it tell us which weight would be larger.