

# Web Crawler Project - Feature List

Author: **Bahy Helmi Hartoyo Putra**  
(bahyhelmi97@gmail.com)

## Description

The goal was to create a set of features for the upcoming fraud detection model used in KYB automation process. Before producing these features, discussions have been done with several stakeholders to determine the characteristics of fraud merchants' website observed to this date. Features were extracted by the crawler machine through the crawling process, receiving the website's URL as an input.

## Summary

No	Feature	Description	Notes
1	<b>broken_link_score</b>	The percentage of link(s) with successful HTTP request (Response 200) out of all sampled hyperlinks.  <b>Values:</b> [0,100]	<ul style="list-style-type: none"><li>• Sampling was done with random sampling with n = 10.</li><li>• Hyperlinks defined as all links available on the homepage of the website.</li></ul>
2	<b>link_contact_us_exist</b>	Flag whether any contact link exists on the website.  <b>Values:</b> 0/1	-
3	<b>cu_email_exist</b>	Flag whether any email address exists on the website.  <b>Values:</b> 0/1	<ul style="list-style-type: none"><li>• Search spaces including homepage and the contact link itself.</li></ul>
4	<b>cu_phone_number_exist</b>	Flag whether any phone number exists on the website.  <b>Values:</b> 0/1	<ul style="list-style-type: none"><li>• Search spaces including homepage and the contact link itself.</li></ul>
5	<b>link_about_us_exist</b>	Flag whether any about us link exists on the website.	-

		<b>Values:</b> 0/1	
6	<code>link_tnc_exist</code>	Flag whether any terms & condition link exists on the website.  <b>Values:</b> 0/1	-
7	<code>tnc_refund_policy_exist</code>	Flag whether any refund policy exists on the website.  <b>Values:</b> 0/1	<ul style="list-style-type: none"> <li>Search spaces including homepage and the terms &amp; condition link itself.</li> </ul>

## Possible Improvements

- Adding more train data, especially for **REJECTED** merchants. It is good for the model to learn more about the variation of a fraud merchant website.
- Implement content prediction on several pages that considered as content-rich pages. The output of the prediction can be joined together with the current model to enhance the performance of the final prediction, using the ensemble learning method.
- Rediscuss the definition of a fraud merchant and the objective of developing a future model.

## Reference

- [Source Code](#)
- [Beautifulsoup](#)
- [Selenium with Python](#)
- [Model Validation Notebook](#)
- [Hyperparameter Tuning Notebook](#)