

In this project, I scrape blockchain data from Wikipedia using `BeautifulSoup`, clean the extracted data, and store it in a structured `pandas` DataFrame. Finally, I saved the data as a CSV file for further exploratory analysis.

```
[1]: #Import the necessary libraries
from bs4 import BeautifulSoup
import requests

[2]: url = 'https://en.wikipedia.org/wiki/List_of_blockchains' #Add the webpage to be scraped
page = requests.get(url)
soup = BeautifulSoup(page.text, 'html')

[5]: print(soup)

<div class="user-links-collapsible-item mw-list-item user-links-collapsible-item" id="pt-createaccount-2"><a class="" data-mw="interface" href="/w/index.php?title=Special:CreateAccount&returnto=List%2Fof%2Fblockchains" title="You are encouraged to create an account and log in; however, it is not mandatory"><span>Create account</span></a>
</div>
<div class="user-links-collapsible-item mw-list-item user-links-collapsible-item" id="pt-login-2"><a accesskey="o" class="" data-mw="interface" href="/w/index.php?title=Special:Userlogin&returnto=List%2Fof%2Fblockchains" title="You're encouraged to log in; however, it's not mandatory. [o]"><span>log in</span></a>
</div>
<div>
<div class="vector-dropdown vector-user-menu vector-button-flush-right vector-user-menu-logged-out" id="vector-user-links-dropdown" title="Log in and more options">
<input aria-haspopup="true" aria-label="Personal tools" class="vector-dropdown-checkbox" data-event-name="ui.dropdown-vector-user-links-dropdown" id="vector-user-links-dropdown-checkbox" role="button" type="checkbox"/>
<label aria-hidden="true" class="vector-dropdown-label cdx-button cdx-button-fake-button cdx-button-fake-button-enabled cdx-button-weight-quiet cdx-button-icon-only" for="vector-user-links-dropdown-checkbox" id="vector-user-links-dropdown-label"><span class="vector-icon mw-ui-icon-ellipsis mw-ui-icon">
[7]: #Find the table in html of the webpage
soup.find('table')

[7]: <table class="plainlinks metadata ambox mbox-small-left ambox-notice" role="presentation" style="width: auto;"><tbody><tr><td class="mbox-image"><span type="file"></span></td><td class="mbox-text" style="width: auto;"><div class="mbox-text-span">This list is <a href="/wiki/Wikipedia:WikiProject_Lists#Incomplete_lists" title="Wikipedia:WikiProject Lists">incomplete</a>; you can help by <a class="external text" href="https://en.wikipedia.org/w/index.php?title=List_of_blockchains&action=edit">adding missing items</a>. <span class="date-container"><div><span class="date"><span></span></div></div></td></tr>

[9]: soup.find_all('table')[1]

[9]: <table class="wikitable sortable">
<caption>
</caption>
<tbody><tr>
<th>Name</th>
<td>Date created</td>
<td>Created by</td>
<td>Native cryptocurrency</td>
<td>Consensus algorithm</td>
<td>Programmable?</td>
<td>Private?<sup class="reference" id="cite_ref-ISO_1-0"><a href="#cite_note-ISO-1"><span class="cite-bracket"></span></span>Note <span class="cite-bracket"></span></td></tr>

[11]: soup.find('table', class_ = 'wikitable sortable')
#table class="wikitable sortable jquery-tablesorter">
#caption>

[11]: <table class="wikitable sortable">
<caption>
</caption>
<tbody><tr>
<th>Name</th>
<td>Date created</td>
<td>Created by</td>
<td>Native cryptocurrency</td>
<td>Consensus algorithm</td>
<td>Programmable?</td>
<td>Private?<sup class="reference" id="cite_ref-ISO_1-0"><a href="#cite_note-ISO-1"><span class="cite-bracket"></span></span>Note <span class="cite-bracket"></span></td></tr>

[13]: table = soup.find_all('table')[1]

[15]: print(table)

<table class="wikitable sortable">
<caption>
</caption>
<tbody><tr>
<th>Name</th>
<td>Date created</td>
<td>Created by</td>
<td>Native cryptocurrency</td>
<td>Consensus algorithm</td>
<td>Programmable?</td>
<td>Private?<sup class="reference" id="cite_ref-ISO_1-0"><a href="#cite_note-ISO-1"><span class="cite-bracket"></span></span>Note <span class="cite-bracket"></span></td></tr>

[17]: soup.find_all('th') #Find the elements that have the "th" label

[17]: [<th>Name</th>,
<th>Date created</th>,
<th>Created by</th>,
<th>Native cryptocurrency</th>,
<th>Consensus algorithm</th>,
<th>Programmable?</th>,
<th>Private?<sup class="reference" id="cite_ref-ISO_1-0"><a href="#cite_note-ISO-1"><span class="cite-bracket"></span></span>Note <span class="cite-bracket"></span></th>]

[19]: block_titles = table.find_all('th')

[21]: print(block_titles)

<th>Name</th>
<th>Date created</th>
<th>Created by</th>
<th>Native cryptocurrency</th>
<th>Consensus algorithm</th>
<th>Programmable?</th>
<th>Private?<sup class="reference" id="cite_ref-ISO_1-0"><a href="#cite_note-ISO-1"><span class="cite-bracket"></span></span>Note <span class="cite-bracket"></span></th>
<th>Permissioned?<sup class="reference" id="cite_ref-ISO_1-1"><a href="#cite_note-ISO-1"><span class="cite-bracket"></span></span>Note <span class="cite-bracket"></span></th>
<th>Finality</th>
<th>Ledger state</th>
<th>Refs.</th>

[23]: block_table_titles = [title.text.strip() for title in block_titles]
print(block_table_titles)

['Name', 'Date created', 'Created by', 'Native cryptocurrency', 'Consensus algorithm', 'Programmable?', 'Private?', 'Permissioned?', 'Finality', 'Ledger state', 'Refs.

[25]: #I want to replace some column titles in the list
block_table_titles = [block_table_titles.replace('Refs.', '') for block_table_title in block_table_titles]
block_table_titles = [block_table_titles.replace('Private?<sup class="reference" id="cite_ref-ISO_1-1"><a href="#cite_note-ISO-1"><span class="cite-bracket"></span></span>Note <span class="cite-bracket"></span></th>', 'Private?') for block_table_title in block_table_titles]
block_table_titles = [block_table_titles.replace('Permissioned?<sup class="reference" id="cite_ref-ISO_1-1"><a href="#cite_note-ISO-1"><span class="cite-bracket"></span></span>Note <span class="cite-bracket"></span></th>', 'Permissioned?') for block_table_title in block_table_titles]
#I also need to remove the links
print(block_table_titles)

['Name', 'Date created', 'Created by', 'Native cryptocurrency', 'Consensus algorithm', 'Programmable?', 'Private?', 'Permissioned?', 'Finality', 'Ledger state', 'Notes', '']

[31]: column_data = table.find_all('tr')

[33]: for row in column_data[1:]: #removed first empty bracket
row_data = row.find_all('td')
# Extract text from each table cell (td) instead of using row tag objects
new_row_data = [td.get_text(strip=True) for td in column_data] #this extracts the text from each cell and shows it as text not a BeautifulSoup tag
print(new_row_data)

individual_row_data = [data.get_text(strip=True) for data in row_data] #Removing this shows error "cannot set rows with mismatched column", adding it back
print(individual_row_data)
# what you get are all individual lists, so we need to find a way to get them into our table
# we cannot stuff everything into our data frame at once. So let's put it in one at a time. We will use Location in pandas

length = len(df) #Looking at length of current df, which is empty right now
df.loc[length] = individual_row_data #we are looping through individual row data and appending each row of info into next Location

[35]: df

2 on Ethereum[7][8]", "StarkNet?StarkWareSTRKZK-rollup?estayer-2 on Ethereum[9]", "Polygon zkEVM?Polygon TechnologyETHZK-rollup?eszkEVM-compatible?layer-2 on Ethereum", "Ethereum ClassicJul 20, 2016ETC?PoW?Yes?No?split from Ethereum due to The DAOhack earlier that month[1]", "Bitcoin CashAug 1, 2017BCH?PoW[10]split from Bitcoin", "CardanoSep 27, 2017Charles Hoskinsonand Jeremy WoodADAPoS?Yes?No?ProbabilisticUTXO[1][11][12]", "TRONJun 24, 2018TRON?PoS?Yes?No[11][13]", "TezosJun 30, 2018Arthur and Kathleen BreitmanXTZ?PoS?Yes?No[1][14]", "Bitcoin SVNov 2018BSV?PoW?Yes (scripts)Hosplit from Bitcoin Cash, itself split from Bitcoin[1]", "Lightning Network[relevant?2018?on?layer-2 on Bitcoin]", "XinFinJune 1, 2019XinFin Fintech, XDC FoundationXDC?PoS?Yes?No?Immediate?oXDC Network is an layer 1 EVM compatible, environmental friendly, near zero transaction cost with high speed settlement blockchain platform.", "AlgorandJun 10, 2019Silvio Micaliand othersALGO?PoS?Yes?No?Immediate?uses verifiable random functionto choose random validators for consensus[15]", "SPC?Sep 16, 2019SPC?PoS?Yes (scripts)No?PoS?PC future well-known blockchain[16]", "SolanaMarch 16, 2020Anatoly Yakovenkoand Raj GokalSOL?PoS?with Proof of History (PoH)?Yes?No?Immediate?Account-balance[17]", "OasisJune 18, 2020Oasis Labs,Oasis Protocol FoundationROSE?PoS?Yes (in Paratimes)No?No?Immediate?Account-balance?Paratime? can useTEEforconfidential computing[18][19]", "PolkadotMay 26, 2020Parity TechnologiesDOT?Started withPoAthen moved toPoS?Yes (in parachains)No?No?bridged?Account-balance[20][21]", "AvalancheSeptember 10, 2020Emin Gün Sirer, Maofan "Ted" Yin andKevin SekniqiAVAX?PoS?Yes (in C-chain)No?No?Immediate?UTXO[22]", "MolochCoinDec 6, 2020MolochCoin Inc. (founded byJosh Goldbergand Shane Glynn)MOC", "Internet Computer?IO?IOITY Foundation(founded byDominic Williams)[23]Computation is very cheap; can host websites", "DESOLan 18, 2021Nader al-Najdi (akadiamondhands) and othersDESOL (formerlyBTC.LT, CLOUT)social media; flagship p appBittCloud; name acquired in Sep 2021[23]", "Terra Classic?Do Kwonaand othersLUNC (formerlyLUNA), USDFormerly Terra until May 2022; ecosystem collapsed in May 2022 (UST depegged to near-zero and LUNA also went to near-zero)", "Terra 2.0May 28, 2022LUNANew blockchain created following the collapse of Terra. [24]", "StellarApr 6, 2016XLM?BFT?Yes?Yes[4]", "EOS.I0Jul 1, 2017EOS?DPoS?Yes?No?[4]", "LBRY?LBC", "RippleJune 2012Ripple LabsXRPBFT?No?No?Immediate?Account-balance?Blockchain is known as XRP Ledger. Smart contract capabilities are being added.[25][26][27]", "Stacks?STX", "VertcoinJan 8, 2014VTC", "Hedera HashgraphHBAR?Yes?Yes?Account-balance?uses directed acyclic graphinstead of a chainper se", "ZcashOct 28, 2016[4]ZEC?PoW?Yes?No?zero-knowledge proofsfor privacy[28]", "MoneroAug 18, 2014XMR?PoW?Yes?Yes[13][14]", "Bitcoin Cash?BCH?PoW?Yes (scripts)BCH forked from Bitcoin[13]", "Dogecoin?DOGE?PoW?Yes?No?No?Initially created as a joke, but gained popularity[29]"

[35]: df

[35]:
```

- We successfully scraped blockchain data from Wikipedia.
- Cleaned the extracted table by handling inconsistent row lengths.
- Stored the data in a structured pandas DataFrame.

- Perform **exploratory data analysis (EDA)** on the dataset.
- Visualize blockchain data trends using **matplotlib** or **seaborn**.