

Interaction & Preference Filtering

A novel approach to Social Recommendation

Riley Kidd

A subthesis submitted in partial fulfillment of the degree of
Bachelor of Software Engineering at
The Department of Computer Science
Australian National University

October 2012

© Riley Kidd

Typeset in Palatino by \TeX and $\text{\LaTeX} 2_{\epsilon}$.

Except where otherwise indicated, this thesis is my own original work.

Riley Kidd
13 October 2012

Abstract

Social networks provide a wide array of user specific interactions, profile information and user preferences. This thesis attempts to decipher which user interactions or preferences are truly indicative of 'likes', this information is then leveraged to allow for binary classification of user specific links with the goal of discovering the ideal combination of information for prediction.

The success of our predictions are evaluated using a number of machine learning algorithms including, *Naive Bayes* (NB), *Logistic Regression* LR and *Support Vector Machines* (SVM), results are compared to previous work using *Matchboxing* (MB) and *Social Matchboxing* (SMB) techniques as baselines. The data set is sourced from a set of over 100 Facebook users and their interactions with over 30,000 friends during a four month period.

Our analysis has shown that user interactions in themselves are not predictive of user likes, while user preferences are. These results, coupled with user like exposure curves offer a useful paradigm for extracting and exploiting user preferences for prediction across our data set.

Contents

Abstract	v
1 Introduction	1
1.1 Objectives	1
1.2 Contributions	2
1.3 Outline	2
2 Background	3
2.1 Facebook	3
2.2 LinkR	5
2.3 Notation	6
2.4 Feature Vectors	7
2.5 Previous Work	7
2.5.1 Social Collaborative Filtering	7
2.6 Training and Testing	8
2.7 Classification Algorithms	8
2.7.1 Constant	8
2.7.2 Social Match Box	8
2.7.3 Naive Bayes	8
2.7.4 Logistic Regression	9
2.7.5 Support Vector Machine	9
2.8 Evaluation Metrics	9
3 Interactions	11
3.1 User Interactions	11
3.2 Conversation	14
3.2.1 Outgoing	16
3.2.2 Incoming	18
3.3 Conclusion	21
4 User Preferences	23
4.1 Demographics	23
4.2 Traits	27
4.3 Groups	37
4.4 Pages	42
4.5 Conclusion	47

5	Feature Combinations	49
5.1	Positive Feature Combination	49
5.2	Summary	52
5.3	Future Work	53
	Bibliography	55

Introduction

An individual's social presence on the web is continually expanding, with the emergence of services such as Facebook, Myspace, LinkedIn, Twitter and Google+ what defines a user and their online social interactions (messages, posting, commenting, etc) and preferences (demographics, group memberships, likes, etc) is an ever expanding graph structure of verbose social content. The Internet is becoming a network of people, providing a myriad of expanding social information and user driven content.

The ultimate question we wish to address in this thesis then becomes: How can we exploit this information to decipher which user interactions or preferences are most indicative of user likes?

We address this question by comparing and contrasting different feature vectors in our data against appropriate baselines and ultimately offer a feature vector combination paradigm which represents improved results which are computationally faster and can more appropriately generalise over a population than previous methods.

1.1 Objectives

The primary objective of this thesis is to contrast and compare differing user feature vectors across user interactions and preferences. Using the machine learning concepts of *Naive Bayes* (NB), *Logistic Regression* (LR) and *Support Vector Machines* (SVM) compared with our appropriate baselines of *Social Matchbox* (SMB) and *Constant Classifiers*. With the goal of discovering which features are most predictive of user likes.

Based on the insight that social influence can play a crucial role in a range of behavioral phenomena [Granovetter 1978; Watts and Strogatz 1998] we will also test using an exposure curve [Romero et al. 2011] hold out technique, where data is only tested if some friend has already liked that datum. [Pantel and Haas 2012] has been shown that positive social annotations on search items adds perceived utility to the worth of a result, implying that a previously liked item will be more predictive.

Finally, we will analyse the effect of combining successful user feature vectors together using a *Feature Combination* approach.

1.2 Contributions

Our specific contributions made during this thesis show:

- Both *Interactions* and *Incoming Outgoing messages* are not more predictive than previously used *Social Matchboxing* techniques.
- Each user preference of *Traits*, *Groups* and *Pages* contributed to a better result than previously used techniques.
- Combining user preferences with an exposure curve for user likes results in a substantial improvement from previously used techniques.
- Combination of beneficial user preferences results in the most advantageous feature vector for this data.

Overall, we provide a methodology which improves upon previous work and offers an approach to combine positively contributing aspects of different feature vectors in our data.

1.3 Outline

The remaining chapters in this thesis are organised as follows:

- **Chapter 2:** We first outline appropriate background information for the reader. Including information pertaining to the source of our data set, mathematical notation used throughout this thesis, previous work in this area and our research approach and methodology.
- **Chapter 3:** In this chapter we discuss different feature vectors for user interactions and the results of applying these feature vectors in comparison with our baselines.
- **Chapter 4:** A similar feature vector analysis as above, however the feature vectors we utilise are for user preferences.
- **Chapter 5:** In this chapter we discuss results from combining different positive feature vectors and propose an ideal combination based on our analysis.
- **Chapter 6:** Finally, we draw the work done throughout this thesis to a conclusion and offer avenues for future work in this area.

All chapters combined, this thesis represents a novel approach to exploiting and analysing user interactions and preferences to ascertain which features are most indicative of user likes and present an approach of combining these useful feature components into an effective classification paradigm.

Background

In the following section, we define the source of our data set, notation used throughout this thesis, our choice of classification algorithms and our testing approach and methodology.

2.1 Facebook

Facebook is the largest and most active social media service in the world. Facebook users can create a profile containing personal preferences and information (age, birthday, group preferences, favourite athlete, etc) and have friendships and interactions between other users. The four main interactions between users are posts (posting an element on a friends' wall), tags (being mentioned in a friends post or comment), comments (written data on a post) and likes (clicking a like button if a user likes a post or comment). The three mediums for these interactions are across links (some URL), posts (some Facebook post), photos (some uploaded Facebook photo) and videos (some uploaded Facebook video).



Figure 2.1: Here we see an example of a link posted to a friends wall, which has subsequently been liked by two friends.

One issue present in this Facebook paradigm is discovering whether a user doesn't like an item, a users Facebook feed is comprised of activity between their friends, content, groups, etc given the enormous scope of potential feed items, Facebook will only show feed items for users who you have recently interacted with using their *Edge-Rank* [Sanghvi and Steinberg 2010] algorithm.

While many Facebook users have a friend count which is close to the human real word limit, known as the Dunbar number [Hill and Dunbar 2003], the *Edge-Rank* algorithm ensures user interactions are focused on a much smaller subset of their friends. Additionally, given the rate of posting, these top feed items are only displayed for a short amount of time. This is coupled with the fact that Facebook allows users to explicitly like an item, but not dislike it - hence distinguishing between what a user does and does not like becomes difficult.

Given this issue, NICTA have developed an app which explicitly determines a users preference for an item, by facilitating explicit like and dislike options, which will be discussed in the following section.

2.2 *LinkR*

NICTA developed a Facebook app named *LinkR*¹ which would make recommendations to users and record whether or not the user liked the recommended item.

The app also collected information about users, their interactions and preferences as well as a subset of available information about their friends. The app tracked and stored information for over 100 app users and their 39,000+ friends over a 4-month time period.

The table below summarises the data collected from both app users and their friends.

App Users	Posts	Tags	Comments	Likes
Wall	27,955	5,256	15,121	11,033
Link	3,974	-	5,757	4,279
Photo	4,147	22,633	8,677	5,938
Video	211	2,105	1,687	710
App Users and Friends	Posts	Tags	Comments	Likes
Wall	3,384,740	912,687	2,152,321	1,555,225
Link	514,475	-	693,930	666,631
Photo	1,098,679	8,407,822	2,978,635	1,960,138
Video	56,241	858,054	463,401	308,763

Table 2.1: Data records for interactions between users. Rows are the type of interaction, columns are the medium.

¹The main developer of the *LinkR* Facebook App is Khoi-Nguyen Tran, a PhD student at the Australian National University.

Pertinent user features we will exploit during this thesis include:

- Gender
- Age
- Hometown
- Locale
- Group Memberships
- Page Likes
- Favourite Activities
- Favourite Books
- Favourite Athletes
- Favourite Teams
- Inspirational People
- Interests
- Favourite Movies
- Favourite Music
- Favourite Sports
- Favourite Television Shows
- School Information
- Work Information
- Messages data

2.3 Notation

The mathematical notation used by our classifiers during this thesis are outlined below.

- A set of users of size N .
- A set of items of size M .
- A user feature vector X of size i , the size of each feature vector varies based on the current user features being analysed and is explicitly defined in each section. In other words $X = \langle x_1, x_2, \dots, x_i \rangle$

- A set of alters A of size j , the size and composition of each alters set varied based on the current user features being analysed and is explicitly defined in each section.
- An exposure of size k , where each k represents the number of some user n 's friends who have liked some item m .
- A data-set D comprised of $D = \{(n, m, x_i) \rightarrow y\}$ with the binary response $y \in \{0, 1\}$ where 0 represents a dislike and 1 represents a like.

2.4 Feature Vectors

Individual feature vectors for X will be discussed and defined in their appropriate sections under *User Interactions* and *User Preferences*. During this thesis we will analyse the results gained from the 7 different user features listed below, additionally a combination of the most predictive individual feature vectors will also be tested.

- Interactions
- Demographics
- Traits
- Groups
- Pages
- Outgoing Messages
- Incoming Messages

2.5 Previous Work

Below we discuss previous work done in this classification area completed in 2011 by Joseph Noel.

2.5.1 Social Collaborative Filtering

Two general approaches to classification prediction are *content-based filtering* (CBF) [Lang 1995] which exploits item features based on items a user has previously liked, or the second approach which is *collaborative filtering* (CF) [Resnick and Varian 1997] which exploits the current user's preferences as well as those of other users.

Previous work defined the term *social CF* (SCF) [Noel 2011] which augments traditional CF methods with additional social network information, the results of this previous work and analysis came to the conclusion that the approach of *Social Matchbox* provided the best results for this data set. In live user trials SMB provided the best performance against all other implemented algorithms and as such will be used as a baseline in this thesis.

2.6 Training and Testing

All evaluation is done using 10 fold cross validation wherein the data is partitioned into 10 complimentary subsets, each subset is composed of two separate parts, one part is used for training (80%) and the other (20%) is used for testing. All training and testing is performed on each distinct fold and the results are averaged along with their standard error in each table and graph.

2.7 Classification Algorithms

Each classification algorithm used in this thesis is passed the training set for each fold as discussed previously. Based on this data and the current feature vector X the classifier builds a model representation of the data and uses this model to classify each component in the test set into either a like or a dislike.

All feature vector analysis carried out in this thesis will be performed on the following classification algorithms.

2.7.1 Constant

The constant (Constant) predictor returns a constant result irrespective of the feature vectors selected. The most common result in our data set is *False* and hence the *False* predictor is displayed in all analysis, tables and graphs.

2.7.2 Social Match Box

SMB is an extension of existing SCF techniques [Yang et al. 2011; Cui et al. 2011] which constrain the latent space to enforce users who have similar preferences to maintain similar latent representations when they interact heavily.

SMB uses the *Social Regularization* method which incorporates user features to learn similarities between users in the latent space which allows us to incorporate the social information of the Facebook data [Noel 2011]. This objective component constrains users with a high similarity rating to have the same values in the latent feature space, which models the assumption that users who are similar socially should also have similar preferences for items.

2.7.3 Naive Bayes

Naive Bayes is a basic probabilistic classifier which involves applying Bayes' theorem using strong conditional independence assumptions between each feature in X . During training each element i in the feature vector is devised to contribute some evidence that this X belongs to either a like or dislike classification, during testing the class with the highest probability when applied to the model is the classification returned.

The NB implementation used during this thesis is an implementation previously devised by Scott Sanner.²

2.7.4 Logistic Regression

Logistic Regression directly estimates parameters based on the training data assuming a parametric form of the distribution. LR predicts the odds of a feature vector X being either a like or a dislike by converting a dependent variable and one or more continuous independent variable(s) into probability odds.

The LR implementation used during this thesis is *LingPipe* [Alias-i. 2008. LingPipe 4.1.0. <http://alias-i.com/lingpipe> (accessed October 1 2011)].

2.7.5 Support Vector Machine

The *Support Vector Machine* is a supervised learning machine based on a set of basis functions which help construct a separating hyperplane between data points. Training involves building the relevant hyperplanes which can then be used for testing. Each data point is classified as a like or dislike depending on which side of the hyperplane it falls.

The SVM implementation used during this thesis is *SVMLibLinear* [Chang and Lin 2011].

2.8 Evaluation Metrics

When evaluating the success of each feature vector at correctly classifying an item, the following metrics will be analysed.

- A *true positive* (TP) prediction refers to when the classifier correctly identifies the class as true.
- A *false positive* (FP) occurs when the prediction is true, but the true class was false.
- A *false negative* (FN) occurs when the prediction is false but the actual class is true.

Accuracy relates to the closeness to the true value. In the context of our results, the accuracy refers to the number of correct classifications divided by the size of the data set.

$$\text{accuracy} = \frac{\text{number of correct classifications}}{\text{size of the test data set}}$$

²Scott Sanner is a Senior Researcher in the Machine Learning Group at NICTA and the supervisor for this research.

Precision relates to the number of retrieved predictions which are relevant. In the context of our results, the precision refers to the number of TP predictions divided by the sum of the TP and FP predictions.

$$\text{precision} = \frac{\text{number of TP}}{\text{number of TP} + \text{number of FP}}$$

Recall refers to the number of relevant predictions that are retrieved. In the context of our results, recall refers to the number of TP predictions divided by the sum of the TP and FN predictions.

$$\text{recall} = \frac{\text{number of TP}}{\text{number of TP} + \text{number of FN}}$$

The f-score combines and balances both precision and recall and is interpreted as the weighted average of both precision and recall.

$$\text{f-score} = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

The main metric we use for analysis, tabulation and graphing in this thesis is accuracy.

Interactions

The user interactions we examine in this thesis can be broken down into two distinct groups, interactions between users and messages sent between users.

3.1 User Interactions

There are a number of potential interaction mediums between users under the Facebook paradigm. These can be summarised into the following categories.

- *Direction*: The manner an interaction is received, either *Incoming* where a message is posted to some user or *Outgoing* where some user posts a message to another user. Interaction directionality has been shown to be highly reflective of user preferences [Saez-Trumper et al. 2011].
- *Modality*: The medium some user employs to interact with another user via either *Links*, *Posts*, *Photos* or *Videos*
- *Type*: The style some user employs to interact with another user via either *Comment*, *Tag* or *Like*.

In this case, the I for our feature vector X is defined as the cross product of the above components where:

$$I = \{Incoming, Outgoing\} \times \{Posts, Photos, Videos, Links\} \times \{Comments, Tags, Likes\}$$

The alters of I can then be defined as all users who have interacted with or been interacted with the current user via some i . Each component of I is set to 1 if any of the alters defined by the current set i have also liked the current item M , otherwise it is set to 0.

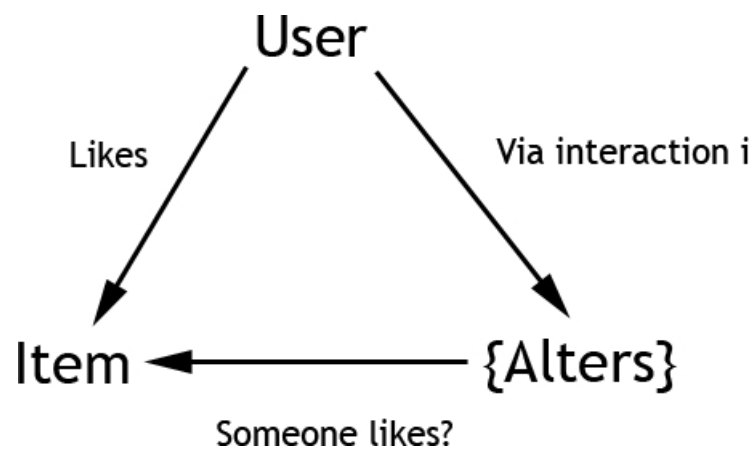


Figure 3.1: *User Interactions* alters paradigm. The alters set is defined by the interaction via i and the feature vector x_i is 1 if some user in the alters set has also liked this item.

Applying these feature and alter sets to our classification algorithms defined above we obtain:

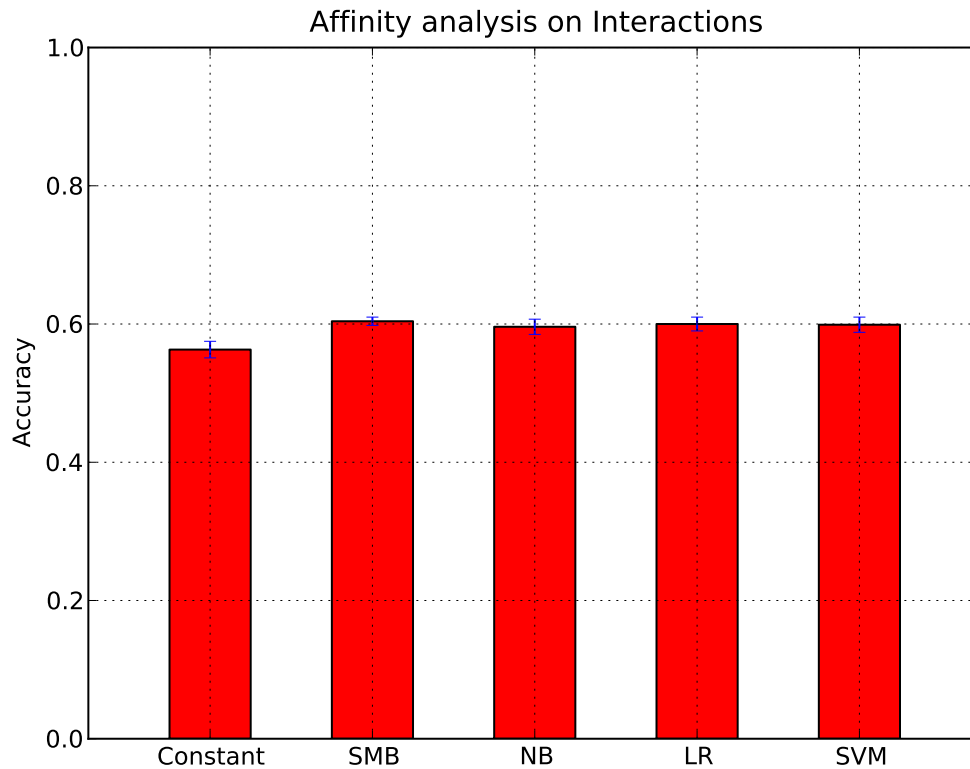


Figure 3.2: Accuracy results using the *User Interactions* feature set

User interactions are marginally outperformed by our SMB baseline, showing that user interactions do not appear to help our classification.

One reason for this result could be we can not track information passing outside of Facebook, users who frequently interact could be real world friends and hence share information via email or word of mouth.

Comparing *User Interactions* against our exposure curve we obtain:

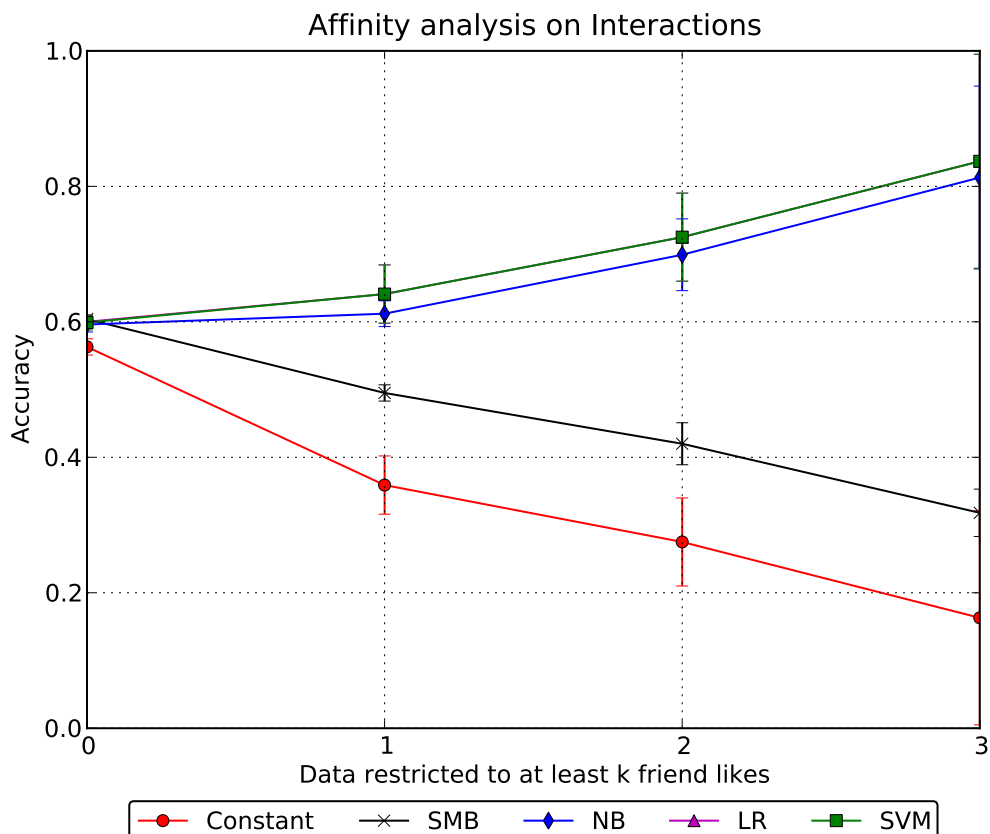


Figure 3.3: Accuracy results for an exposure curve using the *User Interactions* feature set

We glean that as our data is restricted, the performance of our classifiers improves (note LR and SVM obtained the same results in this graph) over time. This graph shows that for *User Interactions* having one user liking an item is enough to improve upon our baseline classifiers.

3.2 Conversation

Given the nature of Facebook, it is possible for users to post or receive messages from other users.

These messages can be broken down based on their directionality, either *Outgoing* which are words sent to other users or *Incoming* which are words received from other users.

Based on our LinkR data set, the most commonly used words occur with a high frequency over our user base and can be seen in the table below:

Rank	Word	Frequency
1	:)	292,733
2	like	198,289
3	good	164,387
4	thanks	159,238
5	one	156,696
6	love	139,939
7	:p	121,904
8	time	106,995
9	think	106,459
10	see	103,690
11	nice	99,672
12	now	94,947
13	well	92,735
14	happy	84,381
15	:d	83,698
16	much	78,719
17	oh	77,321
18	yeah	76,564
19	back	76,032
20	great	70,514

21	going	70,447
22	still	68,245
23	new	67,430
24	day	65,579
25	come	63,837
26	;)	62,936
27	year	61,771
28	look	60,608
29	yes	59,774
30	want	59,514
31	tag	58,633
32	hahaha	57,448
33	also	56,414
34	need	55,921
35	make	54,949
36	sure	54,395
37	thank	54,112
38	people	53,211
39	miss	53,182
40	guys	52,855

Table 3.1: Top conversation content data for all users. We see very common words and online expressions have a high frequency in our data set.

For messages the I of our feature vector X contains an element i for each of the top j most commonly used words based on the conversation content of all users.

The alters of I can then be defined as all users who have liked the current item M . Each component of I is set to 1 if any of the alters have used the current word j where $i = j$ with the user n , otherwise it is set to 0.

3.2.1 Outgoing

The first issue present is to determine the most predictive number of top words j for use by our classifiers. Given the enormous size of potential messages and memory constraints in the testing environment we decided to test within a range of $\{100 - 1000\}$ with an incremental step size of 100 for each test.

The results of testing based on differing sizes of *Outgoing Words* can be seen below:

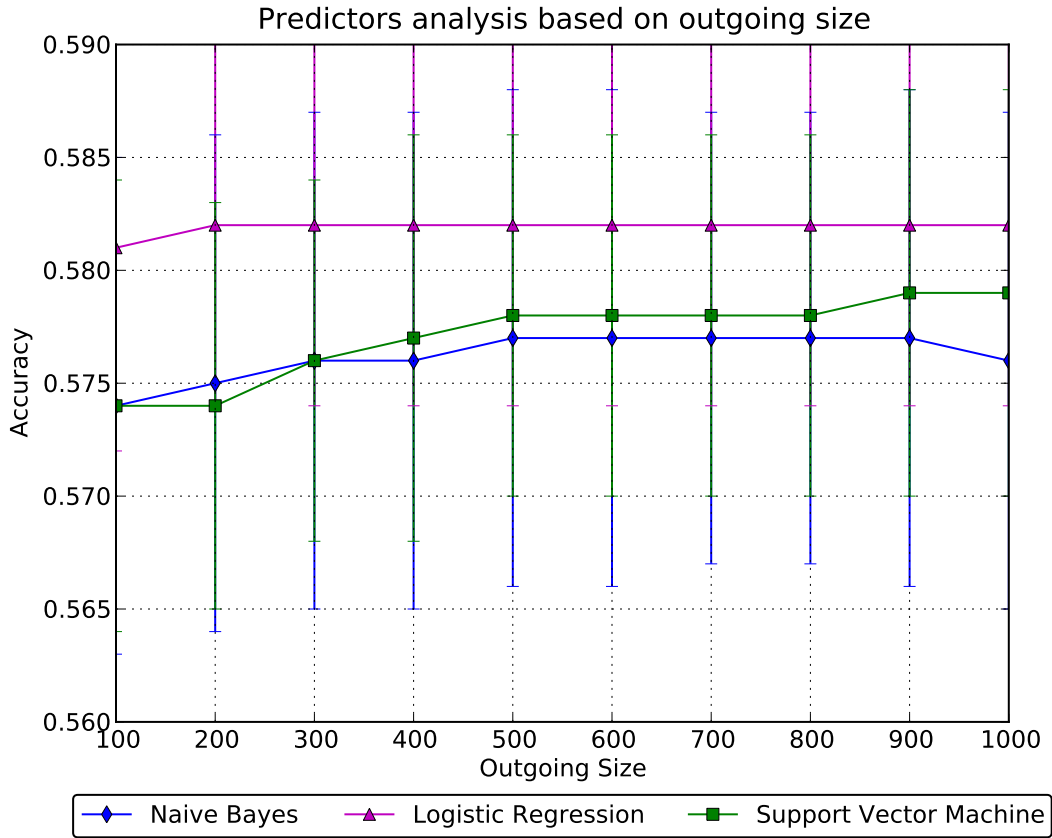


Figure 3.4: Accuracy results for different *Outgoing Words* sizes

The most predictive *Outgoing Words* words sizes j for each of our classifiers are:

- **Naive Bayes:** 500
- **Logistic Regression:** 200

- **Support Vector Machine:** 900

Using the most predictive word sizes j for each of our classifiers as defined above and comparing to our baselines we obtain:

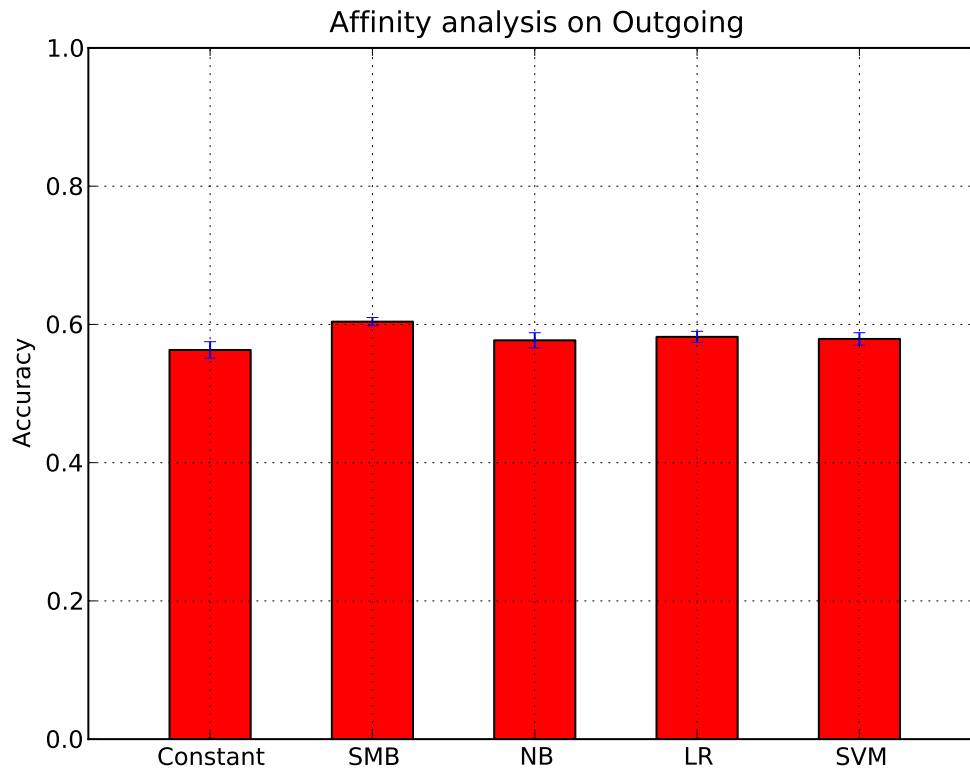


Figure 3.5: Accuracy results using the *Outgoing Words* feature set

These results do not show an improvement over our baselines and in fact are a marginal improvement over the constant baseline. A possible reason for this could be due to the commonality of the words being tested. Highly common and frequently used words would result in poor predictive tendencies, this is eluded to in our graph above which shows an improvement in predictiveness over step sizes for SVM.

Comparing *Outgoing Messages* against our exposure curve we obtain:

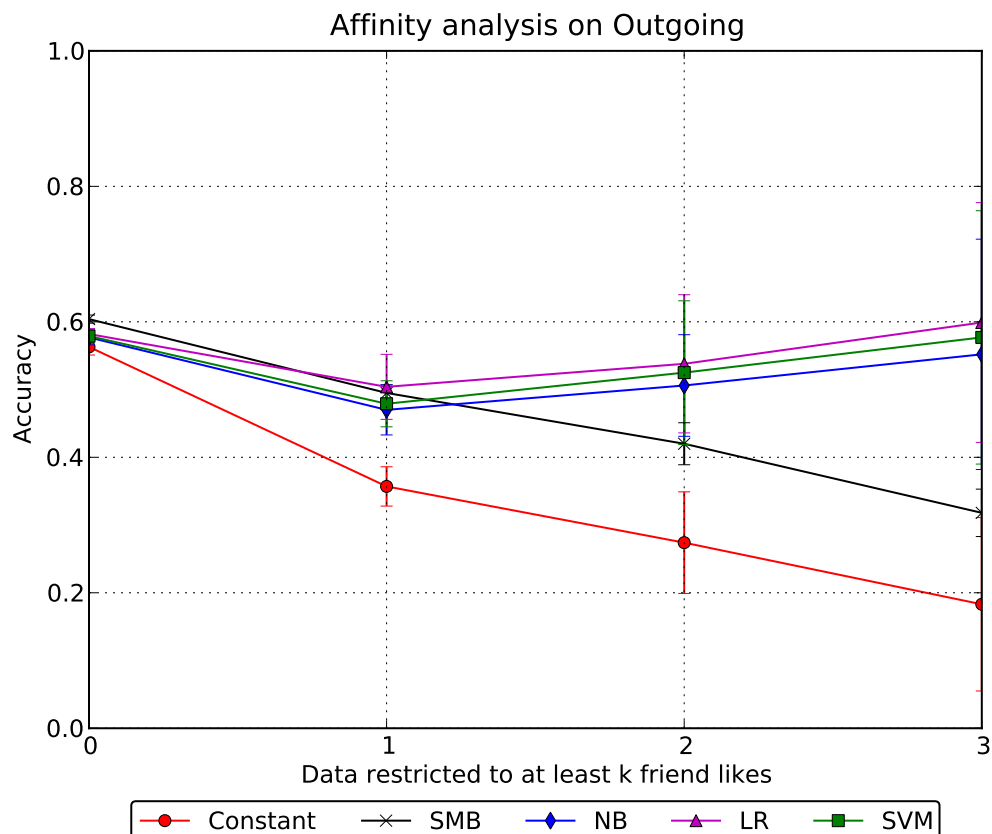


Figure 3.6: Accuracy results for an exposure curve using the *Outgoing Words* feature set

Our exposure curve follows this similar trend of unimproved for $k = 1$ likes, however as k increases there is some improvement from the baselines, but this is negligible in comparison to $k = 0$ for our classifiers.

3.2.2 Incoming

Similarly for *Incoming Words* we need to discover which is the predictive j for use of by our classifiers, using the same methodology as described above for *Outgoing Words* we obtain the following graph:

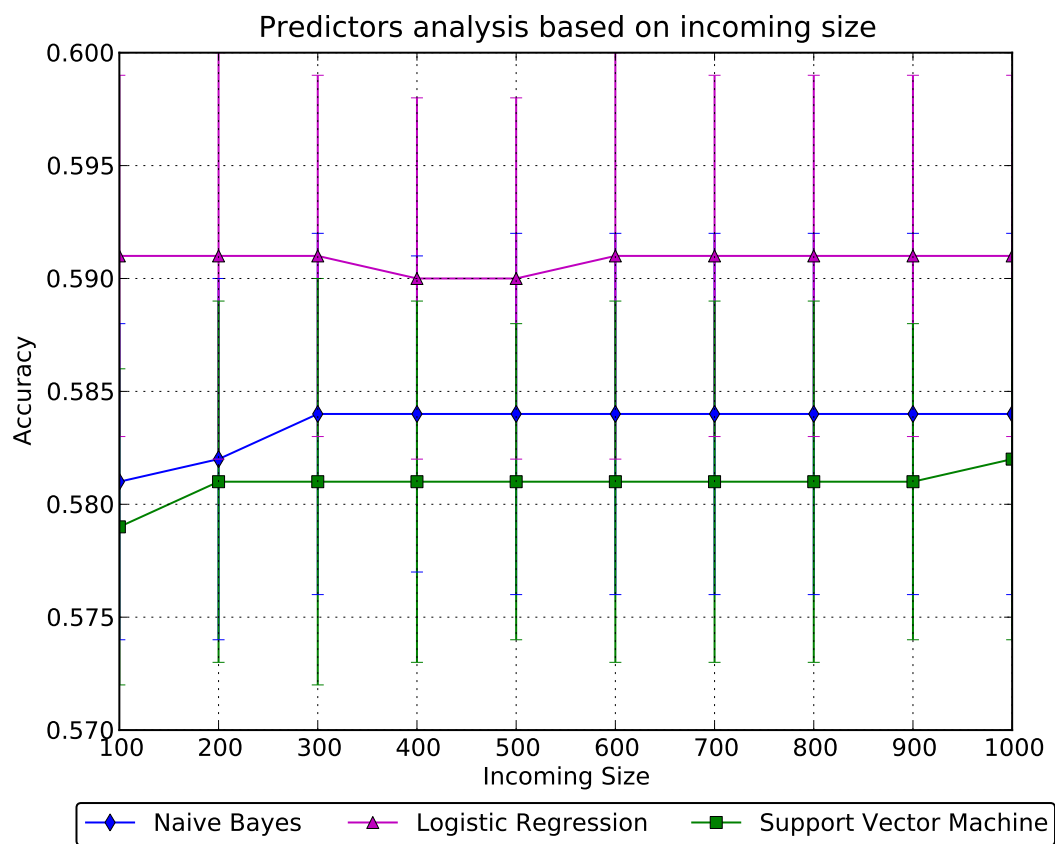


Figure 3.7: Accuracy results for different *Incoming Words* sizes

The most predictive *Incoming Words* words sizes j for each of our classifiers are:

- **Naive Bayes:** 300
- **Logistic Regression:** 100
- **Support Vector Machine:** 1000

Using the most predictive word sizes j for each of our classifiers as defined above and comparing to our baselines we obtain:

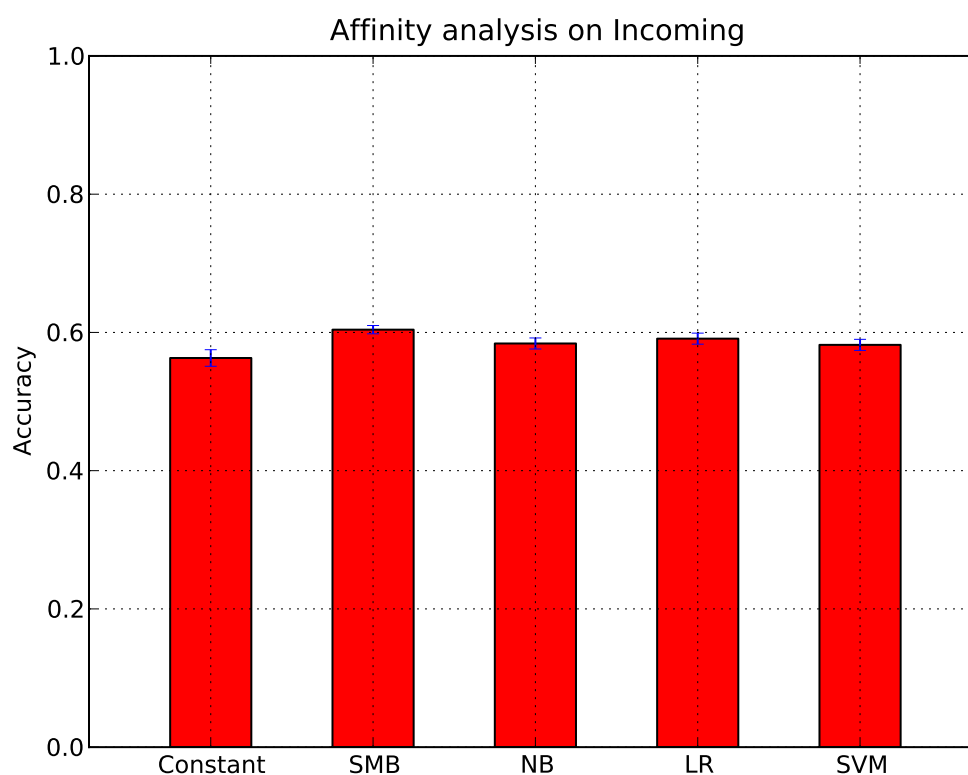


Figure 3.8: Accuracy results using the *Incoming Words* feature set

Again, *Incoming Words* themselves are not predictive in comparison to our baselines, however not to the same extent as *Outgoing Words*.

Comparing *Incoming Messages* against our exposure curve we obtain:

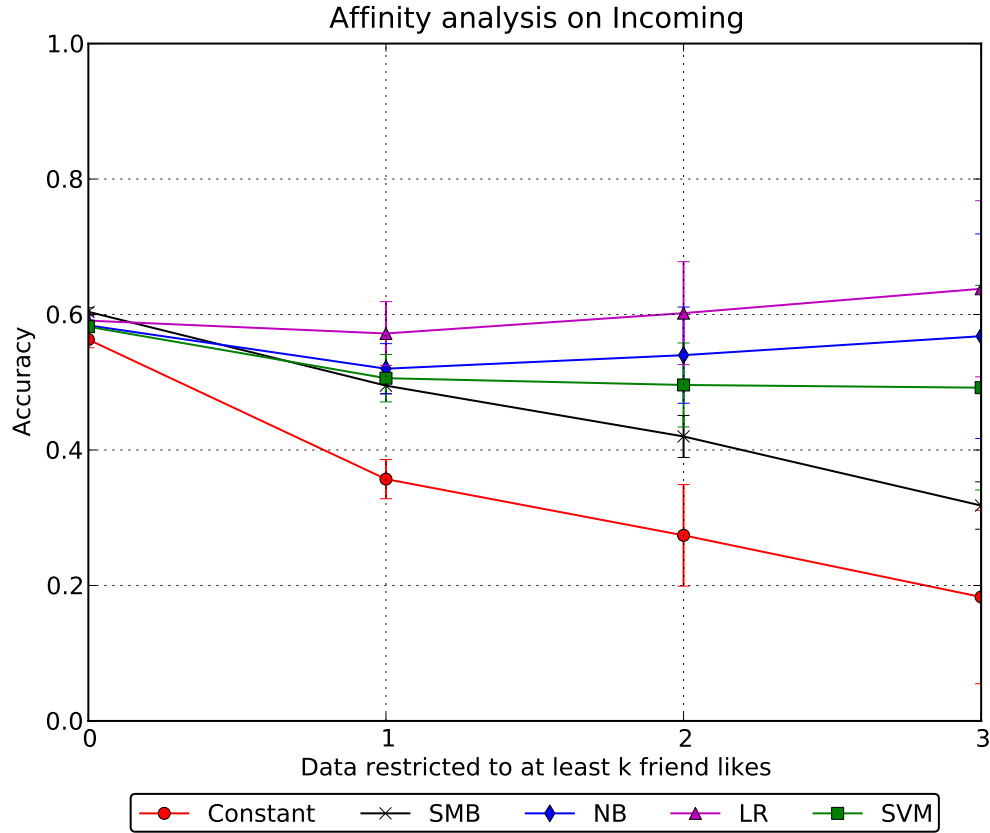


Figure 3.9: Accuracy results for an exposure curve using the *Incoming Words* feature set

Similarly, *Incoming Words* improve upon our baselines as k increases, however this performance increase is negligible in comparison with $k = 0$ and hence *Incoming Words* do not prove to be predictive of user likes.

3.3 Conclusion

Throughout this section we have explored different avenues available to users to maintain interactions between other users.

We have found that words, irrespective of their directionality do not assist in improving predictions. [Anderson et al. 2012] concluded that it is less important what users say, then who they interact with, which we also found in our results our interactions results were comparable to our baselines over $k = 0$ and this improvement continued over the exposure curve as our k increased.

[Brandtzg and Nov 2011] found that virtual interactions help reveal common interests, while real world interactions helps to support friendships.

Our results have shown, that for *User Interactions* it is enough for some user to have previously liked an item to allow our classification methodology to offer an increase in predictiveness as this k increases.

User Preferences

In this section we will discuss the effects of applying different types of *User Preferences* as the feature set and their predictive tendencies within our data set.

4.1 Demographics

The *Demographics* data we are interested in includes:

- Age
- Birthday
- Locale

Below we will give a basic analysis of the *Demographics* data when extracted from our data set.

Gender breakdown:

Male	Female	Undisclosed
85	33	1

Table 4.1: Gender breakdown for app users

Despite this clear male bias [Ugander and Marlow 2011] found that in a social setting, there are no strong gender homophily tendencies. Hence the male skew should not negatively effect our results. Additionally [Backstrom et al. 2011] have shown that different genders have differing tendencies to disperse interactions across genders, implying our male skew should be unimportant. Hence gender information will be used for this feature vector.

Birthday breakdown:

Year	Frequency
Undisclosed	1
1901-1905	1
1906-1910	0
1911-1915	1
1916-1920	0
1921-1925	0
1926-1930	0
1931-1935	0
1936-1940	1
1941-1945	0
1946-1950	0
1951-1955	0
1956-1960	2
1961-1965	1
1966-1970	4
1971-1975	10
1976-1980	12
1981-1985	25
1986-1990	34
1991-1995	25
1996-2000	2

Table 4.2: Birthday breakdown

Birthdays are grouped in a distinct range, most users in this data set are grouped in the age ranges of $\{18 - 30\}$. [Ugander and Marlow 2011] have found that there is a strong effect of age on friendship preferences. Hence birthday information will be used for this feature vector.

Location breakdown:

Location	Frequency
Undisclosed	33
Ahmedabad, India	1
Bangi, Malaysia	1
Bathurst, New South Wales	1
Bellevue, Washington	1
Braddon, Australian Capital Territory, Australia	1
Brisbane, Queensland, Australia	2
Canberra, Australian Capital Territory	56
Culver City, California	1
Frederick, Maryland	3
Geelong, Victoria	1

Table 4.3: Location breakdown

Given the fact that most users are either situated in the ACT (location of the LinkR development and deployment) or are undisclosed, location information will not be used for this feature vector.

For *Demographics* the I of our feature vector X is defined by the following conditions:

- Whether the user is male.
- Whether the user is female.
- Whether the user and any user in the alters set share the same gender.
- Whether the user and any user in the alters set share the same birth range.

The alters of I can then be defined as the set of users who have liked the current item M . Each component of I is set to 1 if any of the alters have meet the conditions described above, in comparison (where required) with the user n , otherwise it is set to 0.

Applying this feature vector to our classifiers we obtain:

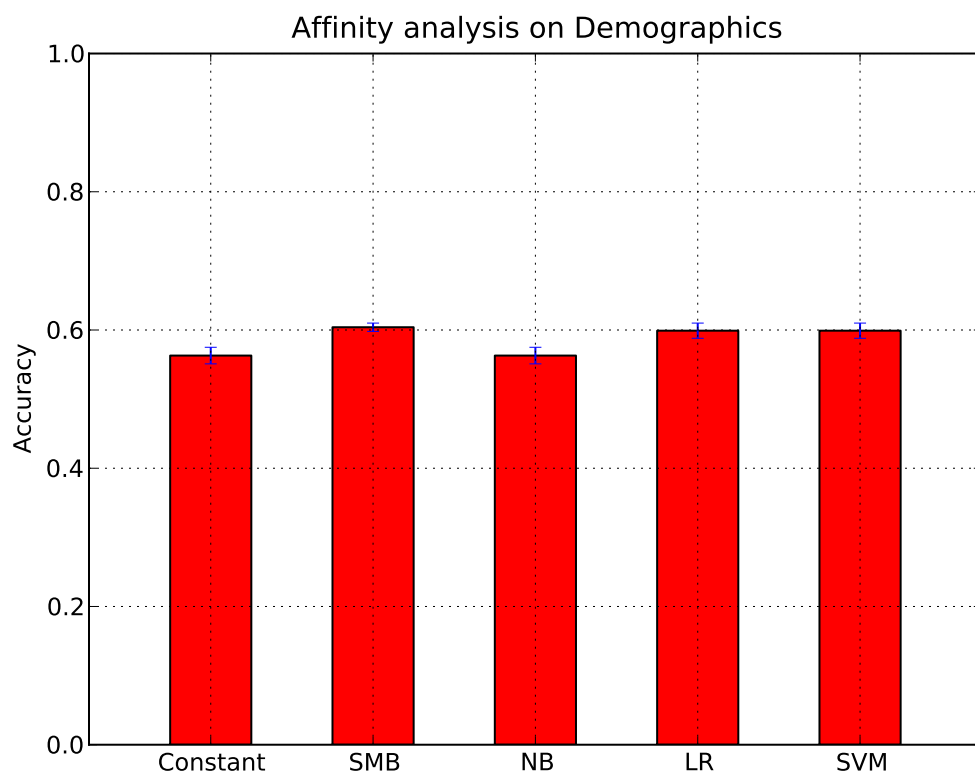


Figure 4.1: Accuracy results using the *Demographics* feature set

The *Demographics* feature vector shows our first positive results, this feature vector almost performs as well as our SMB baseline for the case of $k = 0$.

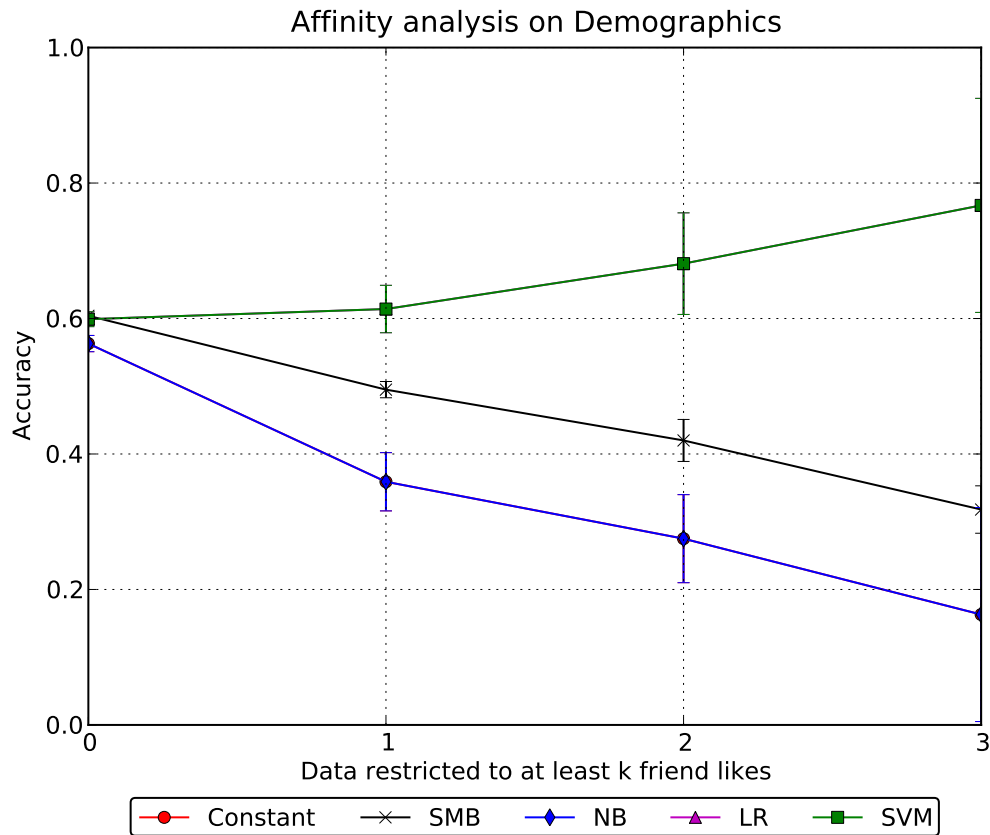


Figure 4.2: Accuracy results for an exposure curve using the *Demographics* feature set. Note in this case Constant = NB and LR = SVM

The exposure curve for *Demographics* shows a sizable improvement over our baselines as our k increases. This demonstrates that as the number of friends who like an item increases, the probability that a user will like that item also increases. This positive correlation between number of likes and user likes increases with k .

4.2 Traits

Facebook facilitates a wide variety of user chosen preferences which we have defined as *Traits*. These *Traits* allow users to define under a specific area of their profile different areas or activities they are interested in.

User *Traits* data we will investigate include:

- Activities
- Books
- Athletes

- Teams
- Inspirational People
- Interests
- Movies
- Music
- Sports
- Television
- School Relationships
- Work Relationships

Below we display graphs for the different *Traits* sets extracted from our data set. Followed by a subsequent analysis. Each table shows only the frequency of app users for each of the *Traits*.

Frequency	Activity
10	Sleeping
5	Eating
5	Reading
4	Running
4	Cycling
4	Minecraft
4	Programming
3	Android
3	Cooking
3	Video Games
3	Xbox 360
3	Piano
3	Guitar
3	Badminton
3	Chocolate

Table 4.4: Top *Activities* for app users

Frequency	Inspirational People
2	Alan Turing
1	Bender
1	Maurice Moss
1	Steve Jobs
1	Sean Parker
1	Pope Benedict XVI
1	Martin Luther
1	Alistair McGrath
1	St Augustine
1	Dennis Ritchie
1	Linus Torvalds
1	Richard Stallman
1	C. S. Lewis
1	Mike Oldfield
1	Ryan Giggs

Table 4.5: Top *Inspirational People* for app users

Frequency	Book
7	Harry Potter
4	The Bible
3	Harry Potter series
3	Discworld
3	That's 3 minutes of solid study, think I've earned 2hrs of Facebook time
3	Freakonomics
3	Tomorrow when the War Began
2	Magician
2	Hitchhiker's Guide To The Galaxy
2	The Discworld Series
2	Terry Pratchett
2	Terry Pratchett
2	George Orwell
2	Lord Of The Rings
2	Goosebumps

Table 4.6: Top *Books* for app users, here we see an example of the non-distinct properties inherent in Facebook, where books can have the same name, yet still be regarded as a different entity.

Frequency	Interest
5	Movies
5	Music
3	Cooking
3	Sports
2	Psychology
2	Internet
2	Video Games
2	Martial arts
2	Literature
2	Economics
2	Tennis
2	Badminton
2	Artificial intelligence
2	Computers
2	Travel

Table 4.7: Top *Interests* for app users

Frequency	Music
9	Daft Punk
9	Muse
8	Michael Jackson
8	Pink Floyd
8	Lady Gaga
7	Linkin Park
7	Avril Lavigne
6	Radiohead
6	Rihanna
6	Coldplay
6	Green Day
6	Katy Perry
6	Taylor Swift
5	Gorillaz
5	Queen

Table 4.8: Top *Music* for app users

Frequency	Movie
9	Inception
8	Avatar
8	Fight Club
7	The Lord of the Rings Trilogy (Official Page)
6	Star Wars
6	I wouldnt steal a car, But i'd download one if i could
6	WALL-E
6	Scott Pilgrim vs. the World
6	Toy Story
6	Shrek
5	Batman: The Dark Knight
5	Harry Potter
4	The Matrix
4	The Social Network Movie
4	Monsters, Inc.

Table 4.9: Top *Movies* for app users

Frequency	Sport
8	Badminton
5	Basketball
3	Cycling
3	Volleyball
2	Starcraft II
2	Football en salle
2	Swimming
2	Towel Baseball
2	Tennis
1	Soccer
1	Taekwondo
1	Rock climbing
1	In The Groove
1	Darts
1	Table tennis

Table 4.10: Top *Sports* for app users

Frequency	Television Show
20	The Big Bang Theory
19	How I Met Your Mother
14	The Simpsons
13	Top Gear
12	Futurama
12	Scrubs
11	Black Books
10	Black Books
10	South Park
10	Family Guy
9	The Daily Show
8	The IT Crowd
8	FRIENDS (TV Show)
7	True Blood
7	MythBusters

Table 4.11: Top *Television* shows for app users

Frequency	Athlete
4	Roger Federer
4	Rafael Nadal
3	Maria Sharapova
2	Leo Messi
1	Andy Schleck
1	Chrissie Wellington
1	Emma Snowsill
1	Emma Moffatt
1	Brbara Riveros
1	The Brownlee Brothers
1	Marie Slamtoinette #1792
1	Wayne Rooney
1	"you are what you eat" " I dont remember eating a Tank."
1	Nemanja Vidic
1	Ryan Giggs

Table 4.12: Top *Athletes* for app users

Frequency	Team
5	Manchester United
2	Bear Grylls cameraman appreciation society
2	Real Madrid C.F.
2	Liverpool FC
1	Leopard Trek
1	British Triathlon
1	TeamCWUK
1	Surly Griffins
1	Canberra Raiders
1	Kolkata Knight Riders
1	Brisbane Roar FC
1	Brisbane Broncos
1	Cricket Australia
1	— Manchester United Fans —
1	Juventus

Table 4.13: Top *Teams* for app users

The *Traits* graphed above can be broken down into three distinct sets based on their locality within the app user base.

- **High Locality:** *Music, Movies, Television* - Showing our app users appear to share similar *Traits* in a media setting.
- **Medium Locality:** *Activities, Books, Interests, Sports* - Showing our app users share some degree of similar preferences across these *Traits*.
- **Low Locality:** *Inspirational People, Athletes, Teams* - Showing our app users do not share many similar preferences across these *Traits*.

For *Traits* the I of our feature vector X is defined by the following conditions:

- Whether the current user and any user in the alters set share the same t where $T \in \{Activities, Books, Athletes, Teams, InspirationalPeople, Interests, Music, Movies, Sports, Television, Work, School\}$.

The alters of I can then be defined as the set of users who have liked the current item M . Each component of I is set to 1 if any of the alters have meet the conditions described above for each t , in comparison with the user n , otherwise it is set to 0.

Applying this feature vector to our classifiers we obtain:

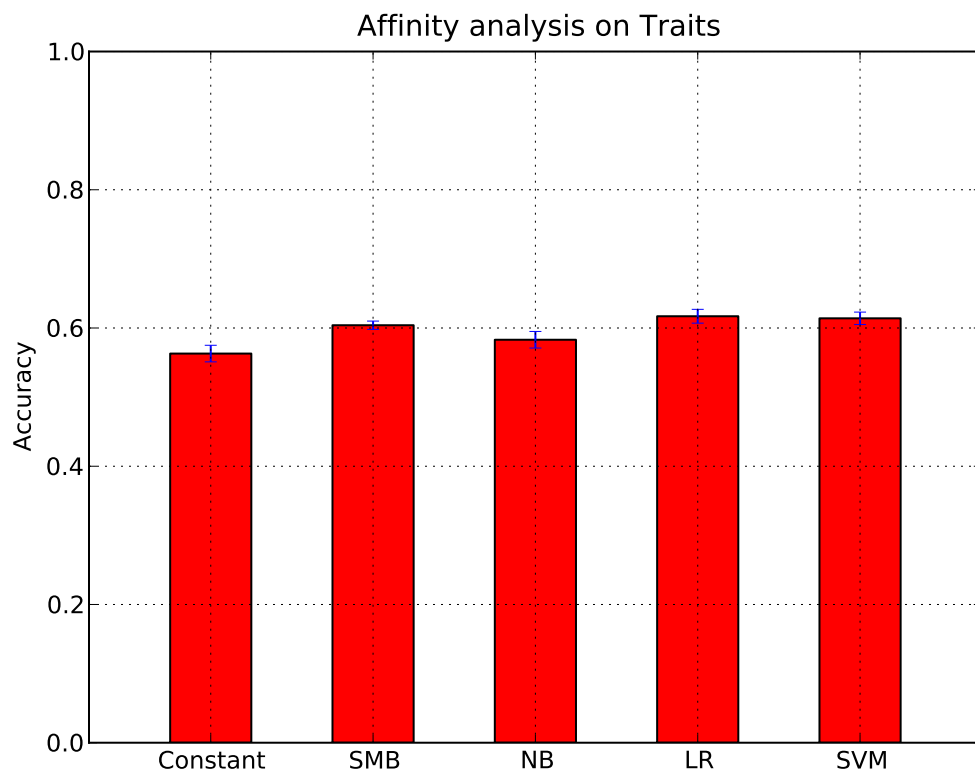


Figure 4.3: Accuracy results using the *Traits* feature set

The *Traits* feature vector shows our first improvement over our SMB baseline in the LR and SVM case for $k = 0$ demonstrating that *Traits* are more predictive than user likes for any previously applied method.

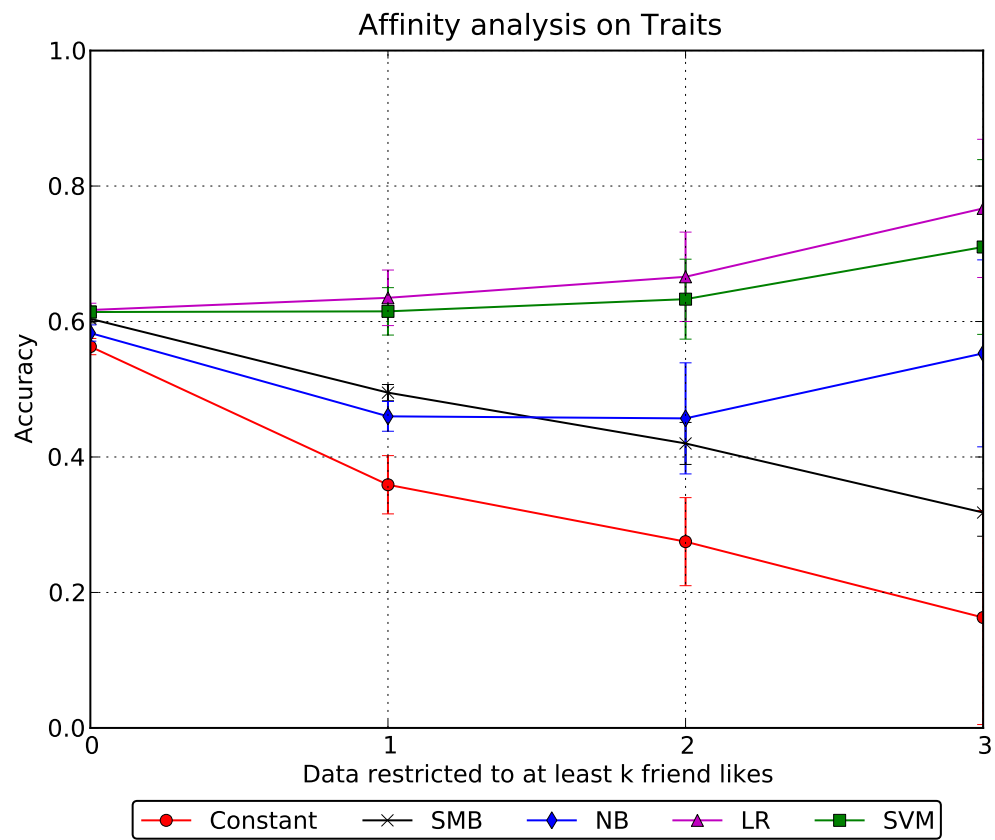


Figure 4.4: Accuracy results for an exposure curve using the *Traits* feature set

This trend continues across the exposure curve where each successive increase of k causes the performance of our classifiers to increase.

By extracting the model weights from the case where $k = 0$ we can see which *Traits* contain the most predictive qualities:

Trait	Weight	Yes'	Distinct Yes'
Activities	-5.927 ± 0.001	281	19
Television	-5.210 ± 0.0	1,029	40
Music	-3.409 ± 0.001	629	31
Movies	-2.668 ± 0.001	454	25
Interests	-1.921 ± 0.001	64	14
Sports	-1.820 ± 0.001	27	7
Books	-1.769 ± 0.0	163	13

Table 4.14: *Logistic Regression* feature weights extracted for the case where $k = 0$. The *Yes'* column displays the number of times this feature vector was set to 1 for a user, the *Distinct Yes'* column displays the number of unique times the feature vector was set to 1.

This shows us that *Traits* which exhibit a medium to high degree of locality have a larger influence during classification.

4.3 Groups

Facebook facilitates users to join *Groups* for a large and varied set of different types ranging from local sports teams, political preferences to computer games.

The most popular groups for our app users are shown below:

Group Name	Frequency
27	ANU StalkerSpace
20	Facebook Developers
15	ANU CSSA
14	CSSA
13	Australian National University
11	ANU - ML and AI Stanford Course
10	iDiscount ANU
10	Our Hero: Clem Baker-Finch
9	Students In Canberra
7	I grew up in Australia in the 90s
7	Grow up Australia - R18+ Rating for Computer Games
7	ANU Engineering Students' Association (ANUESA) 2010
7	ANU Postgraduate and Research Student Association (PARSA)
6	No, I Don't Care If I Die At 12AM, I Refuse To Pass On Your Chain Letter.
6	No Australian Internet Censorship
6	The Chaser Appreciation Society
6	Feed a Child with a Click
6	ANU Mathematics Society
6	ANU International Student Services, CRICOS Provider Number 00120C
6	2011 New & Returning Burton & Garran Hall
5	If You Can't Differentiate Between "Your" and "You're" You Deserve To Die
5	Keep the ANU Supermarket!!!
5	If 1m people join, girlfriend will let me turn our house into a pirate ship
5	The Great Australian Internet Blackout
5	When I was your age, Pluto was a planet.

Table 4.15: App users popular *Groups* breakdown

In comparison with *Traits*, *Groups* show a higher locality for the most popular groups.

Given the quantity of groups on Facebook, we need to find some optimal test size j for our data set. Given memory and time constraints we tested within a range of $\{100 - 1000\}$ with an incremental step size of 100 for each test.

The results for these tests are shown below:

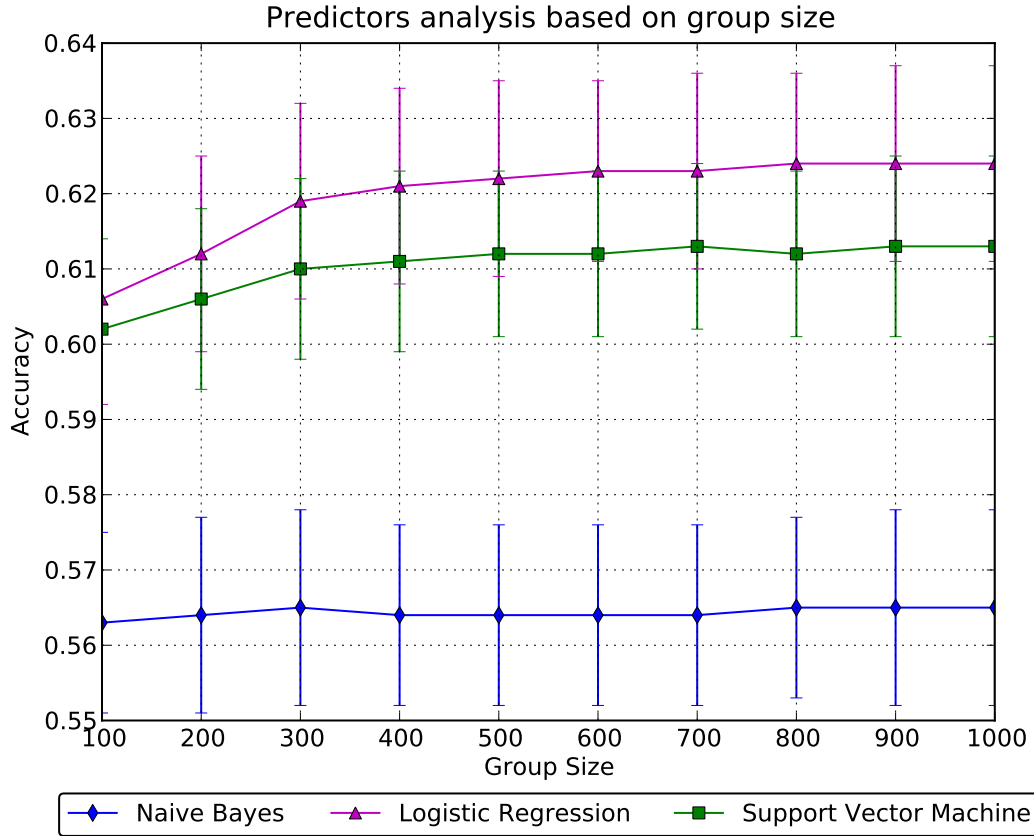


Figure 4.5: Accuracy results for different *Groups* sizes

The most predictive *Group* sizes j for each of our classifiers are:

- **Naive Bayes:** 300
- **Logistic Regression:** 900
- **Support Vector Machine:** 800

LR and SVM show a gradual increase as this group size increases, alluding to the possibility of an even higher group size being optimal.

For *Groups* the I of our feature vector X contains an element i for each of the top j groups sizes defined above.

The alters of I can then be defined as all users who have liked the current item M . Each component of I is set to 1 if any of the alters are a member of the current group j along with the current user n , otherwise it is set to 0.

Using the most predictive *Group* sizes j for each of our classifiers as defined above and comparing to our baselines we obtain:

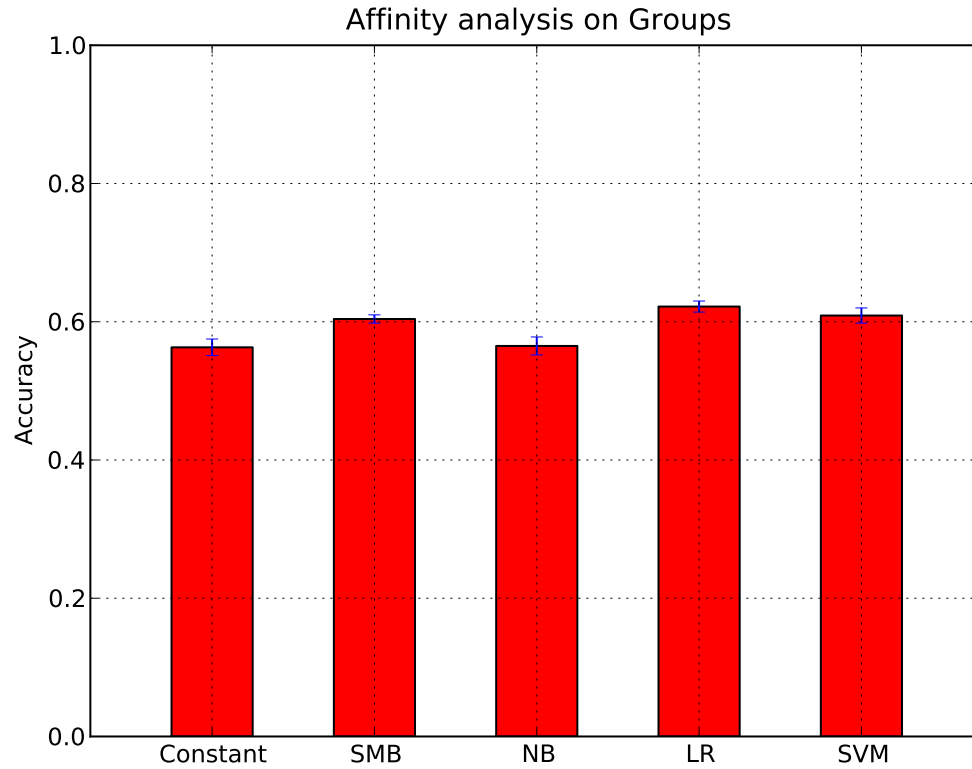


Figure 4.6: Accuracy results using the *Groups* feature set

Both LR and SVM show an improvement over our SMB baseline for the case of $k = 0$ demonstrating that *Groups* are more predictive than previous methods.

Applying the *Groups* feature vector across our exposure curve, we obtain:

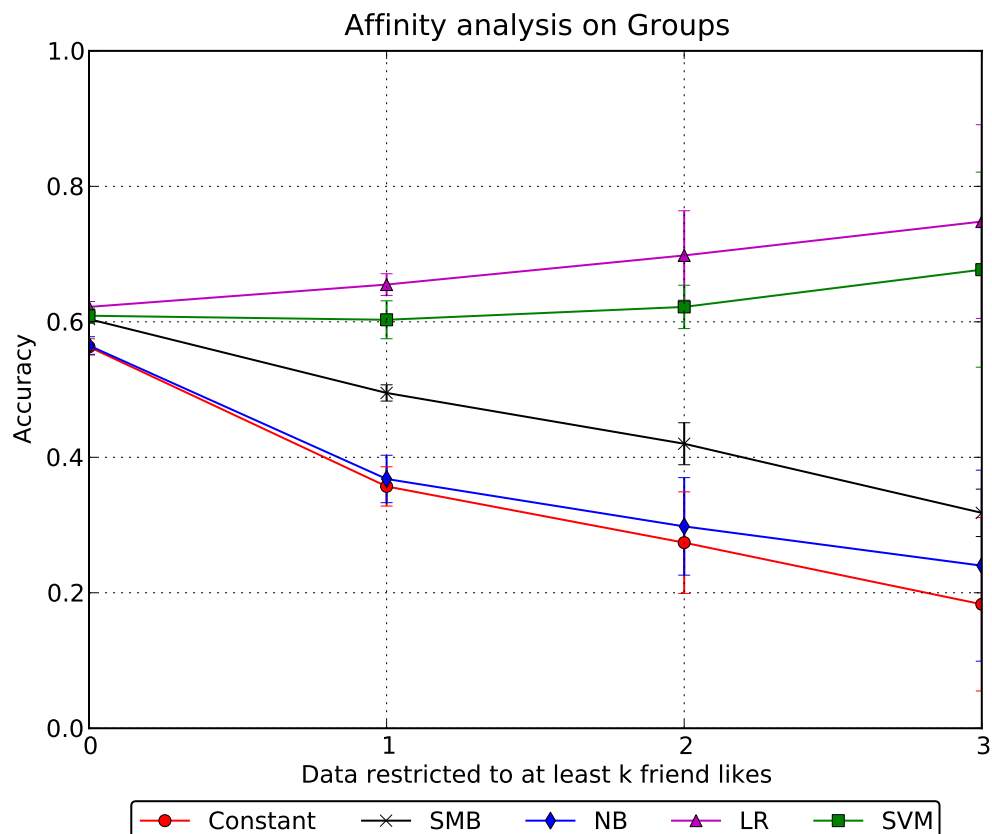


Figure 4.7: Accuracy results for an exposure curve using the *Groups* feature set

This trend continues across the exposure curve where each successive increase of k causes the performance of our LR and SVM classifiers to increase.

By extracting the model weights from the case where $k = 0$ we can see which *Groups* contain the most predictive qualities: This shows us that ...

4.4 Pages

Facebook facilitates users to like *Pages* for 'things' they like across a large and varied set of different areas ranging from web browsers, TV shows to schools.

Group	Weight	Yes'	Distinct Yes'
Activities	-5.927 ± 0.001	281	19
Television	-5.210 ± 0.0	1,029	40
Music	-3.409 ± 0.001	629	31
Movies	-2.668 ± 0.001	454	25
Interests	-1.921 ± 0.001	64	14
Sports	-1.820 ± 0.001	27	7
Books	-1.769 ± 0.0	163	13

Table 4.16: *Logistic Regression* feature weights extracted for the case where $k = 0$. The *Yes'* column displays the number of times this feature vector was set to 1 for a user, the *Distinct Yes'* column displays the number of unique times the feature vector was set to 1.

The most popular pages liked by our app users are shown below:

Page Name	Frequency
33	ANU Computer Science Students' Association (ANU CSSA) 2011
32	The Australian National University
31	ANU Stalkerspace
21	Humans vs Zombies @ ANU
20	The Big Bang Theory
19	Australian National University
19	How I Met Your Mother
18	ANU LinkR
18	ANU ducks
17	Australian National University Students' Association
16	Google
15	Google Chrome
15	ANU XSA
15	Facebook
14	YouTube
14	The Simpsons
13	Portal
13	Top Gear
13	Music
13	ANU Memes
12	Futurama
12	Scrubs
12	ANU O-Week 2012: Escape to the East
12	The Stig
11	Black Books

Table 4.17: App users *Pages* breakdown

In comparison with *Traits* and *Groups*, *Pages* show an higher locality across the most popular pages for app users.

Given the quantity of pages on Facebook, we need to find some optimal test size j for our data set. Given memory and time constraints we tested within a range of $\{100 - 1000\}$ with an incremental step size of 100 for each test.

The results for these tests are shown below:

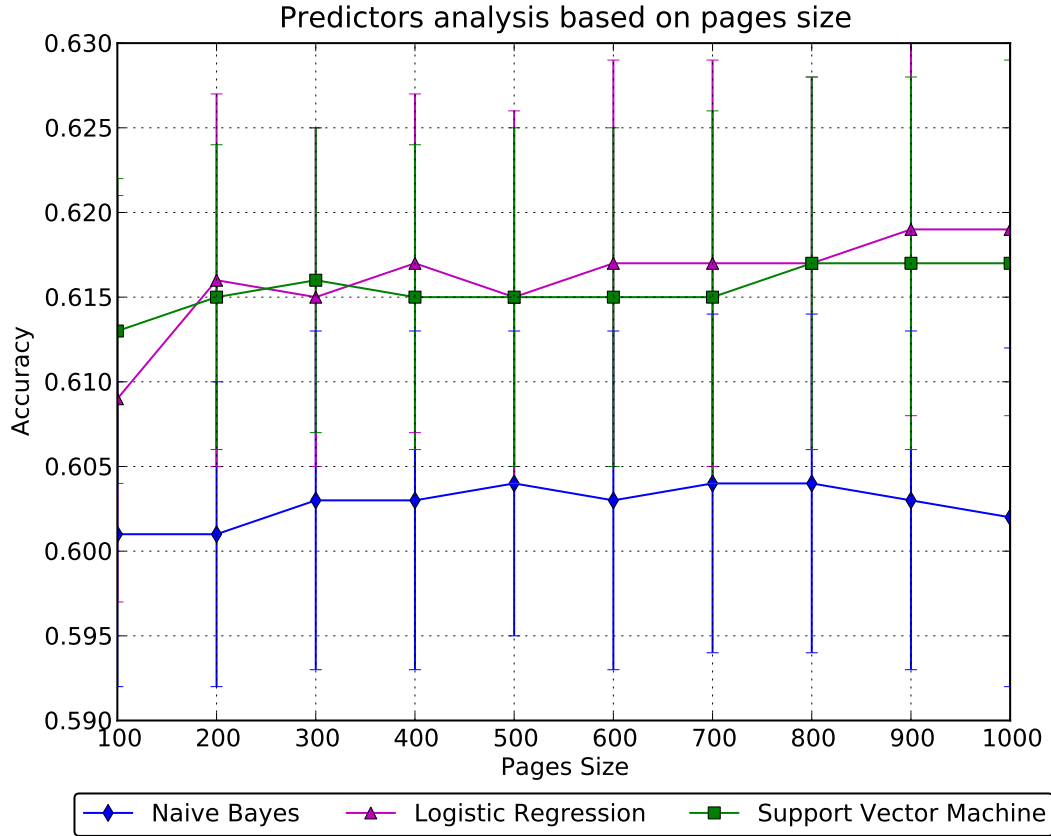


Figure 4.8: Accuracy results for different *Pages* sizes

The most predictive *Page* sizes j for each of our classifiers are:

- Naive Bayes: 500
- Logistic Regression: 900
- Support Vector Machine: 800

LR and SVM show a gradual increase as this group size increases, alluding to the possibility of an even higher page size being optimal.

For *Pages* the I of our feature vector X contains an element i for each of the top j page sizes defined above.

The alters of I can then be defined as all users who have liked the current item M . Each component of I is set to 1 if any of the alters are a have liked the current page j along with the current user n , otherwise it is set to 0.

Using the most predictive *Page* sizes j for each of our classifiers as defined above and comparing to our baselines we obtain:

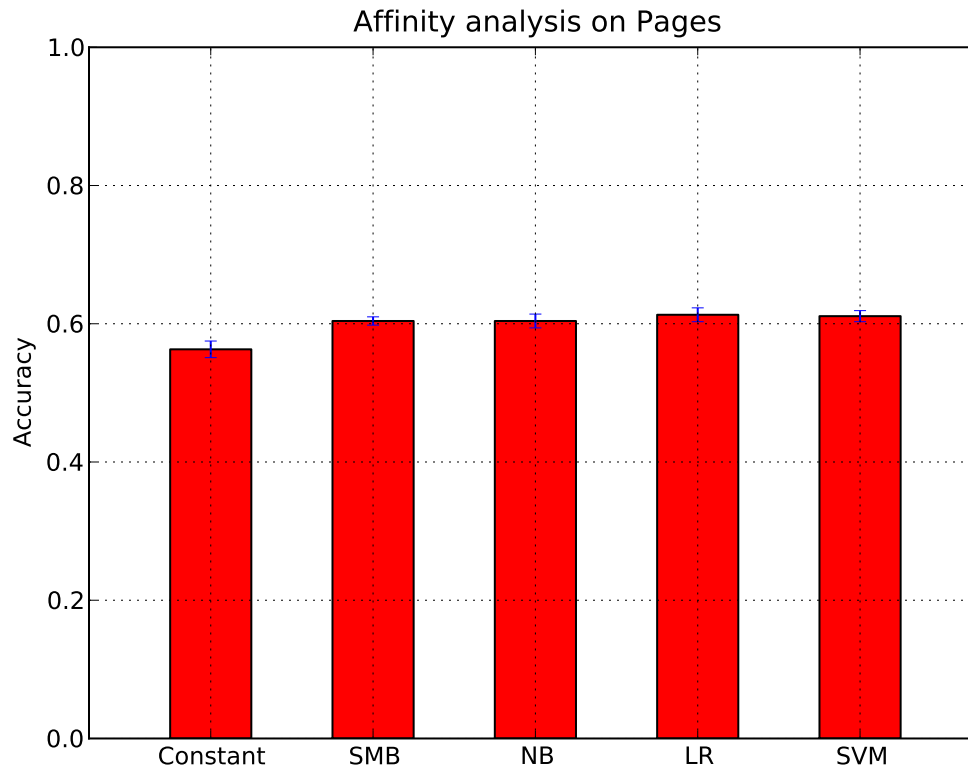


Figure 4.9: Accuracy results using the *Pages* feature set

Both NB, LR and SVM show an improvement over our SMB baseline for the case of $k = 0$ demonstrating that *Pages* are more predictive than previous methods.

Applying the *Pages* feature vector across our exposure curve, we obtain:

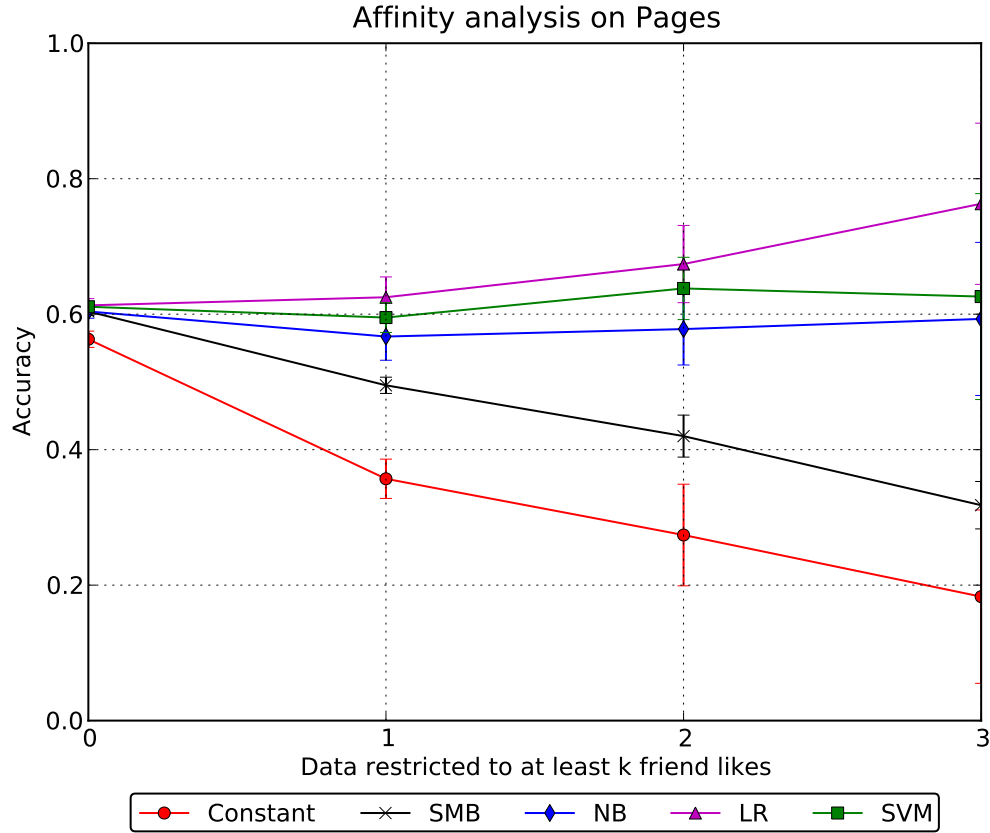


Figure 4.10: Accuracy results for an exposure curve using the *Pages* feature set

This trend continues across the exposure curve where each successive increase of k causes the performance of our LR and SVM classifiers to increase.

By extracting the model weights from the case where $k = 0$ we can see which *Pages* contain the most predictive qualities: PAGES HERE

This shows us that ...

4.5 Conclusion

Throughout this section we have explored different avenues available for users to demonstrate their personal preferences across a range of different mediums.

We have found that *User Preferences* are predictive of user likes, particularly for *Traits*, *Groups* and *Pages*. This holds true for the case of $k = 0$ and continues to improve with k .

Similarly as with *User Interactions*, our results have shown, that it is enough for some user to have liked an item to allow our classification methodology to offer an

increase in predictiveness.

Feature Combinations

As outlined above, features which positively improved classification were from the *Traits*, *Groups* and *Pages* feature vectors. In this section we combine these positive feature vectors together into a larger feature vector comprised of the individual positively contributing elements.

5.1 Positive Feature Combination

Using the combined feature vector X where I is comprised of:

- *Traits*
- *Groups*
- *Pages*

Applying this feature vector to the data set:

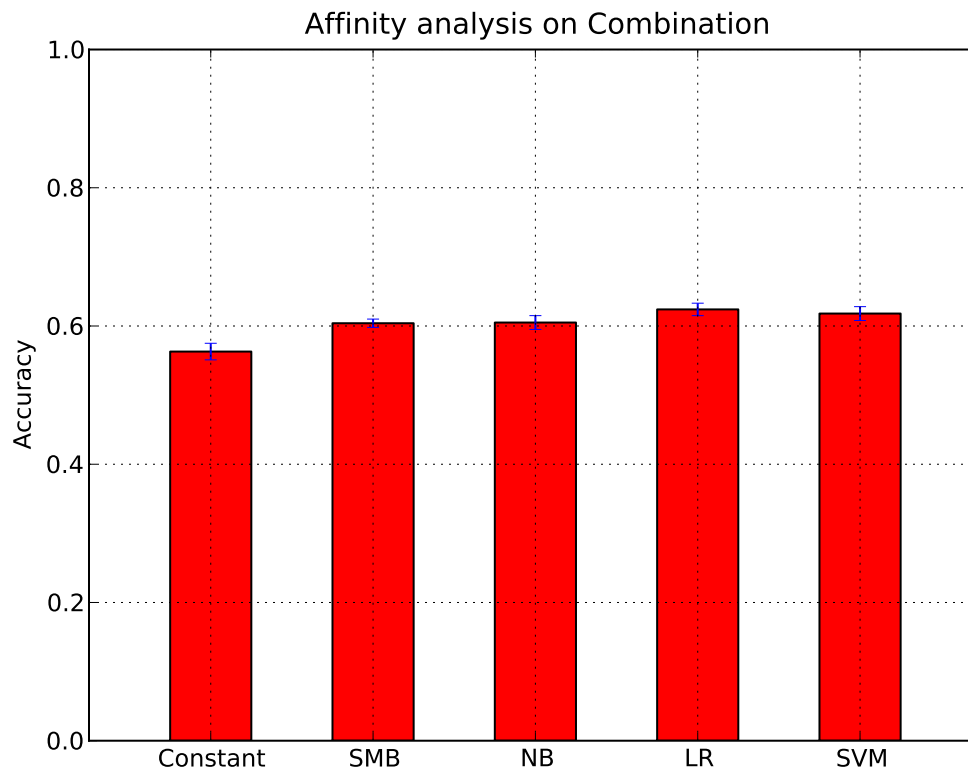


Figure 5.1: Accuracy results using the *Positively Combined* feature set.

We find that the *Combination* feature vector gives better results for our classifiers when compared with our baselines. This holds for all values of k and offers the most predictive feature vector found during this research.

This can be summarised in the table below:

Classifier	Accuracy
SMB	0.604 ± 0.006
NB	0.583 ± 0.012
LR	0.617 ± 0.01
SVM	0.614 ± 0.009

Table 5.1: *Traits* results for $k = 0$

Classifier	Accuracy
SMB	0.604 ± 0.006
NB	0.604 ± 0.01
LR	0.613 ± 0.01
SVM	0.611 ± 0.008

Table 5.2: *Pages* results for $k = 0$

Classifier	Accuracy
SMB	0.604 ± 0.006
NB	0.565 ± 0.013
LR	0.622 ± 0.008
SVM	0.609 ± 0.011

Table 5.3: *Groups* results for $k = 0$

Classifier	Accuracy
SMB	0.604 ± 0.006
NB	0.605 ± 0.01
LR	0.624 ± 0.009
SVM	0.618 ± 0.01

Table 5.4: *Combined* results for $k = 0$

These tables show that based on our results the most predictive feature vector is a combination of the individual best feature vectors found in our *User Preferences* section.

Applying this same feature vector across our exposure curve:

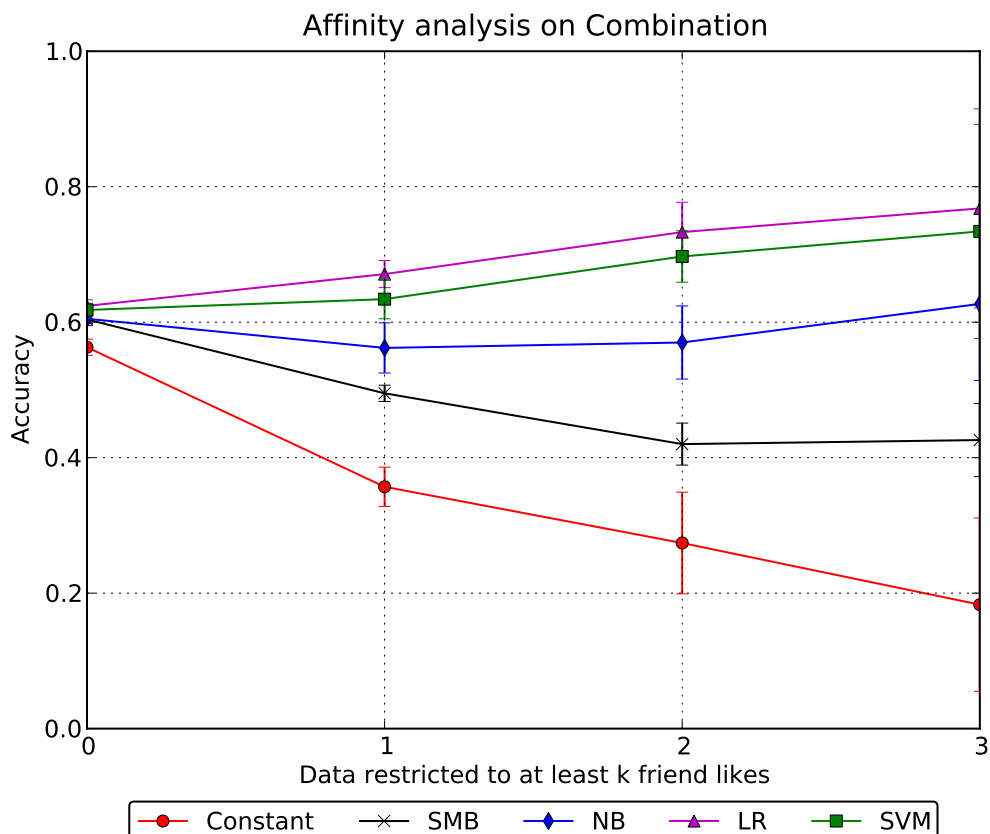


Figure 5.2: Accuracy results for an exposure curve using the *Positively Combined* feature set.

This trend continues over the exposure curve with LR, SVM and NB all improving as k increases. Again, this feature vector combination provides the most predictive results when compared with all other analysis completed during this thesis.

By extracting the model weights from the case where $k = 0$ we can see which components of the *Combination* feature vector were most predictive:

PAGES HERE

This shows us that ...

5.2 Summary

In this thesis we have tested and compared different feature vectors across different exposures of size k . We have shown that *User Interactions* in themselves are not predictive of user likes, however coupled with the likes exposure curve, they do show an improvement over our baselines as k increases.

We have also shown the interesting result that *User Preferences* are predictive of

user likes in the base case of $k = 0$ and this trend continues over the likes exposure curve.

To answer the question initially proposed for this thesis, we have shown the feature vector which provides the highest predictiveness of user likes is the combined vector comprised of *Traits*, *Groups* and *Pages*. These were the highest performing individual feature vectors and combined, represent the most predictive feature vector across our testing scope.

Which is the exciting novel insight examined by this thesis.

5.3 Future Work

Proposed future work can be summarised under the following points:

- **Increase size ranges:** Given our maximum test sizes for *Groups* and *Pages* of 1000 this size could be increased to find the optimal testing range for each of our classifiers.
- **Individual *Traits* analysis:** During our *Traits* analysis the feature vector was set to 1 if the user and any user in the set of alters were part of the same *Traits* group, it could be beneficial to do an individual analysis on each component of the *Traits* data to find which individual elements of each *Trait* are most predictive (similar to the analysis as done for *Groups* and *Pages*).
- **Passive likes:** Given the Facebook model of allowing users to like but not dislike data, explicit dislike data can not be gleaned from Facebook, which is hence why the LinkR active likes data was used for this evaluation. An approach could be developed which can predict whether a user will have seen an item (online timestamps, recent interactions with user) and can infer that if the user did not like the item then they disliked it. This data set could then be applied to the testing methodology outlined above.
- **Cold start:** Leaving out some subset of users when training our models, but including them during testing to explore their effects on results.
- **General user set:** Such as the study done by [Ugander and Marlow 2011] which comprised of the entire active social network of 721 million users as of May 2011, applying these methods to a data set which is more indicative of the general Facebook user population could offer more generalisable results.
- **Bayesian Model Averaging:** Weighting the most successful machine learning models under different feature sets and exposure curves to simulate a new combined classifier, which has learnt from the positive results of each individual classifier.

Bibliography

- ALIAS-I. 2008. LINGPIPE 4.1.0. [HTTP://ALIAS-I.COM/LINGPIPE](http://alias-i.com/lingpipe) (ACCESSED OCTOBER 1, . 2011. LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology* 2, 27:1–27:27. (p.9)
- ANDERSON, A., HUTTENLOCHER, D., KLEINBERG, J., AND LESKOVEC, J. 2012. Effects of User Similarity in Social Media. In *ACM International Conference on Web Search and Data Mining (WSDM)* (2012). (p.21)
- BACKSTROM, L., BAKSHY, E., KLEINBERG, J., LENTO, T., AND ROSENN, I. 2011. Center of attention: How facebook users allocate attention across friends. *ICWSM’11* (2011). (p.23)
- BRANDTZG, P. B. AND NOV, O. 2011. Facebook use and social capital — a longitudinal study. *ICWSM’11* (2011). (p.21)
- CHANG, C.-C. AND LIN, C.-J. 2011. LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology* 2, 27:1–27:27. (p.9)
- CUI, P., WANG, F., LIU, S., OU, M., AND YANG, S. 2011. Who should share what? item-level social influence prediction for users and posts ranking. In *International ACM SIGIR Conference (SIGIR)* (2011). (p.8)
- GRANOVETTER, M. S. 1978. Threshold models of collective behavior. *Am. J. Sociol* 83(6):14201443. (p.1)
- HILL, R. AND DUNBAR, R. 2003. Social network size in humans. *Human Nature* 14, 1, 53–72. (p.5)
- LANG, K. 1995. NewsWeeder: Learning to filter netnews. In *12th International Conference on Machine Learning ICML-95* (1995), pp. 331–339. (p.7)
- NOEL, J. G. 2011. New social collaborative filtering algorithms for recommendation on facebook (2011). (pp.7, 8)
- PANTEL, A., GAMON AND HAAS. 2012. *Proceeding SIGIR ’12 Proceedings of the 35th international ACM SIGIR conference on Research and development in information retrieval*. ACM. (p.1)
- RESNICK, P. AND VARIAN, H. R. 1997. Recommender systems. *Communications of the ACM* 40, 56–58. (p.7)
- ROMERO, D. M., MEEDER, B., AND KLEINBERG, J. 2011. Differences in the mechanics of information diffusion across topics: idioms, political hashtags, and complex contagion on twitter. *WWW ’11* (2011), pp. 695–704. ACM. (p.1)

- SAEZ-TRUMPER, D., NETTLETON, D., AND BAEZA-YATES, R. 2011. High correlation between incoming and outgoing activity: A distinctive property of online social networks? In *Fifth International AAAI Conference on Weblogs and Social Media, ICWSM '11* (2011). (p.11)
- SANGHVI, R. AND STEINBERG, A. 2010. Edgerank: The secret sauce that makes facebook's news feed tick (2010). (p.5)
- UGANDER, B., KARRER AND MARLOW. 2011. The anatomy of the facebook social graph. *CoRR abs/1111.4503*. (pp.23, 24, 53)
- WATTS, D. J. AND STROGATZ, S. H. 1998. Collective dynamics of small-world networks. *Nature* 393(6684):440442. (p.1)
- YANG, LONG, SMOLA, SADAGOPAN, ZHENG, AND ZHA. 2011. Like like alike: Joint friendship and interest propagation in social networks. In *WWW-11* (2011). (p.8)