

Statistical Inference

https://github.com/monastyrskyy/bootcamp_repo/tree/master/Capstone_1_Phishing

The barplot visualization shows that some features stray further from the sample proportion of phishing websites than others. To develop some hypotheses about what features will be influential in model selection, I first filtered down to the 'interesting' features with the following method. For each category of each feature, I ran a 10,000 size bootstrap sample of the 'result' label variable. I then took this information to calculate a 95% confidence interval for the proportion that phishing websites made up of that group. If the overall proportion of phishing websites in the whole sample was outside of the 95% CI for any of the categories of a feature, I labeled that feature as 'interesting'. The number of phishing websites that have these 'interesting' features will either be significantly higher or lower than the sample average, at 95% confidence. This might then suggest that these features will have a lot of influence when building models and may have high predictive power in whether a website is phishing or legitimate. Below is a list of the 'interesting' columns.

```
url_contains_double_slash
favicon
standard_ports
url_contains_https_token
submitting_to_email
redirect
on_mouseover
right_click_disabled
popup_window
uses_iframe
```

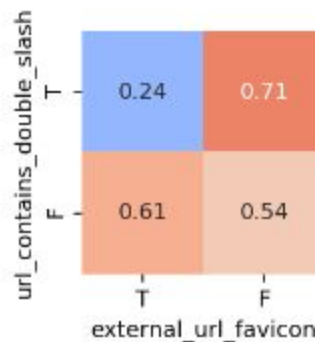
Merely looking at a list of feature names doesn't tell me too much, so the next step would be to visualize the features somehow.

Exploratory Data Analysis - Closer Look

To visualize the data within the 'interesting' columns, I want to display the proportion of phishing websites within each combination of categories for each combination of features. Below is a crosstab example of what I mean.

external_url_favicon	-1	1
url_contains_double_slash		
-1	0.240506	0.710692
1	0.614853	0.535176

In this crosstab, each combination of each category is plotted and the number it contains inside is the proportion of phishing websites of that group. Notice that this crosstab only compares two features: 'external_url_favicon' to 'url_contains_double_slash', whereas in my full analysis I want to compare each combination of features, and visually display the proportions of phishing websites.



The above graph is a visualization of the previously discussed crosstab. The color switches from shades of blue to shades of red, depending on whether the proportion of phishing websites for each specific combination in the visualization is above or below the sample proportion of phishing websites. The darker the color, the bigger the absolute difference between the combination proportion and the sample proportion.

To get the above information for each combination of features, I iterated over each combination of features and created a dictionary of such crosstab data frames as the one above but for every combination of features. Below is a visualization of each combination in a heatmap format. Notice that the combinations with dark colors are easy to see, and the trends are easy to see. For example among the websites that contain a double slash in the url, only 21% of those that also contain a popup window are phishing, as opposed to the 74% of those that don't contain a popup. This example shows that there are a lot of non-intuitive but interesting insights that this visualization can provide.

