

DISTRIBUTED SOCIAL NEWS

by

STEFANOS CHATZAKIS

URN: 6481123

A dissertation submitted in partial fulfilment of the
requirements for the award of

BACHELOR OF SCIENCE IN COMPUTER SCIENCE

May 2020

Department of Computing
University of Surrey
Guildford GU2 7XH

Supervised by: Sotiris Moschoyiannis

I declare that this dissertation is my own work and that the work of others is acknowledged and indicated by explicit references.

Stefanos Chatzakis
May 2020

© Copyright Stefanos Chatzakis, May 2020

Abstract

This report describes the design and implementation process for a web-based system that provides its users with a personalised view of the top social news depending on their interests. Moreover, in order to achieve this, the report will also focus on implementing some of the most popular ranking algorithm which will then be compared with each other in order to select the most appropriate for the project.

This report will loosely be split into two parts. The first being, understanding the strengths and weaknesses of each ranking algorithm and comparing them to find the most suitable. The second part will focus entirely on the implementation using an agile development methodology.

All in all the aim of this report will be to explain the challenges I face during the comparison process as well as the development process.

Acknowledgements

First of all, I would like to thank my parents Fedon and Katerina for all the sacrifices they have made for me to have this opportunity today, but also their continuous support and motivation. Furthermore, I would like to thank my sister and my close friends who have helped me cope with the difficulties I have faced during the past three years.

Contents

1	Introduction	10
1.1	Background	10
1.2	Aims & Objectives	10
1.2.1	Project Limitations	11
1.3	Software Development Life Cycle	11
1.4	Chosen Methodology	12
1.5	Report Structure	13
1.6	Summary	13
2	Literature Review	14
2.1	Technologies Used	14
2.1.1	Django Framework	14
2.1.2	Atom	15
2.2	Existing Solutions	15
2.2.1	Reddit Rank (Hot Sort)	15
2.2.1.1	Weaknesses	16
2.2.2	Hacker News Rank	16
2.2.2.1	Weaknesses	17
2.2.2.2	Conclusion	17

2.2.3	Wilson Confidence Sort	17
2.2.3.1	Conclusion	18
2.2.4	Summary	18
3	The Algorithm	19
3.1	Planning	19
3.2	Designing	20
3.3	Creating	21
3.4	Complexity	22
4	System Analysis	23
4.1	Feasibility and System Limitations	23
4.2	User Experience	24
4.3	System Requirements	25
4.3.1	Hosting Requirements	25
4.3.2	System Functionality	26
5	System Design	27
5.1	Prototype	27
5.2	User Interface Design and Usability	28
5.2.1	Home Page	29
5.2.2	Account Page	29
5.3	Difficulties faced	30
5.3.1	Pre-development	30
5.3.2	Development	30
6	Testing	31

7 Statement of Ethics	33
7.1 Leagal standpoint	33
7.2 Ethical standpoint	34
7.2.1 ACM-IEEE	34
7.2.2 The Computer Misuse Act 1990	34
7.2.3 Copyright	35
7.2.4 Conclusion	35
8 Conclusion	36
8.1 Overview	36
8.2 Potential Improvements	37
8.3 Personal Gains	37

List of Figures

2.1	Reddit Ranking Algorithm in mathematical notation	15
2.2	Hacker News Algorithm in mathematical notation	16
2.3	Edwin B. Wilson confidence sort in mathematical notation	17
3.1	The stages of the final ranking algorithm.	20
4.1	Bootstrap colours	24
4.2	An example post	24
5.1	Django app file layout	27
5.2	Account view methods	27
5.3	Home page from mobile device	28
5.4	Home page from desktop	28
5.5	Account page pt. 1	29
5.6	Account page pt. 2	29
6.1	Demonstration of scoring posts by time	31
6.2	Demonstration of scoring posts by common interests	31
6.3	Demonstration of scoring posts by non-common interests	32
6.4	Demonstration of scoring posts by the use of penalised words	32

Glossary

AI Artificial Intelligence.

NN Neural Network.

Chapter 1

Introduction

1.1 Background

Current social news sites try to produce a list of current important news items by applying algorithms and voting by the community. While more decentralized than a traditional publication, the system still relies on a select few(moderators, administrators) to safeguard it from manipulation (spam, vote-rigging) while maintaining the culture of the website. Further, by providing a centralized view of the current news, they necessarily average the preferences of their members instead of addressing the needs of each one based on their voting patterns. This project aims to develop a new type of social news aggregator that allows each user to see a personalized view of current events while allowing them to customize the algorithm by which the items get sorted.

1.2 Aims & Objectives

The solution, which will first be theorized and made into an algorithm, will first be comprised of the combination of the most critical aspects of the most popular sorting and ranking solutions out there while tailoring the final solution to one that fits a social news aggregator.

1. Develop a Functioning website with users that can create, edit their posts and view all of the users post ranked by timestamp.
2. Compare the most commonly used algorithms, critique the positive and negative attributes and theorize, which parts of each algorithm should be implemented in the final implement-

tation.

3. Compare the relevancy of the results after the final version of the algorithm has been implemented.
4. Create a modern and responsive website that is easy to use and compatible with a variety of screen resolutions, in order for the functionality of the program to be more clear.

1.2.1 Project Limitations

Since the development of this project was initiated, the goal was to create a testing platform on which to compare the results of the implementation of the final algorithm in a real-world environment. Thus, the scalability of the project as a whole was not taken into consideration. Having that said, the issue of scalability would occur not because of the algorithm itself, but due to the reason that the software developed does not address many security issues as well as correct optimisation so that the system would be able to hold on with more users.

1.3 Software Development Life Cycle

Software Development Life Cycle (SDLC) is the process that a software project follows to develop and test the system. A Software Development Life Cycle can be thought of as a master plan that describes how the development, testing, enhancement and maintaining processes should occur during the implementation of the software (Stackify 2017).

The goal of the SDLC is to produce a software with the highest quality at the lowest possible cost. In order to achieve this, SDLC provides six phases that are consistent with the most popular methodologies such as waterfall, spiral and Agile (Stackify 2017).

The SDLC phases include:

- Requirement analysis
- Planning
- Software design such as architectural design
- Software development

- Testing
- Deployment

The most traditional method for software development is the Waterfall methodology. Waterfall, delivers the above phases in a linear fashion, meaning that each phase can begin only after the previous has been completed. Moreover, thanks to the implementation of the phases, the waterfall methodology becomes a fairly straightforward process. However, a process can never be repeated unless a critical error occurs that requires the revision of a previous phase. This makes it incredibly risky to use in a final year project development environment since fully following the methodology, could prove difficult to address the arising errors that are bound to occur.

A more common methodology in the recent years is the Agile development approach. This approach takes the phases previously mentioned and creates the so called sprints. The sprint's duration usually is a couple of weeks long with pre-specified goals set at the beginning of each sprint. Each sprint can be thought of as a many smaller implementation of the waterfall approach, since the phases are consistent with both methodologies but since agile is comprised of sprints, during each sprint each phase will re-occur.

1.4 Chosen Methodology

After considering most of the commonly used methodologies, and due to time limitations, I decided that I was not going to use a methodology that was too complicated or had too many restrictions. Having that said, I chose the Agile development methodology using the Kanban framework. The goal of using this framework is to aid the developer in the categorisation and prioritisation of the work and tasks. Furthermore, the agile methodology was chosen since before the development; the technologies used were previously known. Additionally, the development will be broken down into small increments which will encourage the mastery of the technologies used after each sprint. Lastly, the sprints will allow for overall better testing and feedback since they will be received after each sprint.

1.5 Report Structure

The report structure will mirror the planning and development process as it occurred. First, by reviewing and constructing the ideal ranking algorithm for a social news aggregator following the development of the web application which will implement the constructed algorithm.

1.6 Summary

Overall the introduction chapter has portrayed an overview of the project's definition, flow and goals as well as addressing how the project will deal with the different phases of development through defining the methodology.

Chapter 2

Literature Review

2.1 Technologies Used

During the entire design and implementation process, the decision was made to use the Python programming language. This choice was made originally for the ease of implementing algorithms in a scripting language like python. However, while progressing from the testing phase to the development phase, a dilemma was faced. Either to re-implement most of the progress I had made to a language that was more web and browser-based like JavaScript or, to create a hybrid of both front-end and back-end schemes. After some research online, Django was used, a high-level Python Web framework that encourages rapid development and clean, pragmatic design.

2.1.1 Django Framework

According to the official Django website, "Django is a framework that was built as a tool for front-end developers that needed a simple way to bring their ideas to life without the need of a back-end developer that handles processes such as creating and connecting the server-side with the client-side as well as the creation and handling of the system's database (Editorial-Team 2016). Furthermore, Django is a high-level Python framework that utilizes SQLite as a relational database. Moreover, Django uses an MVT pattern similar to the more widely used MVC that frameworks like Ruby and Ruby on Rails uses. Lastly, Django is secure. It has one of the best out-of-the-box security systems out there, and it helps developers avoid common security issues, including clickjacking, cross-site scripting and SQL injection.

2.1.2 Atom

Since the whole project was entirely in Django, it meant that all the code was compiled through the console. For the text and file editor, Atom was used for its familiarity and simplicity. Moreover, for testing and displaying the application, Google Chrome was used.

2.2 Existing Solutions

2.2.1 Reddit Rank (Hot Sort)

After conducting research, it was clear that one of the most popular social news aggregators currently is without a doubt, Reddit. Reddit, although it has now switched to a different ranking algorithm, since the end of 2010 but have since made their old ranking system available to the public. Their algorithm explained in the purest form takes many parameters such as the timestamp of a post, the difference between a post's upvotes (likes) and downvotes (dislikes) and inputs those parameters in a formula that outputs a final rating that dictates a post's position compared to others.

Given the time the entry was posted A and the time of 7:46:43 a.m. December 8, 2005 B , we have t_s as their difference in seconds

$$t_s = A - B$$

and x as the difference between the number of up votes U and the number of down votes D

$$x = U - D$$

where $y \in \{-1, 0, 1\}$

$$y = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x = 0 \\ -1 & \text{if } x < 0 \end{cases}$$

and z as the maximal value, of the absolute value of x and 1

$$z = \begin{cases} |x| & \text{if } |x| \geq 1 \\ 1 & \text{if } |x| < 1 \end{cases}$$

we have the rating as a function $f(t_s, y, z)$

$$f(t_s, y, z) = \log_{10} z + \frac{yt_s}{45000}$$

Figure 2.1: Reddit Ranking Algorithm in mathematical notation

2.2.1.1 Weaknesses

An issue which is absent currently Reddit, but present in a lot of websites that take average rating as a rating attribute, is the following: Average rating works fine if you always have a ton of ratings, but suppose item 1 has 2 positive ratings and 0 negative ratings. Suppose item 2 has 100 positive ratings and 1 negative rating. This algorithm puts item two (tons of positive ratings) below item one (very few positive ratings).

2.2.2 Hacker News Rank

Hacker News is one of the most popular social news aggregator targeted towards developers and provides its users with mostly technology related news. Their ranking consists of three parameters, penalties, votes and age. The most impactful parameter on the formula is the penalty. The penalty's value is determined by the use of blacklisted words such as "NSA" which drops the story rapidly in the ranking. In order to keep the top stories fresh, Hacker News also issues a severe penalty on stories that reach 40 comments. The impact of a penalty can be calculated with the scoring formula since if an article gets penalty factor of 0.4, each vote will now count as 0.3. However, a factor of 0.1 is equivalent to each vote, counting 0.05. Meaning that although a penalty factor of 0.4 would drop an article 66% faster than usual, a factor of 0.1 would drop an article by 3.6 times than normal. In outline, each item is given a ranking, and the articles are sorted according to the ranking. The simplistic way to think about ranking is the number of votes is divided by time, so more votes results in a higher ranking, but the ranking also drops over time. The votes are raised to a power less than one, while the time is raised to a power greater than one, so time has more effect than votes. Some additional penalties also may be applied to the ranking.

$$rank = \frac{(score - 1)^{.8}}{(age_{hours} + 2)^{1.8}} * penalties$$

Figure 2.2: Hacker News Algorithm in mathematical notation

2.2.2.1 Weaknesses

1. Wall-clock hours penalize an article even if no one is reading (overnight, for example). A time denominated in ticks of actual activity (such as views of the 'new' page, or even upvotes-to-all-submissions) might address this.
2. An article that misses its audience first time through, perhaps due to the aforementioned reason or a bad headline may never recover, even with a later flurry of votes far beyond what the new submissions are getting.

2.2.2.2 Conclusion

Overall, Hacker News's algorithm is quite simple, thus making the implementation of something similar not that difficult, especially since it takes into account many factors. However, there are drawbacks with using timestamps as addressed above, and a solution to this would be to use time denominated in ticks of actual activity (such as views of the 'new' page, or even upvotes-to-all-submissions) which might address this issue. The use of penalty is an interesting concept that would make sense in a user generated content system in which there are no administrators to regulate the content.

2.2.3 Wilson Confidence Sort

The problem that the Wilson confidence score tried to resolve was, developers who need their application's users to sort things and to balance the proportion of positive ratings with the uncertainty of a small number of observations. The math to solve this problem was worked out in 1927 by Edwin B. Wilson (Miller 2009). Considering only positive and negative ratings (i.e. not a 5-star scale), the lower bound on the proportion of positive ratings is given by:

$$\left(\hat{p} + \frac{z_{\alpha/2}^2}{2n} \pm z_{\alpha/2} \sqrt{[\hat{p}(1 - \hat{p}) + z_{\alpha/2}^2/4n]/n} \right) / (1 + z_{\alpha/2}^2/n).$$

Figure 2.3: Edwin B. Wilson confidence sort in mathematical notation

2.2.3.1 Conclusion

Despite a very sought after solution to modern-day sorting, the Wilson Confidence score is to some degree not oftenly used. This very complicated solution only increases the efficiency of the results by a small margin. Furthermore, many problems occur in the voting system that the Wilson Confidence sort is based on. A simple example is that many people who find something mediocre will not bother to rate it at all; viewing or purchasing something and declining to rate; it contains useful information about that item's quality (Miller, Feb. 2009).

2.2.4 Summary

To summarise, most of these algorithms mentioned are tailored for their particular use or website. However, most of the existing solutions are bound by the user's interaction with the content to improve or worsen the overall score of a post. What this project is trying to propose, is a system where the only interaction that is needed by the reader is their interests, after which the algorithm takes care of the rankings. Nonetheless, we can assimilate parts of existing systems that can help polish the final version of the algorithm.

Chapter 3

The Algorithm

3.1 Planning

Most ranking algorithms as have been discussed, mainly revolve around some combination of voting and time. While some have more complicated parameters, they are usually tailored to the website's philosophy, such as the use of the penalty attribute in Hacker Rank. However, since this project aims to provide the user with an automated display of the most relevant news tailored to the user's interests, voting will be replaced with the user's interests.

The user should be able to select interests in their account page, which have a score associated with them and that score should be a significant contributor in the score of a post. When a post gets created is also be a valuable attribute, and it should be beneficial to the algorithm since time can never be the same even if the rest of the attributes are the same, differentiating in that way two potentially similar posts. Another vital aspect to consider is how many interests does the post have, in which the user is not interested. However, the impact of this attribute should not be as negatively equal to a relevant interest in a post. For example, if the score of the common interest between user and post is equal to 10, a non-common interest should impact the overall score by -5 or lower. Lastly, the concept of banning or negatively reducing a ranking score is an essential attribute and can be detrimental in the image and environment that a system portrays.

3.2 Designing

After a couple of failed attempts, the final algorithm was a combination time, common interests, non-common interests and penalties. The algorithm relies on multiplication and division to increase or decrease a rating for the sake of simplicity bundled with the fact that time which when denoted in timestamps usually tends to be a couple of digits long in the first place. The algorithm can be seen below:

$score = t_u ,$
where t_u is the timestamp that the post was updated.

$$score = \left(\sum_{i=0}^{n-1} a_i * 5 \right) * score ,$$

where a_i are the common interests of the post and the user.

$$score = \frac{score}{\left(\sum_{i=0}^{n-1} b_i * 2 \right)} ,$$

where b_i are the non - common interests of the post and the user.

$$score = \frac{score}{\left(\sum_{i=0}^{n-1} c_i * 3 \right)} ,$$

where c_i are the number of unique penalty words in the post's content.

Figure 3.1: The stages of the final ranking algorithm.

Since time was going to be used as timestamps which usually tend to be a couple of digits long, I decided to use multiplication and division to increase or decrease a rating was used. As portrayed in 3.1, if the user is not logged in or, they have no interest selected, the score of each post is set equal to the time that the post was updated. However, if the user is logged in and has selected interests, then the existing score is multiplied by 5 for each common interest of the user and the post.

For the deduction part of the algorithm, the score is halved for each non-common interest. This decision makes sense since, if a potential user is interested in news, a post with only the "news"

interest will rank higher than a post with the news and music interest unless the first post is much older than the second one. The last deduction part of the algorithm is penalties. Penalties, that at the moment are profanity words, essentially counter a common interest, since the score of a post is divided by 5 for each penalized word the post contains.

3.3 Creating

Assigning each post with a score based on the parameters aforementioned happens when the "get_blog_queryset" function is called in "home_screen_view", and the code featured below explains is part of the "get_blog_queryset" function.

```
# Checks if the user is authenticated.

if user.is_authenticated: #1
    pen_post = [] #1

    for penalty in penalty_list: #n
        penalized_posts = BlogPost.objects.filter(
            Q(title__contains=penalty) |
            Q(body__icontains=penalty)) #n
        for q in queryset: #n
            if q in penalized_posts:
                q.bad_words = q.bad_words + 1

    for q in queryset:#n
        # A list of common interests.
        common_result = []#1
        # A list of non-common interests.
        non_common_result = []#1
        # Set's the score to the time a post is updated.
        q.post_score = q.date_updated.timestamp() #1

        for element in q.interests: #n
            # Updates based on common interests.
```

```

if element in user.interests: #1
    common_result.append(element) #1
    q.post_score = q.post_score * (len(common_result) * 5) #1

# Updates based on non-common interests.

if element not in user.interests: #1
    non_common_result.append(element) #1
    q.post_score = q.post_score / (len(non_common_result) * 2) #1

# Updates based on penalty words.

for post in pen_post: #n
    q.post_score = q.post_score / (q.bad_words * 5) #1

```

3.4 Complexity

The time complexity of this algorithm will be calculated with the popular big O notation. Big O notation is used in Computer Science to describe the performance or complexity of an algorithm. Big O specifically describes the worst-case scenario, and can be used to describe the execution time required or the space used (e.g. in memory or on disk) by an algorithm (Bell 2009).

In the code listed in the above section, I have included the running complexity after each line. That way we can calculate that the worst-case complexity will never be more than $O(n^2)$.

Chapter 4

System Analysis

4.1 Feasibility and System Limitations

The project was designed as a demonstration to show how a let-alone, but personalized social news aggregator would work. The system flagship features of competing systems, all while mitigating issues such as the need of moderators to manage common issues such as vote-rigging or out-of-place content that disrupt the culture of the website.

Even though the system achieves the project's scope, that is not to say the system has reached its potential. There are many solutions that could perhaps increase the project's quality. For example, in order to improve sorting over time, the user could set their original interests but, over time a NN could monitor the user's actions and alter their interests judging from the way they interact with posts, such as time spent reading on a post that had other interests that the user did not initially mark as desirable. That way, the AI would create some hidden ranking that would sort the posts by that rank. It would not come with complications; nevertheless, building and running such a network would require not only immense dedicated computing power but also, periodic monitoring and safety regulation.

Changing the project's objectives and looking into lesser-known existing websites and their solutions in the literature review could likely provide more information on how these other systems handle sorting and what intricacies are involved when sorting information. Making the project more complex would defeat the goal of the project and could result in an overwhelming experience for the user. Moreover, in order for this project to yield meaningful results and an enjoyable experience results, responsiveness is key, which by overloading the system could result

in the lack of those aspects.

4.2 User Experience

In order to create a user-friendly interface and experience, the front end of the website is going to use Boostrap. Bootstrap is a free and open-source CSS framework directed at responsive, mobile-first front-end web development. It contains CSS- and (optionally) JavaScript-based design templates for typography, forms, buttons, navigation, and other interface components (Contributors 2019). This allowed me to make a very clean and friendly user interface.

Bootstrap uses default colours for different types of importance as visible:

```
.text-primary  
.text-secondary  
.text-success  
.text-danger  
.text-warning  
.text-info  
.text-light  
.text-dark  
.text-muted  
.text-white
```

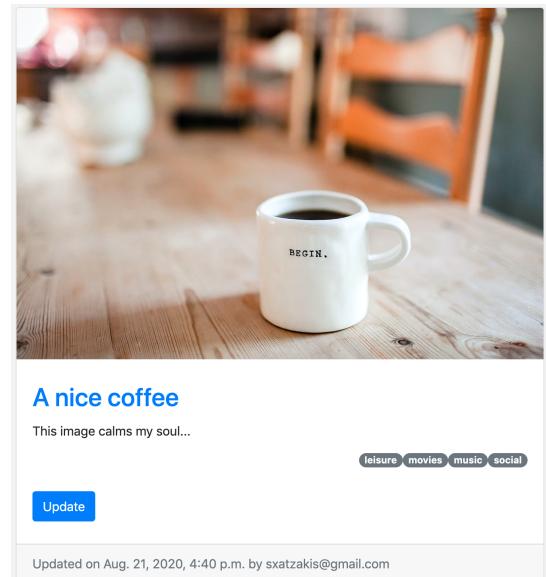


Figure 4.1: Bootstrap colours

Figure 4.2: An example post

After taking a few approaches, the final design I came up with was as shown in 4.2. Even though this is just a design, it should be easily linked to the back-end system. It consists of an image, a title, a body, an update button (if the post is by the user viewing it), several tags for the interests, and mutable text at the bottom with information on the post. The title is an important aspect of the post, meaning the blue colour was used to show importance. Next, the body text is "plain text" meaning it was left as is. The interest tags are not interactable meaning, that the secondary colour of Bootstrap was used. The update button is a call to action; thus, used the primary colour. Lastly, the information below the post should be informative, and

thus the muted colour was used. This style should be kept consistent throughout the website, giving the user a clear picture of what actions they can take at any point.

4.3 System Requirements

The project will be entirely written in Django; the back-end consists of Python and SQL, but the front-end is HTML meaning that the website should be available by any browser regardless of system specifications. Yet, in order to host the project, the host is required to meet a couple of specifications.

4.3.1 Hosting Requirements

Django includes requirements.txt files that includes the versions of the modules that it uses. In this case:

1. Django==2.2.2
2. Pillow==7.1.1
3. pytz==2019.3
4. sqlparse==0.3.1

Specifications

1. A modern operating system (Windows 7 or 10, Mac OS X 10.11 or higher 64-bit, Linux: RHEL 6/7 64-bit)
2. x86 64-bit CPU (Intel/AMD architecture)
3. 4-GB RAM
4. 5GB free disk space

This list includes the requirements in order for the project to run smoothly. The most important requirement for a system to host this project is for Python to be runnable; once this requirement is met, the project will execute regardless of the specifications. Another issue to consider is that the Disk space requirement for this project will scale depending on the number of posts.

4.3.2 System Functionality

The user should be able to view the home page with all the blog posts sorted by time. The user can choose to register an account with an email, username and password. If they are logged in, a personalised view of the home page is displayed to them. Furthermore, the user has the option to view the account details, where they can edit their personal information, and interests, and view all the posts they have made. Also, when creating a post, the user should be able to enter a title, body content, choose an image as well as the post interests. Moreover, the user should be able to view other user's posts but only be able to edit posts made by them. Lastly, the user should be able to log out at any point.

Chapter 5

System Design

5.1 Prototype

The prototype mostly consisted of the user being able to create an account, and create posts that simply included a text. As stated before, since I decided to use Django to develop this app, I had to work with apps. This means that every major aspect of the project such as blog posts and account of is implemented as an app. During this process, the agile design methodology played a big role in the constant reimplementations. In order to design and implement the different parts of the prototype such as user accounts, the agile methodology made it easier to concentrate on one thing at a time through heaps.

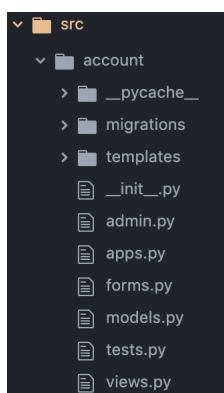


Figure 5.1: Django app file layout

```
def registration_view(request):=#
|#
def logout_view(request):=#
|#
def login_view(request):=#
|#
def account_view(request):=#
|#
def must_authenticate_view(request):=#
|#
```

Figure 5.2: Account view methods

In Figure 5.2 are the methods implemented during the prototype process that handled the account authentication.

5.2 User Interface Design and Usability

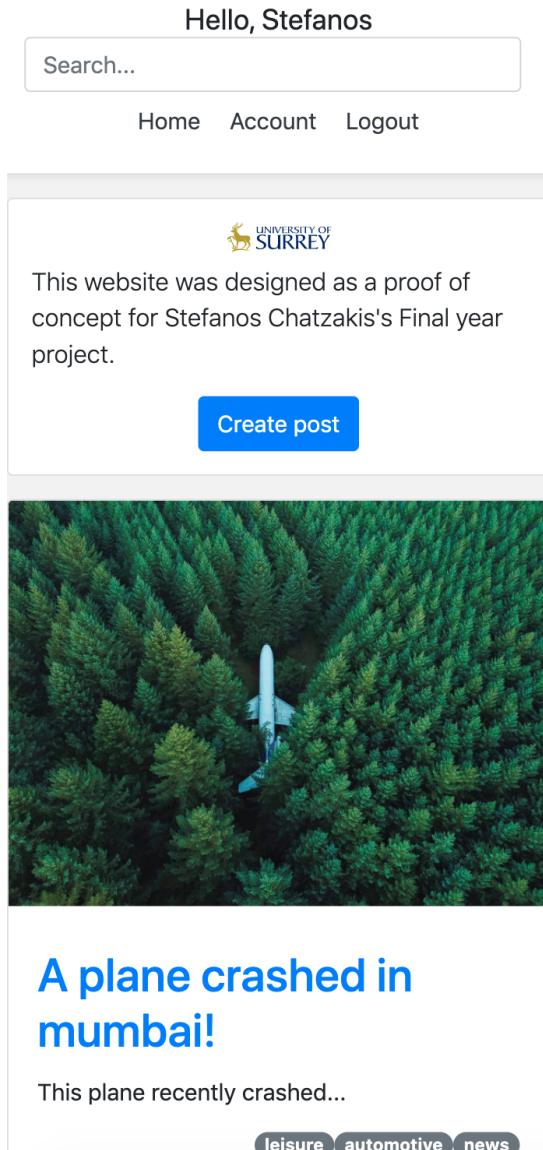


Figure 5.3: Home page from mobile device

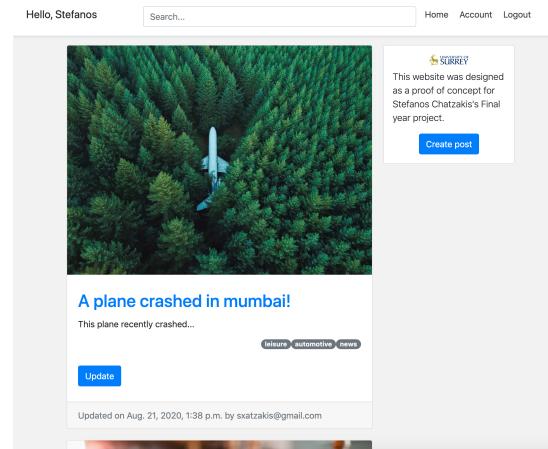


Figure 5.4: Home page from desktop

During the system analysis, I touched upon the user experience, both in how the user should interact with the colours and buttons and when a certain colour should be used. This section explains how those principles were implemented in the final version of the project. The project is designed to work in a plethora of screen sizes as seen on 5.3 and 5.4 making the project responsive giving the algorithm the chance to stand out, while providing an enjoyable experience to the user.

5.2.1 Home Page

5.4 shows the home page view of the website. It consists of the template mentioned in the system analytics, and the rest of the website follows along with those characteristics. The page header, includes a welcoming message, a search bar and three buttons which vary, depending on if the user is logged in or not. These buttons include: "Home", "Account" or "Login" and "Logout" or "Register". Then we can see the main content, which includes the blog post that is implemented as mentioned before and the create post section which again, follows the user experience philosophy mentioned before.

5.2.2 Account Page

Hello, Stefanos

Home Account Logout

Account

Change your Email

Change your Username

GAMES
 SOCIAL
 MUSIC
 NEWS
 FINANCE
 MOVIES
 SPORTS
 TRAVEL
 AUTOMOTIVE

TRAVEL
 AUTOMOTIVE
 LEISURE

Save changes

[Change password](#)

Blog posts:

- [Second Post](#)
- [The first post](#)
- [Test1](#)
- [Yesterday's match](#)
- [Legos!!!](#)
- [A nice coffee](#)
- [Photography with lightbulbs <3](#)
- [Is this orange real?](#)
- [A cold winter night with a movie!](#)
- [A plane crashed in mumbai!](#)

Figure 5.5: Account page pt. 1

Figure 5.6: Account page pt. 2

5.5 and 5.6 display the account page. Apart from the header which hasn't changed, the main page content has changed, and several fields are visible. In this page, the user can edit their personal information (which is securely handled by Django's authentication system); as well

as configure their interests and view all the posts they have made. This page also follows the consistent design decisions mentioned before.

5.3 Difficulties faced

This part will mainly touch upon difficulties faced during the development process of designing and implementing the program.

5.3.1 Pre-development

During the planning process, the choice to use Django was certainly not one without issues. For example, there isn't much information online to create a form that handles the information the way I needed them to and how to render them effectively. These issues mostly have to be learned through trial and error with the Django documentation being the sole guide. This resulted in spending a lot of time to learn all the Django components before the implementation could begin. Still, since many big organisations such as Reddit run their website entirely on the Django framework, it made me confident that the learning curve was worth it.

5.3.2 Development

During the development process, I did come across a few issues, mostly regarding keeping track of versions from different dependencies and programs which sometimes resulted in restarting the project. During the designing process, there were parts of the project that needed reimplementation for a variety of reasons, like UI, and sorting issues. Reflecting back to these issues, it is more apparent that the reimplementation process of different parts of the algorithm, provided a valuable learning opportunity which led to me writing more efficient code overall.

Chapter 6

Testing

In order to test the results of the algorithm and the implementation of the program, I created a test user with two different posts so that we can compare the results. First, to test if the two posts get correctly sorted by time, I will set the account interests to Leisure and Travel. I then will create "Test Post 1" and "Test Post 2" both with only the Leisure interest.

```
Test Post 1  
7991256993.23702  
Test Post 2  
7991255791.706605
```

Figure 6.1: Demonstration of scoring posts by time

As we can see from figure 6.1, Post 2 scores higher since it was created after Post 1.

In order to test the next segment of the algorithm which are common interests, I will change the interests of Post 1 to "news" and "leisure" and keep Post 2 as is.

```
Test Post 1  
79912592046.8192  
Test Post 2  
7991255791.706605
```

Figure 6.2: Demonstration of scoring posts by common interests

The results are visible in figure 6.2, where even though Post 1 is older, it has one more common interest with the user, thus ranking Post 1 higher. This might not have been the case where Post 1 to be much older than Post 2.

Another element to examine is whether a post with one common and one non-common interest ranks lower than a with just one common interest.

Test Post 1	7991262065.31416
Test Post 2	3995630872.8398724

Figure 6.3: Demonstration of scoring posts by non-common interests

Figure 6.3 shows that Post 2 ranks much lower than Post 1 since one of the Post 1's interests are irrelevant to the user.

Lastly, in order to test the final component of the algorithm, which are penalties, I will set the same interests for both Posts but include one of the penalised words in the contents of Post 2 and leave Post 1 as is.

Test Post 1	7991329228.69133
Test Post 2	1598266231.112102

Figure 6.4: Demonstration of scoring posts by the use of penalised words

It is clear from figure 6.4 that the penalty deductions has indeed taken place and Post 2 is now scored lower than Post 1.

Chapter 7

Statement of Ethics

During the creation of this project, all the legal, ethical and social aspects have been considered and analyzed below. It is worth saying that this project is completely legal and created fully in "good faith".

7.1 Leagal standpoint

In the process of creating this project I was extra careful that the information and software I used were under open-source or similar licenses. This is important since many technologies are licensed so that they can not be used in any shape or form. Furthermore, most of the software I used had technologies that handled user information securely by hashing or other ways. Specifically this what Django's authentication system consists of:

1. Users
2. Permissions: Binary (yes/no) flags designating whether a user may perform a certain task.
3. Groups: A generic way of applying labels and permissions to more than one user.
4. A configurable password hashing system
5. Forms and view tools for logging in users, or restricting content
6. A pluggable backend system

7.2 Ethical standpoint

The ethical issues that might arise in Web applications nowadays is not clearly yet defined. The goal of this project is entirely with good intent and is in the public interest. In order for a website to be ethical, the user information should be private, and the user should always know how their information is collected and if so, how is it used (Harris 2009).

7.2.1 ACM-IEEE

This project complies with the four ethical design guidelines as suggested by the ACM-IEEE Joint Software Engineering Code of Ethics (JSECOE), which are:

1. Information should not be hidden
2. Information should not be used or transferred without proper consent
3. Consent should be properly obtained
4. Privacy should be maintained

First and foremost, user information is at all times visible and accessible by them, and only them. Secondly, user information is never transferred. Thirdly, when the user selects a username, email, password or interests, their consent is always asked, whether that would be by the register button or by saving changes in their account profile once they have an account. In fact, when a user creates a post they are warned that their email will be visible underneath the post, so they can decide whether they would like to proceed or not. Where this project to be designed from a more scalable perspective and commercial use, perhaps there would be other issues to consider such as gathering user tendencies which can be used for advertising purposes.

7.2.2 The Computer Misuse Act 1990

The Parliament of the United Kingdom passed three new offences which relate to Computer users. Performing any of the act listed below is considered a crime.

1. Accessing computer material without permission, eg looking at someone else's files.

2. accessing computer material without permission with intent to commit further criminal offences, eg hacking into the bank's computer and wanting to increase the amount in your account.
3. Altering computer data without permission, eg writing a virus to destroy someone else's data, or actually changing the money in an account.

This project also complies with the Computer Misuse Act since all exploits except hacking into the host server have been eliminated. The reason this exploit remains is because no matter how secure the software itself is, if this project was hosted on a non-secure server, some of the website information could potentially be obtained by un-authorized access.

7.2.3 Copyright

This website is copyright free. What the user's decide to share is simply their choice. Moreover, the website does not condone any kind of content copying and is in no way associated with file sharing schemes.

7.2.4 Conclusion

Lastly it is worth mentioning that since the website does not provide many complicated aspects, following an ethical procedure when developing such a website is not difficult. However, in websites with multiple functionalities, actions and a plethora of users, there are a lot of ethical issues that arise and have to be dealt with in "good interest". Even when ethical issues arise, it is usually due to the lack of attention from the developer's side.

Chapter 8

Conclusion

In this section, I will look into how the overall project progressed along with:

1. What this project brings to an oversaturated market of social news aggregators.
2. What parts of the project could be improved in order to scale and provide a more sophisticated solution.
3. Things learned throughout the project.

8.1 Overview

Social news and social media have redefined how people interact with the world, and over the last few years have dominated both the web and mobile market. Applications like Reddit, Pinterest and a whole lot more have become a place where people interact and socialise. Most of these websites operate in their niche category and require the user to adhere to guidelines such as only posting photos, or following content specific subreddits (in the case of Reddit). These websites rely on a select few moderators or administrators that dictate what the content of each category will be. This project provides a solution where all the content is "out there" and its up to the users to decide which type of information they see with not much interaction. This is achieved thanks to the website taking care of sorting the posts and provide content regulation, while the users are able to customise the algorithm.

8.2 Potential Improvements

Although the project achieves the goals and objectives, it set out to achieve. It is not to say, however, that there could not be any improvements. As mentioned before, the implementation of a NN that monitors the interactions of the user with the website components could prove beneficial since people's interests tend to change over time. Nevertheless, the use of an AI would require much computational power and even more testing in order to train the model and provide accurate predictions on the user's interests.

8.3 Personal Gains

While developing this project, I faced many challenges, both in the development process and real-world issues. This project allowed me to venture more into my passion for Web Design and provide a solution to a problem. Furthermore, I was given the opportunity with this project to explore new languages and frameworks that gave me a broader view of what I would like to pursue professionally. The most important lesson I learned during the development, is how the correct and consistent use of the Software Development Lifecycle (SDLC) and even more is the use of the agile methodology leads to the development of high-quality software. All in all, this project was an eye-opening experience in terms of the acquirement of new skills but also, in testing my capabilities and what I can achieve when developing enjoyable software.

Bibliography

Bell, R. (2009), ‘A beginner’s guide to big o notation - rob bell’.

URL: <https://rob-bell.net/2009/06/a-beginners-guide-to-big-o-notation/>

Billard, M. (2019), ‘An introduction to ranking algorithms seen on social news aggregators (example)’.

URL: <https://coderwall.com/p/cacyhw/an-introduction-to-ranking-algorithms-seen-on-social>

Contributors, W. (2019), ‘Bootstrap (front-end framework)’.

URL: [https://en.wikipedia.org/wiki/Bootstrap_\(front-end_framework\)](https://en.wikipedia.org/wiki/Bootstrap_(front-end_framework))

DoorDash (2017), ‘Tips for building high-quality django apps at scale’.

URL: <https://medium.com/@DoorDash/tips-for-building-high-quality-django-apps-at-scale-a5>

Editorial-Team (2016), ‘Why we use django framework and what is django best used for’.

URL: <https://djangostars.com/blog/why-we-use-django-framework/>

Gabriel, R. J. (2020), ‘Robertjgabriel/google-profanity-words’.

URL: <https://github.com/RobertJGabriel/Google-profanity-words>

Harris, J. (2009), ‘Ethical issues in web design’.

Laramee, R. S. (2011), ‘Bob’s project guidelines: Writing a dissertation for a bsc. in computer science’, *Innovation in Teaching and Learning in Information and Computer Sciences* **10**, 43–54.

Majewski, M. (2019), ‘Top 6 software development methodologies - blog | planview’.

URL: <http://leankit.com/blog/2019/03/top-6-software-development-methodologies/>

Miller, E. (2009), ‘How not to sort by average rating’.

URL: <https://www.evanmiller.org/how-not-to-sort-by-average-rating.html>

Salihefendic, A. (2015a), ‘How hacker news ranking algorithm works’.

URL: <https://medium.com/hacking-and-gonzo/how-hacker-news-ranking-algorithm-works-1d9b0c>

Salihefendic, A. (2015b), ‘How reddit ranking algorithms work’.

URL: <https://medium.com/hacking-and-gonzo/how-reddit-ranking-algorithms-work-ef111e33d0d>

Shirriff, K. S. (2013), ‘How hacker news ranking really works: scoring, controversy, and penalties’.

URL: <http://www.righto.com/2013/11/how-hacker-news-ranking-really-works.html>

Stackify (2017), ‘What is sdlc? understand the software development life cycle’.

URL: <https://stackify.com/what-is-sdlc/>

Wilhee, H. (2016), ‘Django introduction’.

URL: <https://developer.mozilla.org/en-US/docs/Learn/Server-side/Django/Introduction>