

UNIVERSITY OF SOUTHAMPTON

A Context-Sensitive Relevance-Based Intelligent Data-Ranking Agent

by

Thomas J. Bell

A project report submitted for the award of
MEng Electronic Engineering

in the
School of Electronics and Computer Science

November 2013

UNIVERSITY OF SOUTHAMPTON

Abstract

School of Electronics and Computer Science

MEng Electronic Engineering

by Thomas J. Bell

The Thesis Abstract is written here (and usually kept to just this page). The page is kept centred vertically so can expand into the blank space above the title too...

Acknowledgements

The acknowledgements and the people to thank go here, don't forget to include your project advisor...

Contents

Abstract	i
Acknowledgements	ii
List of Figures	vi
List of Tables	vii
Abbreviations	viii
Physical Constants	ix
Symbols	x
1 Introduction	1
1.1 The Problem	1
1.2 Project Objective	1
1.3 Goals	1
1.4 Unique Features	2
2 Literature Review	3
2.1 Existing Data-Ranking Implementations	3
Google Now	3
ViralHeat	3
StreamLife	4
2.2 Text analytics libraries	4
AlchemyAPI	4
Semantria	4
Saplo	4
Wingify	4
DatumBox	4
2.3 Data-mining	5
2.3.1 Facebook	5
2.3.2 Twitter	5
2.3.3 Google Calendar	6
2.3.4 Google Tasks	6

2.3.5	Android SMS	6
2.4	Semantic Analysis	7
2.4.1	DatumBox API and usage	7
2.5	Sorting Algorithms	8
2.6	Data Sources and Persistence	8
3	Algorithm Design	10
3.1	Scoring Rule	10
3.1.1	Topic-classification scoring formulae	10
3.1.2	Unweighted Scoring Rule	11
3.1.3	Weighted Scoring Rule	11
3.2	Lexical Analysis and Integration	13
3.3	Ranking Algorithm	13
4	Design Detail and Specification	14
4.1	Data Sources and Acquisition	14
4.1.1	Data Acquisition	14
4.1.2	Data Item Conversion	14
4.1.3	User Profile Acquisition	14
4.2	Scoring Algorithm	14
4.3	Ranking Algorithm	15
4.4	Android Demonstration	15
5	Implementation	16
5.1	API Development	16
5.2	Android App Development	16
6	Planning and Progress	17
6.1	Planning	17
6.1.1	Skills Audit	17
6.2	Time Schedule	17
6.3	Progress	17
6.4	Risk Analysis and Contingency Planning	17
7	Testing Strategy and Results	18
7.1	A Section	18
7.1.1	A Subsection	18
7.2	Another Section	18
8	Critical Evaluation	19
8.1	A Section	19
8.1.1	A Subsection	19
8.2	Another Section	19
9	Conclusion	20
9.1	A Section	20
9.1.1	A Subsection	20
9.2	Another Section	20

10 Further Work	21
10.1 A Section	21
10.1.1 A Subsection	21
10.2 Another Section	21
 A An Appendix	 22
 Bibliography	 23

List of Figures

List of Tables

Abbreviations

LAH List Abbreviations **Here**

Physical Constants

Speed of Light $c = 2.997\,924\,58 \times 10^8 \text{ ms}^{-\text{s}}$ (exact)

Symbols

a	distance	m
P	power	W (Js^{-1})
ω	angular frequency	rads^{-1}

For/Dedicated to/To my...

Chapter 1

Introduction

Context-sensitive techniques, data fusion and ranking agents are frequently used in industry applications and consumer software products. Combined, these techniques are useful for ranking social media according to a users preferences and context. This report describes an agent for the classifying, scoring and ranking of data according to its context-sensitive relevance to a user.

1.1 The Problem

Social media, productivity tools and internet-based information are abundant on mobile devices, leading to users being overwhelmed with information, despite only a small amount of it being of any interest to a particular individual at any given moment. This calls for a means by which such data can be ranked or filtered according to its importance, interest or relevance.

1.2 Project Objective

The objective of this project is to produce a scalable and highly modular context-sensitive mobile-content relevance-based intelligent ranking agent, to order social media, productivity and other web-based information according to its time- and situation-specific relevance to a user.

1.3 Goals

The following are core goals which this project sets out to achieve.

1. Develop a scoring algorithm by which to judge the relevance of an item of data based upon a users
 - (a) Personality profile
 - (b) Historical data such as click-history
 - (c) Environment (conditional upon time-constraints)
2. To perform automatic remote topic analysis to judge the topic of an item of data
3. Develop a sorting/ranking algorithm to sort or insert scored items of data efficiently, into an ordered list
4. Develop a stable and robust data fusion technique to combine a range of data into a user-context and data-context object.
5. Abstract away this agent into an extensible Java API for use in
 - (a) Smartphone apps (Android)
 - (b) Web-apps (Spring MVC)
 - (c) Desktop applications (Java Swing etc.)
6. To develop a consumer smart phone app to demonstrate the working API which automatically ranks a users data according to its relevance

These are the criteria by which the extent of this project's success will be evaluated.

1.4 Unique Features

Many of the aspects of this project have never been seen combined into a single research project before and others (such as query-less context-sensitive scoring of data) have recieved little attention. This project is unique in its endeavour to combine data fusion techniques with context-sensitive scoring in the development of a commercially viable prototype.

Chapter 2

Literature Review

A progressive project in the sphere of cross-platform relevance-based intelligent ranking agents, using text analysis and a mathematically rigorous scoring algorithms, requires research in a range of areas across the entire spectrum of low- to high-level computational theory and existing product research. The following summarises the research undertaken before and during the research and design phases.

2.1 Existing Data-Ranking Implementations

A good number of content aggregators exist at present in various forms, yet all distinctly lack the complementary relationship between social media (and others) and context-sensitive text-analytics for relevance-based ranking. The following is a summary to the key players in this space at present.

Google Now Google Now is a mobile app which combines Google's search feature with useful information which is deemed relevant to the user's environment such as weather, a map to get them home after a night out or nearby events.

ViralHeat ViralHeat is a web-based social media content aggregator and filter, used for commercial uses of social media. It allows the user to filter content from twitter, Facebook and others according to its sentiment (positive/negative). It's not available as a non-commercial social media aggregator and does not perform topic analysis for ranking.

StreamLife This app aggregates Facebook content and tweets, but performs no topic/sentiment analysis or ranking, and provides no capability for including tasks, calendar appointments, SMS messages or emails.

2.2 Text analytics libraries

There are a good number of existing text analytics services available to the end user and the developer for a range of different types of analytics. These services include sentiment analysis, text categorisation, contextual targeting and a range of others. This section highlights some which are relevant to this study.

AlchemyAPI Alchemy provides an API which performs entity extraction, sentiment analysis, concept tagging, relation extraction and most notably, text categorisation (among others). It provides a free licence for up to 1,000 API calls per day. It provides all the core features which this project requires in terms of text analytics, however the text categorisation did not produce strong results for short amounts of text (such as tweets) and often failed to make any categorisation.

Semantria The best solution for sentiment analysis appears to be semantria. It gave consistently precise and accurate sentiment analysis for text containing more than 5 words.

Saplo Saplo was distinguished in its contextual analysis feature which allowed me to define a personalised textual context which could be matched against any type of text. This could have allowed me to define user-specific textual contexts against which to match data items, however it only allows users 2000 API calls per month on their free account.

Wingify This is a beta-stage contextual targeting API which can categorise text from a web page and extract key concepts. The online demo provided accurate results, yet there is not yet a public API available.

DatumBox DatumBox is a free machine learning API which performs sentiment analysis, subjectivity analysis, topic classification, language detection, readability detection, educational detection, document similarity analysis, and gender detection. Many

of these features may be useful for ranking text based upon its relevance to a particular individual. It has a simple API using http POST requests and a JSON response.

2.3 Data-mining

Outline: Facebook and Twitter API (phone, web and desktop), Android API, Android Calendar, Android Tasks, Android SMS, Android Sensors, Google Calendar and Tasks (web-based and desktop API).

2.3.1 Facebook

The Facebook statuses of a users friends can be fetched either using the Facebook API, or through wrapper libraries which simplify common Facebook API usage. The Facebook API is full-featured and well documented, yet for the purposes of demonstrating this ranking agent it's not necessary. Facebook4J is a Java Facebook wrapper API which simplifies the most common Facebook API features into a more minimal library.

```
Facebook facebook = new FacebookFactory().getInstance();
facebook.setOAuthAppId(appId, appSecret);
facebook.setOAuthPermissions(commaSeparatedPermissions);
facebook.setOAuthAccessToken(new AccessToken(accessToken, null));
facebook.postStatusMessage("This is my status.");
//Gets a list of the users feed (friend's status updates)
ResponseList<Post> feed = facebook.getHome();
```

LISTING 2.1: Facebook4J example [1]

This library provides the capability to public messages and links, getting the users news feed, 'liking' a post, publishing a comment, searching for users, groups, events, places or locations and others. It supports pagination and reading options. Altogether it fulfills the needs of this project adequately.

2.3.2 Twitter

Similar to Facebook, Twitter provides an API which is freely available yet overly complex for this project. Twitter4J is a free API which simplifies the Twitter API.

```
Twitter twitter = TwitterFactory.getSingleton();
//Gets and prints the users timeline
List<Status> statuses = twitter.getHomeTimeline();
for (Status status : statuses) {
    System.out.println(status.getUser().getName() + ":" +
        status.getText());
}
```

```
}
```

LISTING 2.2: Twitter4J example [2]

Twitter4J provides sufficient documentation, support and features making it suitable for retrieving a users Twitter feed.

2.3.3 Google Calendar

Calendars from a users Google account can be retrieved on the Android platform using the Calendar Provider. This is a repository of a user's calendar events which can be queried.

```
Cursor cursor = context.getContentResolver()  
    .query(  
        Uri.parse("content://com.android.calendar/events"),  
        new String[] { "calendar_id", "title", "description",  
            "dtstart", "dtend", "eventLocation" }, null,  
        null, null);
```

LISTING 2.3: Calendar Provider example

This query returns a list of events which can be freely processed and sorted as required.

2.3.4 Google Tasks

In Android, Tasks can be retrieved from a Google account using the Google Tasks API by prompting the user for their account credentials, retrieving an AuthenticationToken to create a GoogleCredential and using the Tasks Builder to create a Tasks Service. This is demonstrated in Listing 2.4.

```
List<Task> tasks = service.tasks().list("@default").execute().getItems();
```

LISTING 2.4: Google Tasks example

2.3.5 Android SMS

SMS messages can be received in Android applications using a BroadcastReceiver which collects incoming SMS messages.

```
public class receiver extends BroadcastReceiver {  
    public String str = "";  
    @Override  
    public void onReceive(Context context, Intent intent) {  
        Bundle bundle = intent.getExtras();
```

```
SmsMessage[] msgs = null;
if (bundle != null) {
    Object[] pdus = (Object[]) bundle.get("pdus");
    msgs = new SmsMessage[pdus.length];
    for (int i = 0; i < msgs.length; i++)
    {
        msgs[i] = SmsMessage.createFromPdu((byte[]) pdus[i]);
    }
}
}
```

LISTING 2.5: Android SMS example

This example shows a `BroadcastReceiver` which collects SMS messages when they're received and adds them to an array of `SmsMessage` objects for processing.

2.4 Semantic Analysis

This project requires topic analysis of items of mobile data, in order to compare them to the user's preference and ascribe relevance to them. Other semantic analysis capabilities would prove beneficial for increasing the range of criteria by which relevance may be judged and to eliminate irrelevant data as early on as possible. These include readability, gender, subjectivity and language detection. The DatumBox API is chosen for semantic analysis for its coverage of these requirements; its ease of implementation and the fact that its use is free.

2.4.1 DatumBox API and usage

Each feature provided by DatumBox has a POST Request URL and is retrieved in code by setting up the request headers and URL parameters (including the API key and text), and waiting for the asynchronous response as a JSON object 2.6.

```
//Request URL
http://api.datumbox.com:80/1.0/TopicClassification.json

//Request headers
{"Content-Type":"application/json; charset=UTF-8"}

//Response body
{
    "output": {
        "status": 1,
        "result": "Computers & Technology"
    }
}
```

LISTING 2.6: DatumBox Topic Classification example

This example is for 'Topic Classification', but the process for requesting other available features is similar and will be simplified by creating a 'Semantic Analysis API' which retrieves the response for each of the different requests.

2.5 Sorting Algorithms

Sorting is required for ordering scored items of data into a list whereby the topmost items are the most relevant and the bottommost are the least relevant. Since we are only dealing with relatively small sets of data, efficiency is not paramount in terms of maximising accuracy or usability. Despite this, some consideration is made to using the most appropriate sorting algorithms. Merge sort which uses a divide-and-conquer approach is one of the most suitable for initial scoring of larger numbers of items. It is stable and due to its constant performance of $O(n \log n)$, predictable in terms of execution time. For nearly sorted lists such as scenarios where single or small numbers of data items (≤ 10) may be added to an ordered list, insertion sort is superior in speed and simplicity to others and will have a performance of $O(n)$.

2.6 Data Sources and Persistence

This project attempts to develop a framework for use in a wide range of software applications and as such will require complete flexibility in terms of its data sources. For the purposes of research, development and initial testing, JSON (JavaScript Object Notation) objects in text files are considered the best solution for long-term storage of test data. They are simple to create, read, update and delete; flexible in terms of multi-device usability and fully interchangeable with Java objects. Thus they are suitable for storing items of test data to be loaded into the API. UserContext objects may also be stored as such. Listing 2.7 shows an example of a JSON DataItem representation.

```
//Example JSON DataItem
{
    "id":0,
    "title":null,
    "detail":"",
    "type":"Tweet",
    "score":0,
    "link":null,
    "date":null,
```

```
    "author": "test",  
    "location": null,  
    "from": null,  
    "to": null,  
    "hashtag": "#test",  
    "priority": 0  
}
```

LISTING 2.7: JSON test data example

Android provides a range of ways to save persistent application data. Data storage options include: shared preferences, internal storage, external storage, SQLite databases or on an external server accessed remotely [4]. Any number of these could be used to store information about the user, such as their preferences and `UserContext` which is required by the ranking agent. Internal storage will store text files in the file system. They are private to the user and other applications and can be used to store JSON objects. A database may be used to store test data and user-related data, but this would require classes for managing a database connection and parsing the data fields into their original objects. Either internal or database storage would be perfectly adequate.

Chapter 3

Algorithm Design

This is some sample text.

3.1 Scoring Rule

In order to rank data according to its relevance, it must first be scored according to its relevance. Scoring may be considered a form of classification, in which an item of data is classified as belonging to a particular cluster within a data set, represented as a set of data points in a feature vector. Scoring may be done using supervised machine-learning approaches or using an unsupervised scoring function.

The word 'importance' is used to describe the non-user-specific likely relevance of an item of data. 'Relevance' signifies a user-specific importance of an item of data.

3.1.1 Topic-classification scoring formulae

Scoring functions typically combine key-term statistics into a single score as a measure of the similarity between a query and a document. Term frequency is the most frequently used and widely explored approach to understanding the relevance of a text. Term rank scoring formulae are comprised of a base formula (typically Okapi BM25, Eqn. 3.1) and a multiplicative or additive range variance term R [3].

$$\sum_{t \in d \cap q} \ln \frac{N - df + 0.5}{df + 0.5} \cdot \frac{tf}{0.5 + 1.5 \cdot \frac{dl}{avdl} + tf} \cdot R \quad (3.1)$$

Here df is the document frequency (documents in a collection) and tf is simply the term frequency (to occurrences of a term in a particular document). $avdl$ is the average

document length and R is a component which limits the range over which the term rank has an effect.

Topic classification is done remotely using a similar technique using multiple term frequency using weights which describe how closely each term relates to the specified topic. This classification provides us with information about our item of data which can be used to rank it according to its relevance.

3.1.2 Unweighted Scoring Rule

Having performed topic-classification on an item of data, a scoring algorithm must assign to it a value which is a measure of its relevance. Let us first consider the case of an unweighted rule which scores an item of data.

A classified item of data (that is, one which has been assigned a score for each attribute according to how well that attribute describes the topic of the text), without considering the users preferences, has an important topic if the sum of the scores of each attribute is large. This is related to the average score of an attribute being large. Consequently, the importance of an item of data may be given as the average of the attribute scores

$$f(D) = \frac{\sum_{i=1}^i d_i}{n} \quad (3.2)$$

where D is our item of data, and the elements d_i are its topic-attributes.

3.1.3 Weighted Scoring Rule

Let us now move on to consider the user-specific case, or weighted case. Relevance is defined here as the extent to which the score of each attribute of an item of data, matches the preference score of the user for that attribute. Attribute-specific scores are found by comparing the complementary user and data-item attributes according to a given rule. A final scoring rule must be used to combine scores of individual attributes into an overall score for a particular item of data, to given its relevance. Fagin and Wimmers [?] detail a formula for incorporating weights into scoring rule, no matter the nature of the scoring rule. It takes a function f which describes an unweighted rule for applying a score to a tuple of n entries (attributes in our case) and gives a weighted rule based on f which is compatible with any rule.

They take a set $\Theta = (\theta_1, \theta_2, \dots, \theta_n)$ of weights which relate the effect of an entry x on the score of the tuple $X = (x_1, x_2, \dots, x_n)$. If $f(X)$ is the unweighted scoring rule, then for a tuple of m entries

$$f_{\Theta}(X) = \left(\sum_{i=1}^{m-1} i \cdot (\theta_{\sigma(i)} - \theta_{\sigma(i+1)}) \cdot f(X \upharpoonright \{\sigma(1), \dots, \sigma(i)\}) \right) + m \cdot \theta_{\sigma(m)} \cdot f(X) \quad (3.3)$$

Here $X \upharpoonright \{\sigma(1), \dots, \sigma(i)\}$ is a restriction of X to the domain of a bijection σ which orders the weightings to match the order of entries in the tuple. In our case this bijection would be a mapping from each attribute of an item of data to the weighting associated with that attribute. The weighting Θ would be dependent upon a complementary tuple describing the user, i.e. their attribute preferences. The term $(\theta_{\sigma(i)} - \theta_{\sigma(i+1)})$ is the difference between the weightings of two consecutive entries.

Fagin and Wimmers' equation allows me to include weighed relationships between attributes' sub-scores and the total score of the item of data. An item of data is represented as a tuple $D = (d_1, d_2, \dots, d_n)$ where each entry d_n is an attribute of the item of data. Given the context-sensitive nature of this project, attributes must affect the score of the tuple with differing degrees according to how well they match the preferences of the user.

This idea of weighting will be an analogy for the users preferences, such that the greater the user's preference of an attribute, the greater the weighting will be. Thus the set of weights Θ becomes a tuple of entries $U = (u_1, u_2, \dots, u_n)$ whereby U represents the tuple of attributes depicting the user's preferences. This leads us the following result, showing a general solution for context-sensitive (weighted) scoring rule, for any unweighted rule $f(D)$.

$$f_U(D) = \left(\sum_{i=1}^{m-1} i \cdot (u_{\sigma(i)} - u_{\sigma(i+1)}) \cdot f(D \upharpoonright \{\sigma(1), \dots, \sigma(i)\}) \right) + m \cdot u_{\sigma(m)} \cdot f(D) \quad (3.4)$$

Let us take the simple unweighted scoring rule whereby we compute the average of the the i attributes d_j in a tuple D (Eqn. 3.2). Using the weighting formula (Eqn. 3.4) the weighted equivalent is

$$f_U(D) = \left(\sum_{n=1}^{M-1} n \cdot (u_{\sigma(n)} - u_{\sigma(n+1)}) \cdot \frac{\sum_{i=1}^n d_i}{n} \right) + M \cdot u_{\sigma(M)} \cdot \frac{\sum_{i=1}^M d_i}{M} \quad (3.5)$$

$$= (u_1 - u_2)d_1 + 2(u_2 - u_3)\frac{d_1 + d_2}{2} + 3u_3\frac{d_1 + d_2 + d_3}{3} \quad (3.6)$$

$$= u_1d_1 + u_2d_2 + \dots + u_Md_M \quad (3.7)$$

We are lead therefore to the proof of an intuitive linear weighted scoring rule (Eqn. 3.8).

$$f_U(D) = \sum_{k=1}^{m-1} u_k d_k \quad (3.8)$$

3.1.4 Weighted Scoring Rule using Multiple Weights

This project deals with the relevance of an item of data to a user. An item of data may be social media, news or productivity related. The relevance of productivity related information is not usually based upon its topic, for example the relevance of a calendar appointment for a meeting or a dentist appointment cannot be assessed according to the principles used above. Its topic-classification is not useful in commenting on its relevance. This situation calls for a weighting which is based not upon topic relevance, but upon time relevance. It must complement topic-relevance so as not to affect the ranking order of non-time-specific items of data.

This weighting may be implemented using

1. a discrete binary decision (relevant at this point in time or not), or
2. a continuous function of increasing time-relevance as a deadline approaches

The simplest of the two is the former, whereby an item of data (with a deadline in less than *t_{threshold}* minutes) would be flagged as time-relevant and ranked at the top of the ordered list.

The second implementation would be the preferable choice, since as the deadline approached, the item of data would creep up the ordered list until it reached the top at *t_{threshold}* minutes before the deadline.

3.2 Lexical Analysis and Integration

education/gender/readability etc.

3.3 Ranking Algorithm

i.e. Modified insertion sort

Chapter 4

Design Detail and Specification

Describe the overall specification - what it should do.

4.1 Data Sources and Acquisition

This is some sample text.

4.1.1 Data Acquisition

This is some sample text.

4.1.2 Data Item Conversion

This is some sample text.

4.1.3 User Profile Acquisition

This is some sample text.

4.2 Scoring Algorithm

This is some sample text.

4.3 Ranking Algorithm

This is some sample text.

4.4 Android Demonstration

This is some sample text.

Chapter 5

Implementation

This is some sample text.

5.1 API Development

This is some sample text.

5.2 Android App Development

This is some sample text.

Chapter 6

Planning and Progress

Comment of approach

6.1 Planning

This is some sample text.

6.1.1 Skills Audit

This is some sample text.

6.2 Time Schedule

Inset Gantt chart

6.3 Progress

This is some sample text.

6.4 Risk Analysis and Contingency Planning

This is some sample text.

Chapter 7

Testing Strategy and Results

This is some sample text.

7.1 A Section

This is some sample text.

7.1.1 A Subsection

This is some sample text.

7.2 Another Section

This is some sample text.

Chapter 8

Critical Evaluation

This is some sample text.

8.1 A Section

This is some sample text.

8.1.1 A Subsection

This is some sample text.

8.2 Another Section

This is some sample text.

Chapter 9

Conclusion

This is some sample text.

9.1 A Section

This is some sample text.

9.1.1 A Subsection

This is some sample text.

9.2 Another Section

This is some sample text.

Chapter 10

Further Work

This is some sample text.

10.1 A Section

This is some sample text.

10.1.1 A Subsection

This is some sample text.

10.2 Another Section

This is some sample text.

Appendix A

An Appendix

This is an appendix.

Bibliography

- [1] Facebook4j website. URL <http://facebook4j.org/en/code-examples.html>.
- [2] Twitter4j website. URL <http://twitter4j.org/en/code-examples.html>.
- [3] S Jones M M Hancock-Beaulieu M Gatford S E Robertson, S Walker. Okapi at trec-3. *NIST*, 1995.
- [4] Wei-Meng Lee. Beginning android 4 application development. pages 263–291, 2012.