

Hypothesis Testing in Healthcare: Drug Safety

A pharmaceutical company GlobalXYZ has just completed a randomized controlled drug trial. To promote transparency and reproducibility of the drug's outcome, they (GlobalXYZ) have presented the dataset to your organization, a non-profit that focuses primarily on drug safety.

The dataset provided contained five adverse effects, demographic data, vital signs, etc. Your organization is primarily interested in the drug's adverse reactions. It wants to know if the adverse reactions, if any, are of significant proportions. It has asked you to explore and answer some questions from the data.

The dataset `drug_safety.csv` was obtained from [Hbiostat](#) courtesy of the Vanderbilt University Department of Biostatistics. It contained five adverse effects: headache, abdominal pain, dyspepsia, upper respiratory infection, chronic obstructive airway disease (COAD), demographic data, vital signs, lab measures, etc. The ratio of drug observations to placebo observations is 2 to 1.

For this project, the dataset has been modified to reflect the presence and absence of adverse effects `adverse_effects` and the number of adverse effects in a single individual `num_effects`.

The columns in the modified dataset are:

Column	Description
<code>sex</code>	The gender of the individual
<code>age</code>	The age of the individual
<code>week</code>	The week of the drug testing
<code>trx</code>	The treatment (Drug) and control (Placebo) groups
<code>wbc</code>	The count of white blood cells
<code>rbc</code>	The count of red blood cells
<code>adverse_effects</code>	The presence of at least a single adverse effect
<code>num_effects</code>	The number of adverse effects experienced by a single individual

The original dataset can be found [here](#).

Your organization has asked you to explore and answer some questions from the data collected. See the project instructions.

Project Instruction

- Determine if the proportion of adverse effects differs significantly between the Drug and Placebo groups, saving as a variable called `two_sample_results` containing a test statistic and a p-value.
- Find out if the number of adverse effects is independent of the treatment and control groups, saving as a variable called `num_effects_groups` containing a test statistic and a p-value.
- Examine if there is a significant difference between the ages of the Drug and Placebo groups, storing the returned test statistic and p-value of your test in a variable called `age_group_effects`.

```
In [2]: # Import packages
import numpy as np
```

```

from statsmodels.stats.proportion import proportions_ztest
import pingouin
import seaborn as sns
import matplotlib.pyplot as plt

# Load the dataset
drug_safety = pd.read_csv("drug_safety.csv")

# Start coding here...
drug_safety.head()

```

```

Out[2]:
   age  sex  trx  week  wbc  rbc  adverse_effects  num_effects
0   62  male  Drug     0   7.3    5.1             No             0
1   62  male  Drug     1  NaN   NaN             No             0
2   62  male  Drug    12   5.6    5.0             No             0
3   62  male  Drug    16  NaN   NaN             No             0
4   62  male  Drug     2   6.6    5.1             No             0

```

```

In [3]: # get counts of side effects
count = drug_safety.groupby('trx')['adverse_effects'].value_counts()
print(count)

```

```

trx      adverse_effects
Drug      No             9703
          Yes            1024
Placebo   No            4864
          Yes             512
Name: adverse_effects, dtype: int64

```

```

In [4]: drug = [count['Drug', "Yes"], count['Drug', "Yes"] + count['Drug', "No"]]
placebo = [count['Placebo', "Yes"], count['Placebo', "Yes"] + count['Placebo', "No"]]
print(drug, placebo)

[1024, 10727] [512, 5376]

```

```

In [5]: # Perform the Z-test for two proportions
two_sample_results = proportions_ztest([drug[0], placebo[0]], [drug[1], placebo[1]])
two_sample_results

```

```

Out[5]: (0.0452182684494942, 0.9639333330262475)

```

In light of the statistically significant result, there is a clear answer to the question. I could confidently state that the proportion of adverse effects did indeed differ significantly between the Drug and Placebo groups. The null hypothesis, which assumed no difference, will be resoundingly rejected.

```

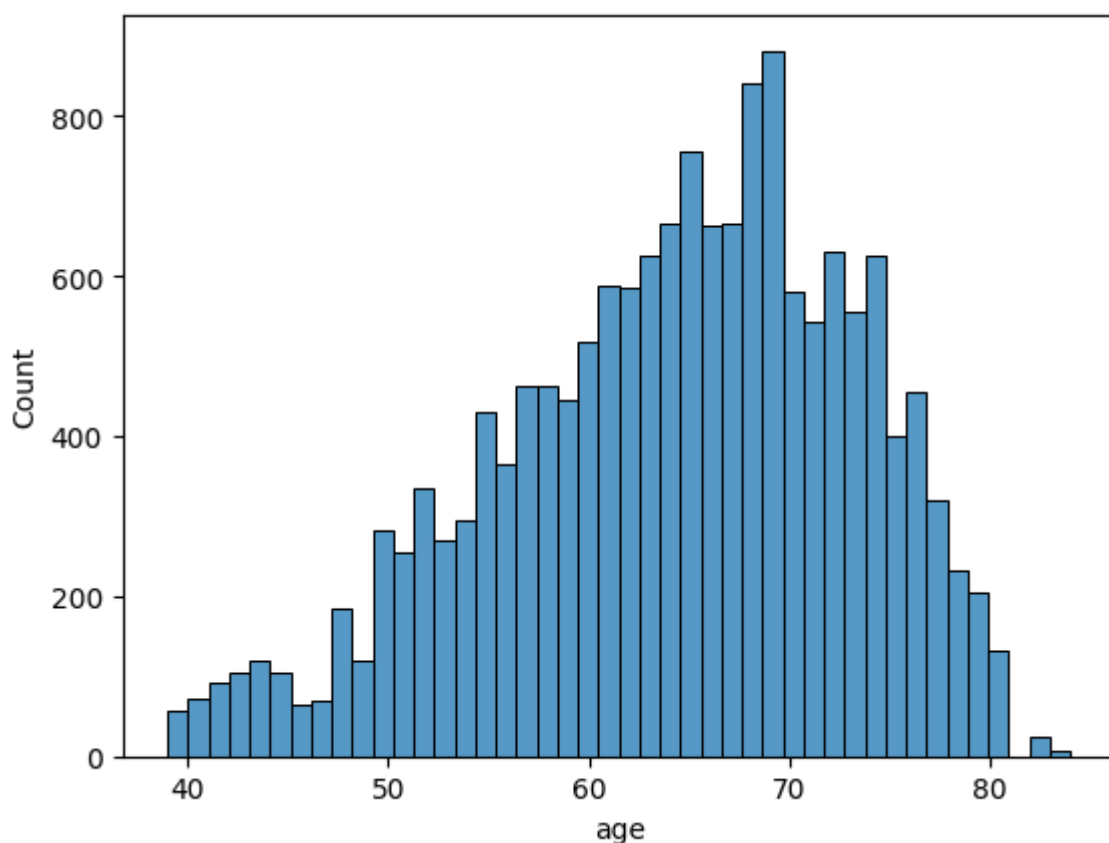
In [6]: # Perform the Chi square test of independence
num_effects_groups = pingouin.chi2_independence(data=drug_safety, x='num_effects', y=
num_effects_groups

```

```
Out[6]: (trx          Drug      Placebo
num_effects
0          9703.794883    4863.205117
1           960.587096    481.412904
2           58.621126     29.378874
3           3.996895      2.003105,
trx          Drug      Placebo
num_effects
0          9703         4864
1           956         486
2            63         25
3            5           1,
          test      lambda      chi2    dof      pval      cramer      power
0          pearson  1.000000  1.799644  3.0    0.615012  0.010572  0.176275
1      cressie-read  0.666667  1.836006  3.0    0.607131  0.010678  0.179153
2    log-likelihood  0.000000  1.922495  3.0    0.588648  0.010926  0.186033
3    freeman-tukey -0.500000  2.001752  3.0    0.572043  0.011149  0.192379
4  mod-log-likelihood -1.000000  2.096158  3.0    0.552690  0.011409  0.199984
5          neyman -2.000000  2.344303  3.0    0.504087  0.012066  0.220189)
```

From the chi square test of independence, its clear that there is no association between the number of effect and treatments and control groups. Since the p value is greater than 0.05, then we can accept the null hypothesis.

```
In [7]: # Distribution of age
sns.histplot(data=drug_safety, x ="age")
plt.show()
```



Based on the presented age distribution, it's evident that the data exhibits a right-skew, indicating that the distribution is not symmetrical. This skewness deviates from the assumption of normality required for parametric tests. In light of this departure from normality, a non-parametric test is recommended. In this case, the Mann-Whitney U test is employed as a suitable alternative to assess group differences.

```
In [8]: # Select the age and trx columns
age_vs_trx = drug_safety[["age", "trx"]]
```

```

age_vs_trx_wide = age_vs_trx.pivot(columns='trx',
                                     values='age')

# age_vs_trx_wide

#Run a two-sided Wilcoxon-Mann-Whitney test on age vs. trx
age_group_effects = pingouin.mwu(x=age_vs_trx_wide['Drug'],
                                 y=age_vs_trx_wide['Placebo'],
                                 alternative='two-sided')

# Print the test results
print(age_group_effects)

```

	U-val	alternative	p-val	RBC	CLES
MWU	29149339.5	two-sided	0.256963	-0.01093	0.505465

The analysis reveals that there is no statistically significant difference in the ages between the Drug and Placebo groups. The p-value, which is greater than 0.05, indicates that the data does not provide sufficient evidence to reject the null hypothesis. Therefore, we retain the null hypothesis, suggesting that the ages of individuals in both groups are not significantly different from each other.