

Activity_Explore hypothesis testing

June 3, 2024

1 Activity: Explore hypothesis testing

1.1 Introduction

You work for an environmental think tank called Repair Our Air (ROA). ROA is formulating policy recommendations to improve the air quality in America, using the Environmental Protection Agency's Air Quality Index (AQI) to guide their decision making. An AQI value close to 0 signals “little to no” public health concern, while higher values are associated with increased risk to public health.

They've tasked you with leveraging AQI data to help them prioritize their strategy for improving air quality in America.

ROA is considering the following decisions. For each, construct a hypothesis test and an accompanying visualization, using your results of that test to make a recommendation:

1. ROA is considering a metropolitan-focused approach. Within California, they want to know if the mean AQI in Los Angeles County is statistically different from the rest of California.
2. With limited resources, ROA has to choose between New York and Ohio for their next regional office. Does New York have a lower AQI than Ohio?
3. A new policy will affect those states with a mean AQI of 10 or greater. Will Michigan be affected by this new policy?

Notes: 1. For your analysis, you'll default to a 5% level of significance. 2. Throughout the lab, for two-sample t-tests, use Welch's t-test (i.e., setting the `equal_var` parameter to `False` in `scipy.stats.ttest_ind()`). This will account for the possibly unequal variances between the two groups in the comparison.

1.2 Step 1: Imports

To proceed with your analysis, import `pandas` and `numpy`. To conduct your hypothesis testing, import `stats` from `scipy`.

Import Packages

```
[1]: # Import relevant packages

    ### YOUR CODE HERE ###
```

You are also provided with a dataset with national Air Quality Index (AQI) measurements by state over time for this analysis. Pandas was used to import the file `c4_epa_air_quality.csv` as a dataframe named `aqi`. As shown in this cell, the dataset has been automatically loaded in for you. You do not need to download the .csv file, or provide more code, in order to access the dataset and proceed with this lab. Please continue with this activity by completing the following instructions.

Note: For purposes of your analysis, you can assume this data is randomly sampled from a larger population.

Load Dataset

```
[2]: # RUN THIS CELL TO IMPORT YOUR DATA.

### YOUR CODE HERE ###
aqi = pd.read_csv('c4_epa_air_quality.csv')
```

1.3 Step 2: Data Exploration

1.3.1 Before proceeding to your deliverables, explore your datasets.

Use the following space to surface descriptive statistics about your data. In particular, explore whether you believe the research questions you were given are readily answerable with this data.

```
[3]: # Explore your dataframe `aqi` here:

### YOUR CODE HERE ###
```

HINT 1

Consider referring to the material on descriptive statistics.

HINT 2

Consider using `pandas` or `numpy` to explore the `aqi` dataframe.

HINT 3

Any of the following functions may be useful: - `pandas`: `describe()`, `value_counts()`, `shape()`, `head()` - `numpy`: `unique()`, `mean()`

Question 1: From the preceding data exploration, what do you recognize? [Write your response here. Double-click (or enter) to edit.]

1.4 Step 3. Statistical Tests

Before you proceed, recall the following steps for conducting hypothesis testing:

1. Formulate the null hypothesis and the alternative hypothesis.
2. Set the significance level.

3. Determine the appropriate test procedure.
4. Compute the p-value.
5. Draw your conclusion.

1.4.1 Hypothesis 1: ROA is considering a metropolitan-focused approach. Within California, they want to know if the mean AQI in Los Angeles County is statistically different from the rest of California.

Before proceeding with your analysis, it will be helpful to subset the data for your comparison.

```
[4]: # Create dataframes for each sample being compared in your test

### YOUR CODE HERE ###
```

HINT 1

Consider referencing the material on subsetting dataframes.

HINT 2

Consider creating two dataframes, one for Los Angeles, and one for all other California observations.

HINT 3

For your first dataframe, filter to `county_name` of **Los Angeles**. For your second dataframe, filter to `state_name` of California and `county_name` not equal to **Los Angeles**.

Formulate your hypothesis: Formulate your null and alternative hypotheses:

- H_0 : There is no difference in the mean AQI between Los Angeles County and the rest of California.
- H_A : There is a difference in the mean AQI between Los Angeles County and the rest of California.

Set the significance level:

```
[5]: # For this analysis, the significance level is 5%

### YOUR CODE HERE
```

Determine the appropriate test procedure: Here, you are comparing the sample means between two independent samples. Therefore, you will utilize a **two-sample -test**.

Compute the P-value

```
[6]: # Compute your p-value here

### YOUR CODE HERE ###
```

HINT 1

Consider referencing the material on how to perform a two-sample t-test.

HINT 2

In `ttest_ind()`, `a` is the `aqi` column from our “Los Angeles” dataframe, and `b` is the `aqi` column from the “Other California” dataframe.

HINT 3

Be sure to set `equal_var = False`.

Question 2. What is your P-value for hypothesis 1, and what does this indicate for your null hypothesis? [Write your response here. Double-click (or enter) to edit.]

1.4.2 Hypothesis 2: With limited resources, ROA has to choose between New York and Ohio for their next regional office. Does New York have a lower AQI than Ohio?

Before proceeding with your analysis, it will be helpful to subset the data for your comparison.

```
[7]: # Create dataframes for each sample being compared in your test

    ### YOUR CODE HERE ###
```

HINT 1

Consider referencing the materials on subsetting dataframes.

HINT 2

Consider creating two dataframes, one for New York, and one for Ohio observations.

HINT 3

For your first dataframe, filter to `state_name` of New York. For your second dataframe, filter to `state_name` of Ohio.

Formulate your hypothesis: Formulate your null and alternative hypotheses:

- H_0 : The mean AQI of New York is greater than or equal to that of Ohio.
- H_A : The mean AQI of New York is **below** that of Ohio.

Significance Level (remains at 5%)

Determine the appropriate test procedure: Here, you are comparing the sample means between two independent samples in one direction. Therefore, you will utilize a **two-sample -test**.

Compute the P-value

```
[8]: # Compute your p-value here

### YOUR CODE HERE ###
```

HINT 1

Consider referencing the material on how to perform a two-sample t-test.

HINT 2

In `ttest_ind()`, `a` is the `aqi` column from the “New York” dataframe, and `b` is the `aqi` column from the “Ohio” dataframe.

HINT 3

You can assign `tstat`, `pvalue` to the output of `ttest_ind`. Be sure to include `alternative = less` as part of your code.

Question 3. What is your P-value for hypothesis 2, and what does this indicate for your null hypothesis? [Write your response here. Double-click (or enter) to edit.]

1.4.3 Hypothesis 3: A new policy will affect those states with a mean AQI of 10 or greater. Will Michigan be affected by this new policy?

Before proceeding with your analysis, it will be helpful to subset the data for your comparison.

```
[9]: # Create dataframes for each sample being compared in your test

### YOUR CODE HERE ###
```

HINT 1

Consider referencing the material on subsetting dataframes.

HINT 2

Consider creating one dataframe which only includes Michigan.

Formulate your hypothesis: Formulate your null and alternative hypotheses here:

- H_0 : The mean AQI of Michigan is less than or equal to 10.
- H_A : The mean AQI of Michigan is greater than 10.

Significance Level (remains at 5%)

Determine the appropriate test procedure: Here, you are comparing one sample mean relative to a particular value in one direction. Therefore, you will utilize a **one-sample -test**.

Compute the P-value

```
[10]: # Compute your p-value here

### YOUR CODE HERE ###
```

HINT 1

Consider referencing the material on how to perform a one-sample t-test.

HINT 2

In `ttest_1samp`, you are comparing the `aqi` column from your Michigan data relative to 10, the new policy threshold.

HINT 3

You can assign `tstat`, `pvalue` to the output of `ttest_1samp`. Be sure to include `alternative = greater` as part of your code.

Question 4. What is your P-value for hypothesis 3, and what does this indicate for your null hypothesis? [Write your response here. Double-click (or enter) to edit.]

1.5 Step 4. Results and Evaluation

Now that you've completed your statistical tests, you can consider your hypotheses and the results you gathered.

Question 5. Did your results show that the AQI in Los Angeles County was statistically different from the rest of California? [Write your response here. Double-click (or enter) to edit.]

Question 6. Did New York or Ohio have a lower AQI? [Write your response here. Double-click (or enter) to edit.]

Question 7: Will Michigan be affected by the new policy impacting states with a mean AQI of 10 or greater? [Write your response here. Double-click (or enter) to edit.]

2 Conclusion

What are key takeaways from this lab?

What would you consider presenting to your manager as part of your findings?

What would you convey to external stakeholders?

Congratulations! You've completed this lab. However, you may not notice a green check mark next to this item on Coursera's platform. Please continue your progress regardless of the check mark. Just click on the "save" icon at the top of this notebook to ensure your work has been logged.