

Sure! Below is a **short explanation of each column** in your dataset and how to obtain the data for a given URL:

## ◊ Basic URL Features

Column Name	Explanation	How to Get Data?
<b>URL</b>	The full website link	Direct input
<b>LengthOfURL</b>	Total characters in the URL	<code>len(url)</code> Extract using <code>tldextract.extract(url).domain</code>
<b>Domain</b>	Main domain name	<code>tldextract.extract(url).domain</code>
<b>DomainLengthOfURL</b>	Length of domain name only	<code>len(domain)</code> Extract using <code>tldextract.extract(url).suffix</code>
<b>TLD</b>	Top-Level Domain (e.g., .com, .edu)	<code>tldextract.extract(url).suffix</code>
<b>TLDLength</b>	Length of the TLD	<code>len(tld)</code>

## ◊ Complexity & Special Characters

Column Name	Explanation	How to Get Data?
<b>URLComplexity</b>	Number of unique characters in the URL	<code>len(set(url))</code>
<b>CharacterComplexity</b>	Ratio of unique characters to total characters	<code>len(set(url)) / len(url)</code>
<b>dot_count</b>	Number of . in the URL	<code>url.count('.')</code>
<b>hyphen_count</b>	Number of - in the URL	<code>url.count('-')</code>
<b>slash_count</b>	Number of / in the URL	<code>url.count('/')</code>
<b>digit_count</b>	Number of digits (0-9) in URL	<code>sum(c.isdigit() for c in url)</code>
<b>URLLetterRatio</b>	Ratio of letters to total characters	<code>sum(c.isalpha() for c in url) / len(url)</code>
<b>URLDigitRatio</b>	Ratio of digits to total characters	<code>digit_count / len(url)</code>

<b>EqualCharCntInURL</b>	Count of = characters	<code>url.count('=')</code>
<b>QuesMarkCntInURL</b>	Count of ? in URL	<code>url.count('?')</code>
<b>AmpCharCntInURL</b>	Count of & in URL	<code>url.count('&amp;')</code>
<b>OtherSpclCharCntInURL</b>	Count of other special characters (%,_,#,etc.)	<code>sum(url.count(c) for c in ['%', '_', '#', '!', '\$'])</code>
<b>URLOtherSpclCharRatio</b>	Ratio of special characters in URL	<code>OtherSpclCharCntInURL / len(url)</code>

## ◊ Path & Query Analysis

Column Name	Explanation	How to Get Data?
<b>NumberOfHashtags</b>	Count of # in URL	<code>url.count('#')</code>
<b>NumberOfSubdomains</b>	Number of subdomains	<code>len(tldextract.extract(url).subdomain.split('.'))</code>
<b>HavingPath</b>	Whether the URL has a path (1 = Yes, 0 = No)	<code>1 if urlparse(url).path else 0</code>
<b>PathLength</b>	Length of the URL path after domain	<code>len(urlparse(url).path)</code>
<b>HavingQuery</b>	Whether the URL has a query string	<code>1 if urlparse(url).query else 0</code>
<b>HavingFragment</b>	Whether the URL has a fragment (#...)	<code>1 if urlparse(url).fragment else 0</code>
<b>HavingAnchor</b>	Whether the URL has an anchor link	<code>1 if '#' in url else 0</code>

## ◊ Security Features

Column Name	Explanation	How to Get Data?
<b>HasSSL</b>	Whether the URL uses HTTPS (1 = Yes, 0 = No)	<code>1 if url.startswith('https') else 0</code>

<b>IsUnreachable</b>	Whether the website is unreachable	Check response using `requests.get(url, timeout=5)`
<b>IsDomainIP</b>	Whether the domain is an IP address	<pre>1 if domain.replace('.','').isdigit() else 0</pre>

## ◊ Website Content & HTML Analysis

Column Name	Explanation	How to Get Data?
<b>LineOfCode</b>	Number of lines in the page source	<code>len(response.text.split('\n'))</code>
<b>LongestLineLength</b>	Length of the longest line in source code	<code>max(len(line) for line in response.text.split('\n'))</code>
<b>HasTitle</b>	Whether the page has a <code>&lt;title&gt;</code> tag	<code>1 if '&lt;title&gt;' in response.text else 0</code>
<b>HasFavicon</b>	Whether the page has a favicon	Check <code>&lt;link rel="icon"&gt;</code> in HTML
<b>HasRobotsBlocked</b>	Whether robots.txt is blocking access	Check <code>/robots.txt</code> response
<b>IsResponsive</b>	Whether the page is mobile-friendly	Check viewport meta tag in <code>&lt;head&gt;</code>
<b>HasDescription</b>	Whether meta description exists	Check <code>&lt;meta name="description"&gt;</code>
<b>HasPopup</b>	Whether the page has popups	Search for JavaScript popup functions ( <code>window.alert</code> , etc.)
<b>HasIFrame</b>	Whether the page contains <code>&lt;iframe&gt;</code> elements	<code>1 if '&lt;iframe&gt;' in response.text else 0</code>

## ◊ Redirection & Social Engineering

Column Name	Explanation	How to Get Data?
<b>IsURLRedirects</b>	Whether the URL redirects to another page	Check <code>response.history</code> in <code>requests.get(url)</code>
<b>IsSelfRedirects</b>	Whether the page redirects to itself	Compare final URL with initial URL

<b>IsFormSubmitExternal</b>	Whether a form submits data to an external domain	Extract action attribute in <form>
-----------------------------	---	------------------------------------

## ◊ Banking, Payment, and Phishing Indicators

Column Name	Explanation	How to Get Data?
<b>HasSocialMediaPage</b>	Whether the site has social media links	Check for facebook.com, twitter.com in HTML
<b>HasSubmitButton</b>	Whether the page has a submit button	Search <input type="submit">
<b>HasHiddenFields</b>	Whether forms contain hidden input fields	Search <input type="hidden">
<b>HasPasswordField</b>	Whether the page has a password field	Search <input type="password">
<b>HasBankingKey</b>	Whether banking-related words exist	Check for bank, account, secure, etc.
<b>HasPaymentKey</b>	Whether payment-related words exist	Check for pay, checkout, credit card, etc.
<b>HasCryptoKey</b>	Whether crypto-related words exist	Check for bitcoin, crypto, wallet, etc.
<b>HasCopyrightInfoKey</b>	Whether copyright information exists	Check for ©, copyright, etc.

## ◊ External Resources & Linking Behavior

Column Name	Explanation	How to Get Data?
<b>CntImages</b>	Number of images on the page	Count <img> tags
<b>CntFilesCSS</b>	Number of external CSS files	Count <link rel="stylesheet"> tags
<b>CntFilesJS</b>	Number of JavaScript files	Count <script> tags with src
<b>CntSelfReferf</b>	Number of self-referencing links	Count <a href> with same domain
<b>CntEmptyRef</b>	Number of empty links (href="#")	Count <a href="#">

<b>CntExternalRef</b>	Number of links pointing to external domains	Compare href domains to site domain
<b>CntPopup</b>	Number of JavaScript popups	Search for window.alert, window.prompt
<b>CntIFrame</b>	Number of iframe elements	Count <iframe> tags

## ◊ Advanced Security Metrics

Column Name	Explanation	How to Get Data?
<b>UniqueFeatureCnt</b>	Number of unique security features	Count nonzero security-related values
<b>WAPLegitimate</b>	Whether the site is a well-known web application	Compare against a list of trusted apps
<b>WAPP phishing</b>	Whether the site is a known phishing app	Compare against phishing database
<b>ShannonEntropy</b>	Measures randomness in URL (higher = more phishing)	Compute entropy formula
<b>FractalDimension</b>	Measures complexity of page structure	Compute fractal dimension of DOM tree
<b>KolmogorovComplexity</b>	Measures compressibility of HTML source	Compute using Lempel-Ziv complexity
<b>HexPatternCnt</b>	Number of hexadecimal patterns in URL	Count occurrences of [0-9A-Fa-f]{4,}
<b>Base64PatternCnt</b>	Number of Base64 encoded patterns	Count occurrences of [A-Za-z0-9+/]{10,}
<b>LikelinessIndex</b>	Probability of being phishing based on features	Model-derived score

## ◊ Label (Target)

Column Name	Explanation
<b>Label</b>	The actual classification (Legitimate/Phishing)

## Next Steps

- Use BeautifulSoup, requests, tldextract, and re to extract most of these features automatically.
- Train your **LightGBM** and **CatBoost** models on this dataset.

This will **improve accuracy and reduce false positives** for login-based legitimate websites! 