



MATES Computer Science

Senior Capstone Project Bi-Weekly Progress Report

| | |
|-------------------------|-----------------------------------------------------------------------------------------------------|
| Project Title | MorningBread |
| Team Members | Max Bradshaw, Jake Dorick |
| Dates Covered by Report | March 9th - March 22nd |
| Link to Github | https://github.com/maxbshaw17/MorningBread |

1. **Summary of Project** (Provide a one paragraph summary of your project. You can largely copy/paste this from one progress report to the next, unless there are significant changes.)

To provide a personalized morning financial report to our users by utilizing a web scraping tool for financial website articles to find the general topics of the articles. In doing so, we are creating and designing a website to give the public a general grasp of financial news for the day, over the course of the day with hourly, highly summarized news headlines/stock tickers. An account system will allow users to follow certain companies and manage tickers.

2. **Summary of Progress this Period** (Provide a high-level, one paragraph overview of what was accomplished this progress period collectively by the team.)

Max -

I implemented the NLP text clustering using the DBSCAN density-based algorithm. I also worked on the SQL injection code and cleaned up the scraping methods again, as well as completely rewriting the yahoo scraping function again (in progress). Finally, I structured the back-end databases and created a test database for Jake to work off of for now.

Jake -

Creating the background for what is to come to be dynamically-added soon was my main priority. The design of how the articles and tickers will be laid out on the front page was my top priority. The structure of my code is bound to change as the JavaScript is implemented, allowing divs to be dynamically-added as they are fetched from the API, but the CSS will remain the same.

3. **Detailed Progress this Period, separated by Team Member** (Provide detailed information on the progress that you made in the reporting weeks. Include screenshots of code, your game or website, etc. Each team member should have a separate subsection covering their accomplishments. Not including screenshots, this section should be 1-2 pages.)

Max -

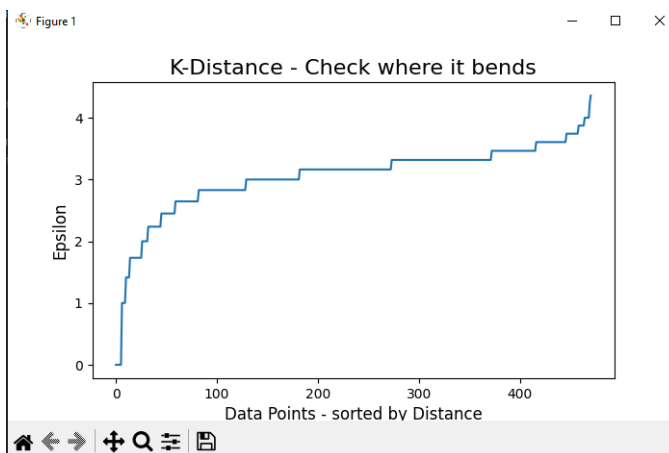
- Wrote functions for ML so I can call in the file that I will be running automatically
 - This took the longest as I tried several models before choosing this one

```
2 from sklearn.feature_extraction.text import CountVectorizer
3 from sklearn.cluster import DBSCAN
4 from sklearn.neighbors import NearestNeighbors
5 import pandas as pd
6 import numpy as np
7 import matplotlib.pyplot as plt
8 import string
9 from nltk.tokenize import word_tokenize, sent_tokenize
10 from nltk.corpus import stopwords
11 from nltk.stem import PorterStemmer
12
13 default_stemmer = PorterStemmer()
14 default_stopwords = stopwords.words('english') # or any other list of your choice
15
16 # check for repeats - need to fix on sql side
17 def remove_repeats(headlines):
18     text_no_repeats = []
19     for sent in headlines:
20         if sent not in text_no_repeats:
21             text_no_repeats.append(sent)
22     return (text_no_repeats)
23
24 > def clean_text(text, ):
25     return text
26
27
28
29 def clean_text_list(text_list):
30     text_list = remove_repeats(text_list)
31     return list(map(lambda text: clean_text(text), text_list))
32
33
34 def knn_plot(text_array):
35     neigh = NearestNeighbors(n_neighbors=2)
36     nbrs = neigh.fit(text_array)
37     distances, indices = nbrs.kneighbors(text_array)
38     distances = np.sort(distances, axis=0)
39     distances = distances[:,1]
40
41     plt.figure(figsize=(7,4))
42     plt.plot(distances)
43     plt.title('K-Distance - Check where it bends',fontsize=16)
44     plt.xlabel('Data Points - sorted by Distance',fontsize=12)
45     plt.ylabel('Epsilon',fontsize=12)
46     plt.show()
47
48
49 def text_vectorizer(text_list):
50     text = remove_repeats(text_list)
51     text_cleaned = clean_text_list(text_list)
52
53     cv = CountVectorizer()
54     x = cv.fit_transform(text_cleaned)
55     return x.toarray(), text
56
57
58 def fit_dbscan_text(text_array, text, ep=1, min_s=2):
59     dbscan_opt=DBSCAN(eps = ep,min_samples = min_s)
60     dbscan_opt.fit(text_array)
61     df_data = {"group": dbscan_opt.labels_, "sentence" : text}
62     df = pd.DataFrame(data=df_data)
63
64     return df
```

- Added a duplicate removal call to the database

```
13 def remove_dupes_sql(table, column):
14     try:
15         command = f"""DELETE FROM
16         {table}
17         WHERE
18         id IN (
19             SELECT
20                 id
21             FROM
22                 (
23                     SELECT
24                         id,
25                         ROW_NUMBER() OVER (
26                             PARTITION BY {column}
27                             ORDER BY
28                                 {column}
29                             ) AS row_num
30                     FROM
31                         {table}
32                 ) t
33             WHERE
34                 row_num > 1
35         );"""
36
```

- Because DBSCAN is based on the multidimensional distance between each vectorized text, I have to fine-tune the distance parameter
- I made a quick function that plots the number of groups vs each epsilon (distance) input
 - 2.5 seems to be the sweet spot



- Script to run ML and example of groupings

```

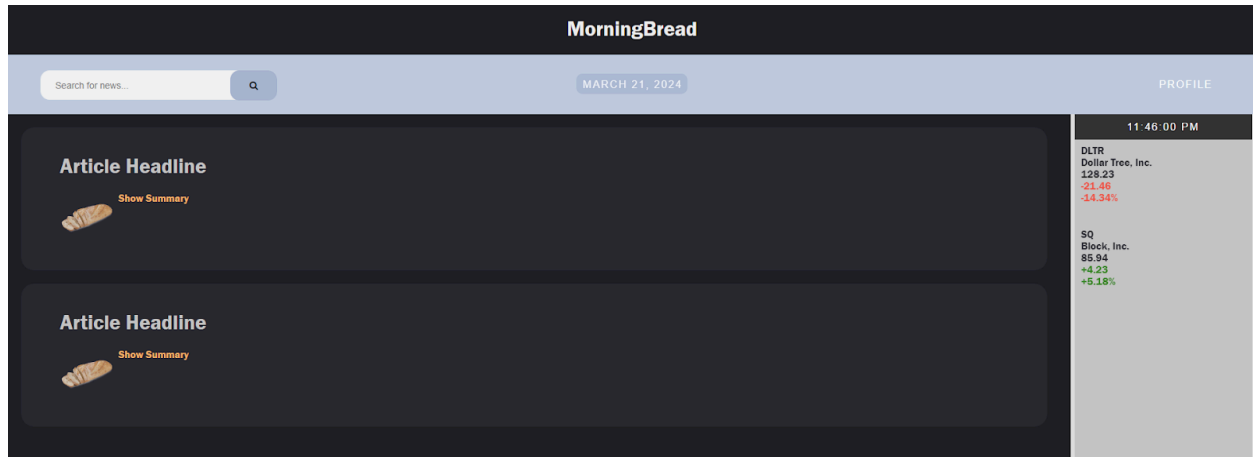
1 from db_update.db_connector import *
2 from nlp_algo.nlp_functions import *
3
4 headlines = get_all_headlines()
5
6 array, sents = text_vectorizer(headlines)
7
8
9 print(knn_plot(array))
10
11
12 fit_df = fit_dbscan_text(array, sents, ep = 2.5, min_s=2)
13
14 print(fit_df.sort_values(by='group', ascending=False).head(80).to_string())
15
16 #for i in fit_df.sort_values(by='group', ascending=False).index:
17 #    if fit_df['group'][i] != -1:
18 #        print(f"{fit_df['group'][i]}, {fit_df['sentence'][i]}")

```

| PROBLEMS | OUTPUT | DEBUG CONSOLE | TERMINAL | PORTS | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|----------|----------|---------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------|----------|-----|----|-----|----|-----|----|-----|----|-----|----|-----|----|-----|----|-----|----|-----|----|-----|----|-----|----|-----|----|-----|----|-----|----|-----|----|-----|----|-----|----|-----|----|-----|----|-----|----|-----|----|-----|----|-----|----|-----|----|-----|----|-----|----|-----|----|-----|----|-----|----|-----|----|-----|----|-----|----|--|
| | | | <table border="1"> <thead> <tr> <th>group</th> <th>sentence</th> </tr> </thead> <tbody> <tr><td>467</td><td>24</td></tr> <tr><td>388</td><td>24</td></tr> <tr><td>386</td><td>24</td></tr> <tr><td>385</td><td>23</td></tr> <tr><td>458</td><td>23</td></tr> <tr><td>399</td><td>22</td></tr> <tr><td>378</td><td>22</td></tr> <tr><td>418</td><td>21</td></tr> <tr><td>376</td><td>21</td></tr> <tr><td>375</td><td>20</td></tr> <tr><td>396</td><td>20</td></tr> <tr><td>407</td><td>19</td></tr> <tr><td>366</td><td>19</td></tr> <tr><td>401</td><td>18</td></tr> <tr><td>361</td><td>18</td></tr> <tr><td>381</td><td>17</td></tr> <tr><td>348</td><td>17</td></tr> <tr><td>339</td><td>16</td></tr> <tr><td>351</td><td>16</td></tr> <tr><td>242</td><td>15</td></tr> <tr><td>236</td><td>15</td></tr> <tr><td>323</td><td>14</td></tr> <tr><td>250</td><td>14</td></tr> <tr><td>315</td><td>14</td></tr> <tr><td>233</td><td>14</td></tr> <tr><td>228</td><td>13</td></tr> <tr><td>308</td><td>13</td></tr> <tr><td>225</td><td>12</td></tr> <tr><td>231</td><td>12</td></tr> <tr><td>207</td><td>11</td></tr> <tr><td>198</td><td>11</td></tr> <tr><td>192</td><td>10</td></tr> </tbody> </table> | group | sentence | 467 | 24 | 388 | 24 | 386 | 24 | 385 | 23 | 458 | 23 | 399 | 22 | 378 | 22 | 418 | 21 | 376 | 21 | 375 | 20 | 396 | 20 | 407 | 19 | 366 | 19 | 401 | 18 | 361 | 18 | 381 | 17 | 348 | 17 | 339 | 16 | 351 | 16 | 242 | 15 | 236 | 15 | 323 | 14 | 250 | 14 | 315 | 14 | 233 | 14 | 228 | 13 | 308 | 13 | 225 | 12 | 231 | 12 | 207 | 11 | 198 | 11 | 192 | 10 | |
| group | sentence | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 467 | 24 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 388 | 24 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 386 | 24 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 385 | 23 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 458 | 23 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 399 | 22 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 378 | 22 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 418 | 21 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 376 | 21 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 375 | 20 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 396 | 20 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 407 | 19 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 366 | 19 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 401 | 18 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 361 | 18 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 381 | 17 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 348 | 17 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 339 | 16 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 351 | 16 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 242 | 15 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 236 | 15 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 323 | 14 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 250 | 14 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 315 | 14 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 233 | 14 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 228 | 13 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 308 | 13 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 225 | 12 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 231 | 12 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 207 | 11 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 198 | 11 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 192 | 10 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |

- You can see that sentences with similar topics are grouped together

Jake -



- The layout for the articles (and adding the areas for the tickers) has been changed
- JavaScript is implemented in both the articles section and the tickers
 - The articles section has a “Show Summary” button that will give the article’s information when clicked
 - This requires reading the button’s status based off of its original status when the page is first loaded



- Above shows the general layout the articles will have for website (probably going to remove the goofy bread icon)
- The tickers on the right are scrolling, they will also be dynamically-added, but that is my second priority
- What’s currently there is just an example of what I hope to implement after I get the articles onto the site
- Time function because that’s important for stocks, as so I’m told

```
from flask import Flask, jsonify
import mysql.connector

app = Flask(__name__)

# Connect to the MySQL database (Not yet created)
articles_db = mysql.connector.connect(
    host = "mysql-2ed0e70f-morningbread.a.aivencloud.com",
    user = "avnadmin",
    password = "AVNS_-1y1cgAxePfkqdPTpji",
    port = 25747,
    database = "morningbread",
)

c = articles_db.cursor()

@app.route('/articles_tickers_api/articles_api', methods=['GET'])
def get_articles():
    # Query the database for article headlines and summaries
    c.execute("SELECT headline, summary FROM articles")
    articles = c.fetchall()

    # Convert the articles to a list of dictionaries
    article_list = [{'headline': row[0], 'summary': row[1]} for row in articles]

    return jsonify(article_list)

if __name__ == '__main__':
    app.run(debug=True)
```

- I have used a similar concept of the Flask API for my research project (PneumoVision), so I was able to recycle a lot of this code
 - articles_db connects to the MySQL database to begin fetching the information that will be added to the website
 - I have a similar file to the one above for the tickers despite there being no ticker data in the database yet...
-

```
from flask import Flask, jsonify
import mysql.connector

app = Flask(__name__)

# Connect to the MySQL database (Not yet created)
tickers_db = mysql.connector.connect(
    host = "mysql-2ed0e70f-morningbread.a.aivencloud.com",
    user = "avnadmin",
    password = "AVNS_-1y1cgAxePfkqdPTpji",
    port = 25747,
    database = "morningbread"
)

c = tickers_db.cursor()

@app.route('/articles_tickers_api/tickers_api', methods=['GET'])
def get_tickers():
    # Query the database for the ticker's symbol, name, change, and percent change
    c.execute("SELECT symbol, name, change, percent_change FROM tickers")
    tickers = c.fetchall()

    # Convert the tickers to a list of dictionaries
    tickers_list = [{'headline': row[0], 'summary': row[1]} for row in tickers]

    return jsonify(tickers_list)

if __name__ == '__main__':
    app.run(debug=True)
```

- Only difference is variable names and what the cursor is selecting for in the database
- An external JavaScript file (all of my current JavaScript is in the HTML files [probably going to want to change]) was created to call the API and summon it into the website

```
fetch('/articles_tickers_api/articles_api')
    .then(response => response.json())
    .then(articles => {
        const dynamicArticlesContainer = document.getElementById('dynamic-articles');

        // Loop through the articles and display them on the page
        articles.forEach(article => {
            const articleElement = document.createElement('div');
            articleElement.classList.add('article-preview'); // Add a class for styling
            articleElement.innerHTML = `
                <h1>${article.headline}</h1>
                <p>${article.summary}</p>
            `;
            dynamicArticlesContainer.appendChild(articleElement);
        });
    })
    .catch(error => console.error('Error fetching articles:', error));
```


4. **Difficulties Encountered this Progress Period** (Provide detailed information on the difficulties and issues that you encountered in the reporting weeks. Discuss mitigation strategies for how you got around or plan to get around these issues.)

Max -

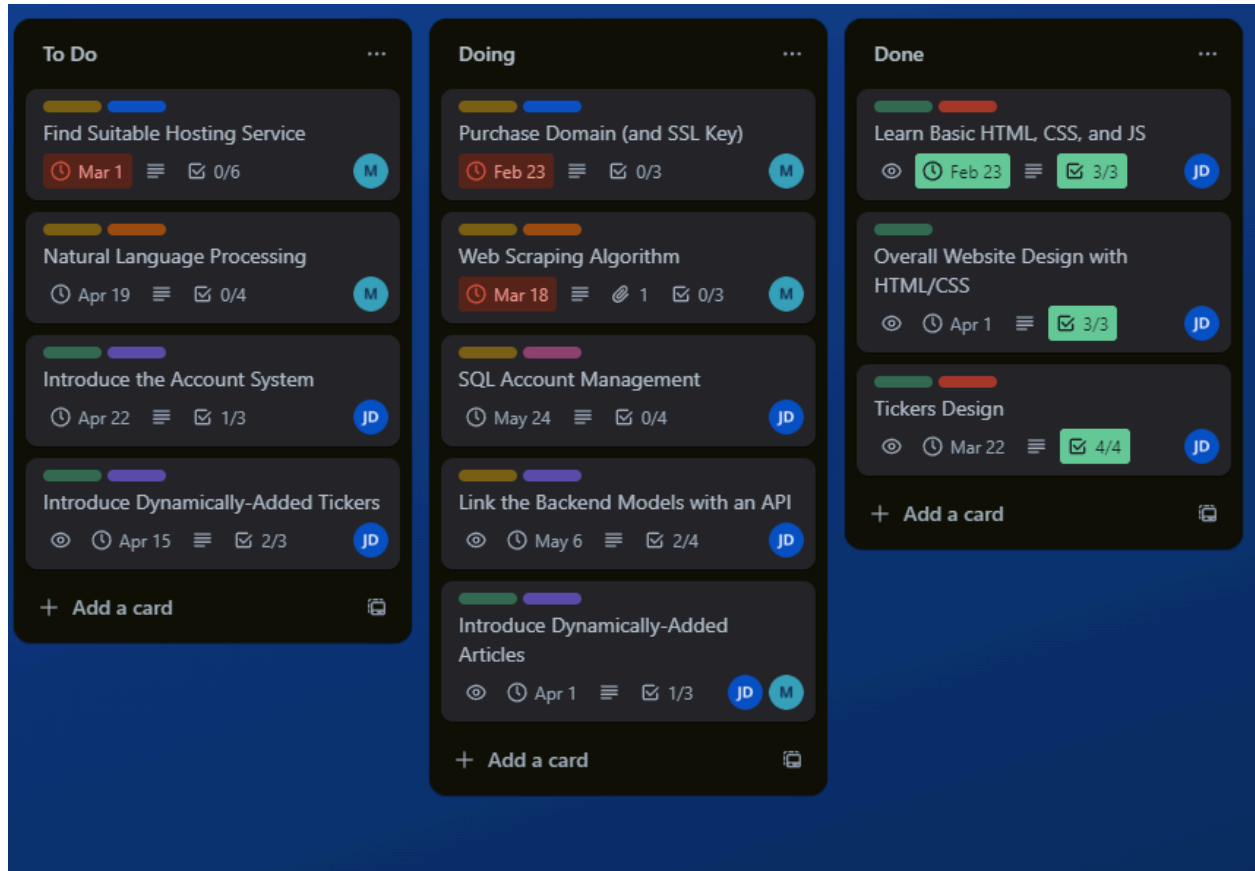
- Database stuff was difficult before I had mysqlworkbench
- Lots of random version issues
- It was tricky choosing which model I would use

Jake -

- Styling decisions have become difficult to determine as that is no longer in the foreground of my time
- Understanding the interactions between the API and how the dynamically-added divs are going to build on top of one another
- Sign up form is still going all the way to the bottom of the page instead of cutting off like the profile page despite them having the exact same CSS? It just looks a bit weird but it is by no means a top priority

5. **Updated Trello Board and Discussion** (Provide screenshot of and link to updated Trello board. Discuss any changes made to board since last progress report and why.)

<https://trello.com/b/sSuMPdYn/morningbread>



I (Jake) have been updating the tasks and marking things as done when they are done, as well as adding more detail to the cards that already exist.

6. **Tasks to Be Worked on in Next Progress Period** (Discuss the tasks to be worked on in the following two weeks. Discuss who is working on each.)

Max -

- Finally finish the yahoo scrape function (almost there for real this time) and clean up scraping
- Automate upload to SQL
 - Upload articles automatically
 - Delete old article on upload
- Use chat gpt to combine grouped headlines into one

Jake -

- If everything goes smoothly, articles stored in the MySQL database will be presented on the front page of our website
- Ensuring hourly updates is the next task as I'm not even sure if Max's database updates at all (I need to talk to him, I know)
- Keeping all aesthetics of the website when information is added

- I really want to make an accordion-styled About Us page, but I know that shouldn't be worthy of my attention right now
 - I'll probably do it over the weekend/Spring Break just for fun

7. **Additional Information** (Provide any additional information that you want to provide in this section; for example, one of your teammates is going away next week, your Github account is gone, etc. It could be good news as well.)

Max -

I feel like I am very ahead of where I thought I would be. I think it is likely that I finish my original goals in the next two sprints.

Jake -

No comments. My research project didn't win anything at the science fair, but I think the skills in API that I gained from doing it will translate well to MorningBread.
