SPECIALE

# FRA HADES
## OF
## DOOM

BJARKE MØNSTED

PRETENTIOUS
QUOTE.

- FAMOUS PERSON, BORN-DIED

# Fra Hades

## of
# Doom

|  |  |
|---|---|
| Author | My Name |
| Advisor | His Name |
| Co-Advisor | Her Name |

*Ting, som KU siger der skal stå her*

## CMOL

Center for Models of Life

Submitted to the University of Great Justice
June 14, 2014

# Acknowledgements

Thank you! Thank you all!

# Contents

# Englishabstract

WORDS! SOOOOO MANY WORDS!

# DANSK SAMMENFATNING

ORD! SAAAAAAAAAA MANGE ORD

Part I

WIKIPEDIA-BASED EXPLICIT SEMANTIC ANALYSIS

# 1

# WIKIPEDIA-BASED EXPLICIT SEMANTIC ANALYSIS

ATURAL language processing has long been both a subject of interest and a source of great challenges in the field of artificial intelligence. The difficulty varies greatly depending with the different language processing tasks; certain problems, such as text categorization, are relatively straightforward to convert to a purely mathematical problem, which in turn can be solved by a computer, whereas other problems, such as computing semantic relatedness, necessitates a deeper understanding of a given text, and thus poses a greater problem. This sections aims firstly to give a brief introduction to some of the most prominent techniques used in language processing in order to explain my chosen method of explicit semantic analysis (ESA), and secondly to explain in detail my practical implementation of an ESA-based text interpretation scheme.

ref

## 1.1   Methods

This section outlines a few methods used in natural language processing, going into some detail on ESA while touching briefly upon related techniques.

### 1.1.1   Bag-of-Words

An example of a categorization problem is the 'bag of words' approach, which has seen use in spam filters.Here, text fragments are treated as unordered collections of words drawn from various bags, which in the case of spam filters would be undesired mails (spam) and desired mails (ham). By analysing large amounts of regular mail and spam, the probability of drawing each of the words constituting a given text from each bag can be computed, and the probability of the given text fragment representing draws from each bag can be computed using Bayesian statistics.

ref

More formally, the text $T$ is represented as a collection of words $T = \{w_1, w_2, \ldots, w_n\}$, and the probability of $T$ actually representing draws from bag $j$ is hence

$$P(B_j|T) = \frac{P(T|B_j)P(B_j)}{P(T)}, \tag{1.1}$$

$$= \frac{\prod_i P(w_i|B_j)P(B_j)}{\sum_j \prod_i P(w_i|B_j)P(B_j)}, \tag{1.2}$$

for an arbitrary number of bags labelled by $j$. This method is simple and powerful whenever a text is expected to fall in one of several discrete categories (such as spam filters or language detection). However, for more complex tasks it proves lucrative to attempt instead to assign some kind of meaning to text fragments rather than to consider them analogous to lottery numbers or marbles. This notion of meaning will be elaborated on shortly, as it varies depending on the method of choice, but the overall idea is to ascribe to words a meaning which depends not only on the word itself, but also on the connection between the word and and existing repository of knowledge. The reader may think of this as mimicking the reading comprehension of humans. In itself, the word 'dog' for instance, contains a mere 24 bits of information if stored with a standard encoding, yet a human reader immediately associates a rich amount of existing knowledge to the word, such as dogs being mammals, related to wolves, being a common household pet, etc. The objective of both explicit and latent semantic analysis is to establish a high-dimensional 'concept space' in which words and text fragments are represented as vectors. The difference between explicit and latent semantic analysis is the method used to obtain said concepts, a explained in the following sections.

ref

### 1.1.2  Semantic Analysis

Salton et al proposed in their 1975 paper *A Vector Space Model for Automatic Indexing*[15] an approach where words and text fragments are mapped

with a linear transformation to vectors in a high-dimensional concept space,

$$T \rightarrow |V\rangle = \sum_i v_i |i\rangle, v_i \in \mathbb{R}, \tag{1.3}$$

where a similarity measure of two texts can be defined as the inner product of two normalized such vectors,

$$S(V, W) = \langle \hat{V} | \hat{W} \rangle = \frac{\sum_i v_i w_i}{\left( \sum_i v_i^2 \right) \left( \sum_i w_i^2 \right)}, \tag{1.4}$$

and the cosine quasi-distance can be considered as a measure of semantic distance between texts:

$$D(V, W) = 1 - S(V, W). \tag{1.5}$$

This approach has later seen use in the methods of Latent Semantic Analysis (LSA) and Explicit Semantic Analysis (ESA). Both methods can be said to mimic human cognition in the sense that the transformation from (1.3) is viewed as a mapping of a text fragment to a predefined *concept space* and thus, processing of texts relies heavily on external repositories of knowledge.

The difference between LSA and ESA is how the concept space is established. Although I have used solely ESA for this project, I will give an extremely brief overview of LSA for completeness following (Landauer 1998 [11]). LSA constructs its concept space by first extracting every unique word encountered in a large collection of text corpora and essentially uses the leading eigenvectors (i.e. corresponding to the largest eigenvalues) of the word-word covariance matrix as the basis vectors of its conceptual space. This is the sense in which the concept are latent - rather than interpret text in terms of explicit concepts, such as 'healthcare', LSA would discover correlations between words such as 'doctor', 'surgery' etc. and consider that a latent concept. Owing to the tradeoff between performance and computational complexity, only about 400 such vectors are kept[13]. In psychology, LSA has been proposed as a possible model of fundamental human language acquisition as it provides computers a way of estimating e.g. word-word relatedness (a task which LSA does decently) using nothing but patterns discovered in the language it encounters[11].

In contrast, the concepts in ESA correspond directly to certain parts of the external text corpora one has employed to construct a semantic analyser. Concretely, the matrix playing the role of the reduced covariance

matrix in LSA has columns corresponding to each text corpus used and rows corresponding to individual terms or words, with the value of each matrix element denoting some measure of relatedness between the designated word and concept. I have used the English Wikipedia, so naturally each concept consists of an article, although the process could easily be tuned to be more or less fine-grained and associate instead each concept with e.g. a subsection or a category, respectively. Of course, a wholly different collection of texts could also be used - for instance a version of ESA more suited to compare the style or period of literary works could be constructed using a large collection of literature such as the Gutenberg Project. However, with no prior knowledge of the subject matter of the text to be analysed, Wikipedia seems like a good all-round solution considering its versatility and the massive numbers of volunteers constantly keeping it up to date.

I wish to point out two advantages of ESA over LSA. First, it is more successful, at least using the currently available text corpora and implementation techniques. A standard way of evaluating the performance of a natural language processing program is to measure the correlation between relatedness scores assigned to pairs of words or text fragments by the computer and by human judges. In both disciplines, ESA has outperformed LSA since it was first implemented ([8], p 457).

Second, the concepts employed in ESA are directly understandable by a human reader, whereas the concepts in LSA correspond to the leading moments of the covariance matrix. For example, to test whether the first semantic analyser built by my program behaved reasonably, I fed it a snippet of a news article on CNN with the headline "*In Jerusalem, the 'auto intifada' is far from an uprising*". This returned an ordered list of top scoring concepts as follows: "Hamas, Second Intifada, Palestinian National Authority, Shuafat, Gaza War (2008-09), Jerusalem, Gaza Strip, Arab Peace Initiative, Yasser Arafat, Israel, West Bank, Temple Mount, Western Wall, Mahmoud Abbas", which seems a very reasonable output.

## 1.2   Constructing a Semantic Analyser

The process of applying ESA to a certain problem may be considered as the two separate subtasks of first a very computationally intensive construction af the machinery required to perform ESA, followed by the application of said machinery to some collection of texts. For clarity, I'll limit the present section to the details of the former subtask while description its application and results in section 1.3.

The construction itself is divided into three steps which are run in succession to create the desired machinery. The following is a very brief overview of these steps, each of which is elaborated upon in the following subsections.

1. First, a full Wikipedia XML-dump[1] is parsed to a collection of files each of which contains the relevant information on a number of articles. This includes the article contents in plaintext along with some metadata such as designated categories, inter-article link data etc.

2. Then, the information on each concept (article) is evaluated according to some predefined criteria, and concepts thus deemed inferior are purged from the files. Furthermore, two list are generated and saved, which map unique concepts and words, respectively, to an integer, the combination of which is to designate the relevant row and column in the final matrix. For example, the concept 'Horse' corresponded to column 699221 in my matrix, while the word 'horse' corresponded to row 11533476.

3. Finally, a large sparse matrix containing relevance scores for each word-concept pair is built and, optionally, pruned (a process used to remove 'background noise' from common words as explained in section 1.2.3)

These steps are elaborated upon in the following.

## 1.2.1   XML Parsing

The Wikipedia dump comes in a rather large (~50GB unpacked for the version I used) XML file which must be parsed to extract each article's contents and relevant information. This file is essentially a very long list of nested fields where the data type in each field is denoted by an XML tag, such as `<text> blabla </text>`. A very simplified example of a field for one Wikipedia article is shown in 1.1 The content of each field has already been sanitised by Wikipedia so that if for instance the symbol '<' is entered into an article, it is instead represented as '&lt' in the XML file. To this end, I wrote a SAX parser, which processes the dump sequentially to accommodate its large size. When running, the parser walks through the file and sends the various elements it encounters to a suitable processing function depending on the currently open XML tag. For example, when

---

[1]These are periodically released at `http://dumps.wikimedia.org/enwiki/`

```
  <page>
    < title >Horse</title>
    <ns>0</ns>
    <id>12</id>
    <revision>
      <id>619093743</id>
      <parentid>618899706</parentid>
      <timestamp>2014−07−30T07:26:05Z</timestamp>
      <contributor>
        <username>Eduen</username>
        <id>7527773</id>
      </contributor>
      <text xml:space="preserve">
===Section title===
[[Image:name_of_image_file] Image caption]]]
Lots of educational text containing, among other things links to  [[other  article | text  to  display ]].
</text>
      <sha1>n57mnhttuhxxpq1nanak3zhmmmcl622</sha1>
      <model>wikitext</model>
      <format>text/x−wiki</format>
    </revision>
  </page>
```

Snippet 1.1: A simplified snippet, of a Wikipedia XML dump.

a 'title' tag is encountered, a callback method is triggered which assigns a column number to the article title and adds it to the list of processed articles. For the callback method for processing the 'text' fields, I used a bit of code from the pre-existing Wikiextractor project to remove Wiki markup language (such as links to other articles being displayed with square brackets and the like) from the content. This code is included in section A.1.5. The remainder of the parser is my work, and is included in section A.1.1. Throughout this process, the parser keeps a list of unique words encountered as well as outgoing link information for each article. These lists, along with the article contents, are saved to files each time a set number of articles have been processed. The link information is also kept in a hashmap with target articles as keys and a set articles linking to the target as values. The point of this is to reduce the computational complexity of the link processing as detailed in the following section.

### 1.2.2   Index Generation

The next step reads in the link information previously saved and adds it as ingoing link information in the respective article content files. The point of this approach is that link information is initially saved as a hashmap so the link going to a given article can be found quickly, rather that having

to search for outgoing links in every other article to determine the ingoing links to each article, which would be of $O(n^2)$ complexity.

Following that, articles with sufficiently few words and/or ingoing/outgoing links are discarded an index lists for the remaining articles and words are generated to associate a unique row/column number with each word/concept pair. The code performing the step described here is included in section A.1.2.

### 1.2.3 Matrix Construction

This final step converts the information compiled in the previous steps into a very large sparse matrix. The program allows for this to be done in 'chunks' in order to avoid insane RAM usage. Similarly, the matrix is stored in segments with each file containing a set number of rows in order to avoid loading the entire matrix to interpret a short text.

The full matrix is initially constructed using a DOK (dictionary of keys) sparse format in which the $i, j$th element simply counts the number of occurrences of word $i$ in the article corresponding to concept $j$. This is denoted $\text{count}(w_i, c_j)$. The DOK format as a hashmap using tuples $(i, j)$ as keys and the corresponding matrix elements as values and is the fastest format available for element-wise construction. The matrix is subsequently converted to CSR (compressed sparse row) format, which allows faster operations on rows which performs much quicker when computing TF-IDF (term frequency - inverse document frequency) scores and extracting concept vectors from words, i.e. when accessing separate rows corresponding to certain words.

Each non-zero entry is then converted to a TF-IDF score according to

$$T_{ij} = \left(1 + \ln\left(\text{count}(w_i, c_j)\right)\right) \ln\left(\frac{n_c}{df_i}\right), \tag{1.6}$$

where $n_c$ is the total number of concepts and

$$df_i = |\{c_k, w_i \in c_k\}| \tag{1.7}$$

is the number of concepts whose corresponding article contains the $i$th word. Thus, the first part of (1.6), $1 + \ln\left(\text{count}(w_i, c_j)\right)$ is the *text frequency* term, as it increases with the frequency of word $i$ in document $j$. Similarly, $\ln\left(\frac{n_c}{df_i}\right)$ in (1.6) is the *inverse document frequency* term as it decreases with the frequency of documents containing word $i$. Thus, the TF-IDF score as somewhat complement to entropy in that it goes to zero as the fraction

Giver det mening?

of documents containing word $i$ goes to 1, and takes its highest values if word $i$ occurs with high frequency in only a few documents[12]. While (1.6) is not the only expression to have those properties, empirically it tends to achieve superior results in information retrieval[16].

Each row is then $L^2$ normalized (divided by their Euclidean norm):

$$T_{ij} \rightarrow \frac{T_{ij}}{\sqrt{\sum_i T_{ij}^2}}. \tag{1.8}$$

Finally, each row is pruned to reduce spurious associations between concepts and articles with a somewhat uniform occurrence rate. This was done in practice by following the pragmatic approach of Gabrilovich [8] of sorting the entries of each row, move a sliding window across the entries, truncating when the falloff drops below a set threshold and finally reversing the sorting. The result of this step is the matrix which computes the interpretation vectors as described in 1.1.2. The code is included in section A.1.3.

Hearst nævner noget offentligt tilgængeligt Reuters-data som folk øver tekstklassifikation på. Det kunne være ret sjovt. Det kunne være sjovt at lave 'semantisk nearest neighbor'

## 1.3   Applications & Results

Having constructed a necessary machinery, I wrote a small Python module to provide an easy-to-use interface with the output from the computations described earlier. The code for this is included in section A.1.4. The module consists mainly of a SemanticAnalyser class, which loads in the previously mentioned index lists and provides methods for various computations such as estimating the most relevant concepts for a text, determining semantic distance etc. For example, the following code will create a semantic analyser instance and use it to guess the topic of the input string:

```
sa = SemanticAnalyser()
sa.interpret_text("Physicist from Austria known for the theory of relativity")
```

This returns a sorted list of the basis concepts best matching the input string, where the first element is of course 'Albert Einstein'. The SemanticAnalyser class contains equally simple methods to interpret a target text file or keyboard input, to calculate the semantic similarity or cosine distance between texts, and to compute interpretations vectors from a text.

The same module contains a TweetHarvester class which I wrote in order to obtain a large number of tweets to test the semantic analyser on, as tweets are both numerous and timestamped, which allows investigations

of the temporal evolution of tweets matching a given search term. The
TweetHarvester class provides an equally simply interface - for instances,
the 100 most recent tweets regarding a list of companies can be mined and
printed by typing

```python
terms = ['google', 'carlsberg', 'starbucks']
th = TweetHarvester()
th.mine(terms, 100)
print th.harvested_tweets
```

in addition to actively 'mining' for tweets matching a given query, the
class can also passively 'listen' for tweets while automatically saving its
held tweets to a generated date-specific filename once a set limit of held
tweets is exceeded:

```python
th = TweetHarvester(tweets_pr_file = 100)
th.listen(terms)
```

The downloaded tweets are stored as tweet objects which contain a built in
method to convert to a JSON-serializable hashmap, an example of which is
provided in 1.2. As can be seen in the example, the tweet object contains
not only the tweets textual content but also a wide range of metadata such
as the hashtags contained in the tweet, users mentioned, time of creation,
language etc. Indsæt fine grafer OSV!!! DO IT NAW!

```
{u'contributors': None,
 u'coordinates': None,
 u'created_at': u'Wed Jun 03 12:16:23 +0000 2015',
 u'entities': {u'hashtags': [],
               u'symbols': [],
               u'urls': [],
               u'user_mentions': [{u'id': 630908729,
                                   u'id_str': u'630908729',
                                   u'indices': [0, 12],
                                   u'name': u'Alexandra White ',
                                   u'screen_name': u'lexywhite86'}]},
 u'favorite_count': 0,
 u'favorited': False,
 u'geo': None,
 u'id': 606071930282745857L,
 u'id_str': u'606071930282745857',
 u'in_reply_to_screen_name': u'lexywhite86',
 u'in_reply_to_status_id': 605991263129714688L,
 u'in_reply_to_status_id_str': u'605991263129714688',
 u'in_reply_to_user_id': 630908729,
 u'in_reply_to_user_id_str': u'630908729',
 u'is_quote_status': False,
 u'lang': u'en',
 u'metadata': {u'iso_language_code': u'en', u'result_type': u'recent'},
 u'place': None,
 u'retweet_count': 0,
 u'retweeted': False,
 u'source': u'<a href="http://twitter.com/download/iphone" rel="nofollow">Twitter for iPhone
     </a>',
 u'text': u'@lexywhite86 if carlsberg did mornings \U0001f602',
 u'truncated': False,
 u'user': {u'contributors_enabled': False,
           u'created_at': u'Thu Apr 05 13:48:00 +0000 2012',
           u'description': u'',
           u'favourites_count': 169,
           u'follow_request_sent': False,
           u'followers_count': 110,
           u'friends_count': 268,
           u'geo_enabled': False,
           u'id': 545986280,
           u'id_str': u'545986280',
           u'lang': u'en',
           u'listed_count': 0,
           u'location': u'',
           u'name': u'Robert Murphy',
           u'notifications': False,
           u'protected': False,
           u'screen_name': u'7Robmurphy',
           u'statuses_count': 1650,
           u'time_zone': None,
           u'url': None,
           u'utc_offset': None,
           u'verified': False}}
```

Snippet 1.2: Example of a downloaded tweet.

Part II

SOCIAL FABRIC PROJECT

# 2

# Phone Activity and Quantitative Data Extraction

He main objective of part II of this thesis is to investigate behavioural patterns in the phone activities of the participants in the Social Fabric Project, and to predict various traits of the users based only on their phone logs. Throughout the part, I'll provide brief examples of usage for the software I've written to process the large amount of data available and to apply various prediction schemes to it, while the source code itself is included in the appendix.

My first objective was to investigate how various phone activities correlate with each other temporally, i.e. how a given user's probability for e.g. receiving a call increases or decreases around other activities such as moving around physically. This is the topic of section 2.1.

Next, I set out to replicate some recent research results claiming that people's phone activities predict certain psychological traits. In the most general terms, then, the task consists of predicting a collection of numbers or labels denoted $Y$ based on a set of corresponding data points $X$. The topic of section 2.2 is the extraction of the many-dimensional data points or *feature vectors* $X$ from the phone logs of the participants, while section 2.3 gives a brief description of the psychological traits $Y$. Finally, section 2.4 derives an often used linear classification method known as Linear Discriminant Analysis or Fisher's Discriminant and provides a discussion of why it fails for the present dataset, which serves to motivate the more sophisticated prediction schemes introduced in chapter 3.

kilder!!!

15

## 2.1   Temporal Correlations in Activity

One category of interesting quantities is the predictability of mobile phone behaviour from recorded behaviour at different times, i.e. the influence of certain events deduced from a user's Bluetooth, GPS or call log data on the tendency of some event to happen in the near past or future. A simple example would be to determine how much placing or receiving a call increases or decreases the probability of a user placing or receiving another call in the period following the first call.

This analysis was performed by comparing an 'activity' signal with a 'background' signal in the following fashion: For each user, the time period around each call is sliced into bins and the each of the remaining calls placed in the bin corresponding to the time between the two calls. Once divided by the total number of calls for the user, this is the activity signal. The background is obtained in a similar fashion but comparing each call in a given user's file with calls in the remaining users' files.

This involves repeatedly binning the time around certain events and then determining in which bin to place other events; a situation in which confusion may arise easily and errors may be hard to identify. To accommodate this, I started out by writing a custom array class designed to greatly simplify the binning procedure. This class called features the following:

- Methods to bin the time around a given event and determine determine which bin a given event falls into. This is useful to implement in the class itself as one then avoids having to continually worry about which bin an event fits into, and as it ensures that bin placement errors can only arise in one small piece of code which can then be tested rigorously.

- Attributes that keep track of the number of events that didn't fit into any bin, and of the current centre of the array, which can then be manipulated to move the array and a method to use this to return a normalized list of the bins.

In short, the binarray can be visualized as a collection of enumerated buckets that can be moved so as to center it on some event and then let other events 'drip' into the buckets. The code for this class is included in A.2.1. In general, objects can be converted to byte streams and stored using Python's pickle module, but as that tends to be both slow and insecure, I generally used json to save my objects. This poses a slight problem as some data types, such as tuples, and custom classes in general are not json

serializable. I got around this by writing some recursive helper methods to help store the relevant information about arbitrary nested combinations of some such objects and to help reconstruct said objects again. These are also included in section A.2.1. As an example of usage, the following code constructs a Binarray, centers it around the present time, and generates a number of timestamps which are then placed in the event. It is then saved to a file using the helper method previously described.

```python
from time import time
from random import randint
#Create Binarray with interval +/− one hour and bin size ten minutes.
ba = Binarray(interval = 60*60, bin_size = 10*60)
#Center it on the present
now = int(time())
ba.center = now
#Generate some timestamps around the present
new_times = [now + randint(−60*60, 60*60) for _ in xrange(100)]
for tt in new_times:
    ba.place_event(tt)

#Save it
with open('filename.sig', 'w') as f:
    dump(ba, f)
```

This data will be visualized by plotting the relative signal from the activity of some event, such as in- or outgoing calls or texts, over the background (simply $A/B$) around another type of event hypothesized to trigger the activity.

### 2.1.1 Influence of phone calls

I first investigated the effects of incoming and outgoing calls as triggers for other phone activities. The call logs were stored in a format where each line represents a hashmap with quantities such as call time or call duration mapping to their corresponding value. Below is an example of one such line, were any personal data have been replaced by a random number or hexadecimal string of equal length.

```json
{
    "timestamp": 6335212287,
    "number": "c4bdd708b1d7b82e349780ee1e7875caa600c579",
    "user": "ea42a1dbe422f83b0178d158f154f4",
    "duration": 483,
    "type": 2,
    "id": 45687
}
```

the text logs are similar except for the missing duration entry. Computing the relative signal in Binarrays centered on each incoming and outgoing

(a) Relative activity of events triggered by incoming calls.



(b) Relative activity of events triggered by outgoing calls.

Figure 2.1: Comparison of the increased activity caused by incoming and outgoing calls over an interval of ± 3 hours around an event with bins of three minutes.

call using bin sizes of three and thirty minutes resulted in the plots shown in figures 2.1 and 2.2, respectively. As the figures clearly show, the all four activities increase significantly for the average user around incoming and outgoing calls.

(a) Relative activity of events triggered by incoming calls.



(b) Relative activity of events triggered by outgoing calls.

Figure 2.2: Comparison of the increased activity caused by incoming and outgoing calls over an interval of ± 12 hours around an event with bins of thirty minutes.

## 2.1.2   Influence of GPS activity

The raw format of the users' GPS logs looks similar to those of the call
and text logs:

```
{
    "timestamp": 8058876274,
    "lon": 6.45051654,
    "user": "0c28e8f4ad9619bca1e5ea4167e10a",
    "provider": "gps",
    "lat": 28.20527041,
    "id": 6429902,
    "accuracy": 39.4
}
```

An analysis similar to that of described in section 2.1.1 was carried out using
GPS phone data as triggers. I chose to define a user as being 'active' if they
travelled at an average speed of $0.5\,\mathrm{m/s}$ between two consecutive GPS log
entries, while discarding measurements closely following each other. The
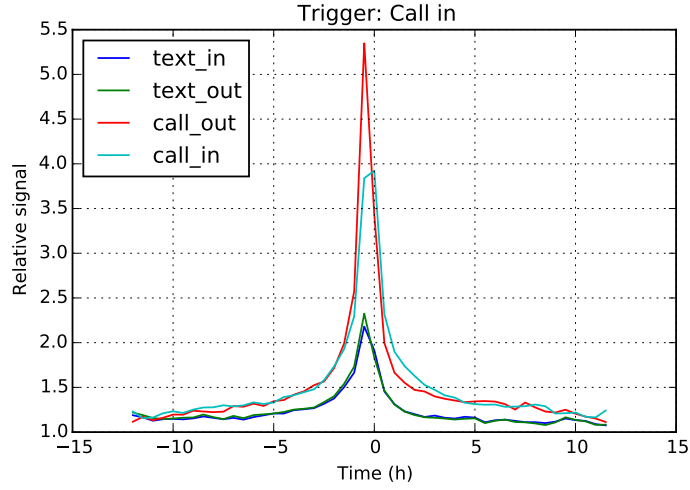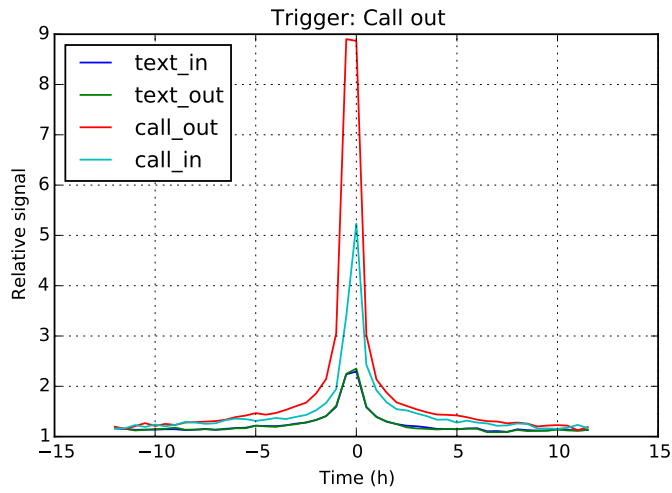reason for this is that the uncertainty on the location measurements could
yield false measurements of high average speeds when the measurements
are not temporally separated. A lot of the measurements turned out to
be grouped somewhat tightly - for instance, approximately 80% of the
time intervals were below $100\,\mathrm{s}$. This occurs because the Social Fabric
data not only actively records its users' locations with some set interval,
but also passively records the location when another app requests it, so
when users spend time on apps that need to continually update their
position such as Google Maps, a location log entry is written every second.
The distribution of intervals between consecutive GPS measurements is
shown in figure 2.3. A typical uncertainty on civilian GPS devices is at
most $100\,\mathrm{m}$[18], so because I choose to consider a user active if they travel
at a mean speed of $0.5\,\mathrm{m/s}$, and based on the time spacings shown in figure
2.3, I chose to discard measurements separated by less than $500\,\mathrm{s}$.

An analysis like that of section 2.1.1 reveals that a user's phone activity
is significantly increased around times when they are on the move, as
shown in figure 2.4. Note the asymmetry of the signal, especially visible
in figure 2.4(a). After a measurement of a user being active, the signal
dies off about two and a half hours into the future, whereas it persists
much longer into the past. Concretely, this means that people's phone
activity (tendency to call or text) becomes uncorrelated with their physical
activity after roughly two and a half hours, whereas their tendency to
move around is increased for much longer time after calling or texting.

The relative signal in figure 2.4(b) appears to be increasing at around
$\pm24\,\mathrm{h}$, which would seem reasonable assuming people have slightly

Figure 2.3: Plot of typical temporal spacings between consecutive GPS measurements.

different sleep schedules - if a person is on the move and hence more likely to place a call at time $t = 0$, they're slightly more likely than the general user to be on the move around $t = \pm 24$ h. Figure 2.5 shows the same signal extended to $\pm 36$ h where slight bumps are visible 24 hours before and after activity.

(a) Interval: 5 hours. Bin size: 5 minutes.



(b) Interval: 24 hours. Bin size: 15 minutes.

Figure 2.4: Relative increase of activities triggered by GPS activity.

Figure 2.5: GPS-triggered activity increase over an interval of 36 hours using a bin size of 10 minutes.

## 2.1.3   Influence of Bluetooth signal

The following is a randomized entry in a user's bluetooth log.

```
{
    "name": "d5306a3672b7a0b8f9696d294ec4b731",
    "timestamp": 6870156680,
    "bt_mac": "1f158ae269d69efa5bb4794ee2a0b2dd68bd3a9badfeaf70f258ad3c74b0c09b",
    "class": 1317046,
    "user": "41cdb7ecaaaec3d33391ed063e7fa2",
    "rssi": -76,
    "id": 4139043
}
```

The 'bt_mac' entry is the MAC-adress of the device which the Bluetooth receiver in the user's phone has registered, so it is reasonable to assume several different MAC addresses occur at several consecutive timestamps. I call the number of repeated MAC adresses needed for a user to be considered social the 'social threshold'. Figures 2.6, 2.7 and 2.8 show the increased activity around times when users were considered social with a threshold of 1, 2 and 4 repeated pings.

Contrary to the previous analyses, phone activities decreased somewhat when users were social. As stated, each of these analyses were fairly similar, I've only explicitly included the code used to extract and save

(a) Interval: 12 hours, Bin size: 10 minutes.



(b) Interval: 36 hours, Bin size: 15 minutes.

Figure 2.6: The effect on phone activity of sociality as measured by the user's Bluetooth signal. The threshold used for being considered social as one repeated signal.

(a)  Interval: 12 hours, Bin size: 10 minutes.



(b)  Interval: 36 hours, Bin size: 15 minutes.

Figure 2.7: The effect on phone activity of sociality as measured by the user's Bluetooth signal. The threshold used for being considered social as two repeated signals.

(a) Interval: 12 hours, Bin size: 10 minutes.



(b) Interval: 36 hours, Bin size: 15 minutes.

Figure 2.8: The effect on phone activity of sociality as measured by the user's Bluetooth signal. The threshold used for being considered social as four repeated signals.

Bluetooth data, as well as the code used to load the data and generate figures 2.6 through 2.8. This code is included in section A.2.2.

## 2.2 Extraction of Input Data

The predictive powers of mobile phone behaviour on the user's psychological profile is currently an area of active research. As part of my thesis work, I have tried to predict the psychological profiles of the SFP participants using various machine learning methods on the available phone logs.

1000 kilder!!!

The software I've written first preprocesses the phone logs to extract various relevant parameters, then collects the parameters and psychological profile scores for each user to serve as input and output, respectively, for the various learning methods. Many of the parameters are chosen following a recent article by de Montjoye et al[5]. The following contains an outline and brief explanation of the extracted parameters.

This section contains a list of the extracted parameters used for psychological profiling along with a brief description of the extraction process when necessary. The preprocessing code is included in section A.2.3.

### 2.2.1 Simple Call/Text Data

The most straightforward data type is the timestamps from a given user's call/text logs. Six of the parameters used were simply the standard deviation and median of the times between events in the logs for each user's call log, text log, and the combination thereof, excluding time gaps of more than three says on the assumption that it would indicate a user being on vacation or otherwise having a period of telephone inactivity. The entropy $S_u$ of each of the three was also included simply by computing the sum

$$S_u = -\sum_c \frac{n_c}{n_t} \ln_2 \frac{n_c}{n_t}, \tag{2.1}$$

where $c$ denotes a given contact and $n_t$ the total number of interactions, and $n_c$ the number of interactions with the given contact. The number of contacts, i.e. the number of unique phone numbers a given user had contacted by means of calls, texts, and the combination thereof, was also extracted along with the total number of the various kinds of interactions and the contact to interaction ratios. The response rates, defined as the rate of missed calls and incoming texts, respectively, that a given user replied to within an hour, where also determined along with the text

latency defined as the median test response time. Finally the percentage calls and texts that were outgoing was determined as well as the fraction of call interactions that took places during the night, defined as between 22-08.

### 2.2.2   Location Data

A number of parameters based on the spacial dynamics of the user were also extracted. Among these is the radius of gyration, meaning simply the radius of the smallest enclosing circle enclosing all the registered locations of the user on the given day, and the distance travelled per day. I chose to extract the median and standard deviation of each, filtering out the radii that exceeded 500km so as to keep information about long distance travels in the distance parameter and information about travel within a given region in the radius of gyration parameter.

#### Cluster Analysis

One parameter which has strong links[5] to psychological traits is the number of locations in which the user typically spends time, and the entropy of their visits to that location. Hence, the task at hand is to identify dense clusters of GPS coordinates for each user. This is a typical example of a task which is very intuitive and quickly approximated by humans, but is extremely computationally expensive to solve exactly. Concretely, the problem of finding the optimal division of $n$ data points into $K$ clusters is formulated as minimizing the 'score' defined as

$$S = \sum_{K} \sum_{x_n \in C_k} |x_n - c_k|^2, \tag{2.2}$$

where $c_k$ denotes the centroid of the cluster $C_k$. Each point $x_n$ is assigned to the cluster corresponding to the nearest centroid. The usual way of approaching this problem is to use Lloyd's algorithm, which consists of initializing the centroids randomly, assigning each point to the cluster corresponding to the centroid which is nearest, then moving each centroid to the center of its points and repeating the last two steps until convergence. As this isn't guaranteed to converge on the global minimum of (2.2), the process can be repeated a number of times, keeping only the result with the lowest value of $S$. I accomplished this by writing a small Python module to perform various variations of Lloyd's algorithm and to produce plots of the resulting clusters. The code is included in section A.2.5.

   This allows one to implement Lloyd's algorithm and visualize its result easily, as the code allows automatic plotting of the result from the

algorithm while automatically selecting different colors for the various clusters. As an example, the following code snippet generates 1000 random points, runs Lloyd's algorithm to determine clusters and saves a plot of the results.

```
points = [[random.uniform(−10,10), random.uniform(−10,10)] for _ in xrange(10**3)]
clusters = lloyds(X = points, K = 6, runs = 1)
draw_clusters(clusters = clusters, filename = 'lloyds_example.pdf')
```

This results in the following visualization:



I chose to modify the algorithm slightly on the following basis: Usually, the algorithm takes as its initial centroids a random sample of the data. I'll call this 'sample' initialization. This leads to a greater number of clusters being initialized in the areas with an increased density of data points, meaning that centroids will be highly cluttered at first, 'fighting' over the dense regions of data points then slowly spreading out. A few such iterations are shown in figure 2.9. However, this method is dangerous: The goal is to identify locations in which a user spends much of their time, i.e. in which more than some threshold of their GPS pings originated, and this initialization is likely to 'cut' the more popular locations into several clusters, neither of which contains more data points than the threshold. One example might be the DTU campus, which is a risk of being divided into several locations with too few data points in each, giving the false impression that user doesn't visit the campus that often. To avoid this effect, I implemented another initialization, 'scatter', in which the clusters start out on points select randomly from the entire range of $x, y$-values in the user's dataset. This turned out to not only solve the problem described above, but also converge much quicker and reach a slightly lower score as define in (2.2). A few such iterations are shown in figure 2.10. The difference in end results for the two methods is exemplified in figure 2.11. While this works great for users who stay in or around Copenhagen, it will cause problems for people who travel a lot. A user who has visited Australia, for instance, will have their initial clusters spread out across

Figure 2.9: A few iterations of Lloyd's algorithm using 'sample' initialization. The axes denote the distance in km to some typical location for the user. Note that clusters are initially cluttered, then slowly creep away from the denser regions.

the globe, and it's highly likely that one them will end up representing all of Denmark. I ended up simply running both versions and keeping the result yielding the highest amount of locations.

Figure 2.10: A few iterations of Lloyd's algorithm using 'scatter' initialization. The axes denote the distance in km to some typical location for the user. Note that clusters are initially randomly spread across the entire range of $x, y$-values and converge quickly to a local minimum for (2.2).

(a) Sample initialization.



(b) Scattered initialization.

Figure 2.11: Comparison of the final results of the two initialization methods using 100 initial clusters, a threshold of 5% of the data points before a cluster is considered a popular location and running the algorithm 10 times and keeping the best result. Clusters containing more than 5% of the total amount of data points are in color, whereas the remaining points are black dots.

### 2.2.3 Time Series Analysis

Another interesting aspect to include is what one somewhat qualitatively might call behavioural regularity - some measure of the degree in which a user's phone activities follow a regular pattern. Quantifying this turns out to take a bit of work. First of all, any user's activity would be expected to closely follow the time of day, so the timestamps of each user's outgoing texts and calls are first converted into 'clock times' meaning simply the time a regular clock in Copenhagen's time zone would display at the given time. This process is fairly painless when using e.g. the UTC time standard, which does not observe daylight saving time (DST), but some subtleties arise in countries that do use DST, as this makes the map from Unix/epoch time to clock time 'almost bijective' - when changing *away* from DST, two consecutive hours of unix time map to the same clock time period (02:00 to 03:00), whereas that same clock period is skipped when changing *to* DST. The most commonly used Python libraries for datetime arithmetic accommodate this by including a dst boolean in their datetime objects when ambiguity might arise, however I simply mapped the timestamps to clock times and ignored the fact twice a year, one time bin will artificially contain contributions from one hour too many or few. One resulting histogram is shown in figure 2.12.

Tilføj lidt om AR-serier når du har bogen!!!

### 2.2.4 Facebook Data

Unfortunately, the only available Facebook data was a list of each user's friends, so the only contribution of each user's Facebook log was the number of friends the user had.

### 2.2.5 Bluetooth Data

I extracted a number of different features from each user's Bluetooth log file. First, I set a threshold for when a given user is considered social, as described in section 2.1.3. I chose to use a threshold of two. I then tried to estimate how much time each user spends in the physical company of others in the following way: for each time stamp in the user's Bluetooth log, I checked if the user was social or not and assumed that this status was the same until the following log entry, unless the delay was more than two hours. The rationale behind this is to avoid skewing the measurements if a user turns off their phone for extended periods of time. Otherwise, e.g. studying with a few friends at DTU, turning off your phone and going on vacation for two weeks would give the false impression that the user

Figure 2.12: Histogram of a user's outgoing calls and texts with a bin size of six hours.

were highly social for a long period of time. I then recorded the fraction of times the user was estimated as being social in this fashion.

Finally, I also wanted some measure of the degrees to which a user's social behaviour follows a pattern. I looked for temporal patterns by fitting AR-series and computing autocorrelation coefficients for each user's social behaviour as described in section 2.2.3. I also chose to compute a 'social entropy' much like (2.1), but weighted by the time the user spends with each acquaintance:

$$E = -\sum_i f_i \ln_2 (f_i), \tag{2.3}$$

$$f_i = \frac{\text{time spent with } i}{\sum_j \text{time spent with } j}. \tag{2.4}$$

Note that the denominator of (2.4) is not equal to the total amount of time spent being social, as the contribution from each log entry is multiplied by the number of people present.

## 2.3 Output Data

The main emphasis of this part of the thesis is on predicting so-called *Big Five* personality traits. This section contains a brief description of those, following[6]. **Extraversion** signifies how extroverted and sociable a person is. People with high extraversion scores are supposed to be more eager to seek the company of others. **Agreeableness** is supposed to be a measure of how sympathetic or cooperative a person is, whereas **conscientiousness** denotes constraint, self discipline, level of organization etc.. **Neuroticism** signify the tendency to experience mood swings, and is complementary to emotional stability. Finally, **Openness**, also called 'openness to experience', or 'inquiring intellect' in earlier works, signifies thoughtfulness, imagination and so on. These five are collectively referred to as the 'big five' or 'OCEAN' after their initials.

I addition to the above, I also had access to a range self-explanatory traits about the participants such as their gender, whether they smoke etc.

## 2.4 Linear Discriminant Analysis & Gender Prediction

Linear discriminant analysis is basically a dimensionality reduction technique developed by Fisher in 1936 [7] for separating data points into two or more classes. The general idea is to project a collection of data points in $n$-dimensional variable space, onto the line or hyperplane which maximizes the separation between classes. Representing data points in $n$-space by vectors denoted $x$, the objective is to find a vector $\omega$ such that separation between the projected data points on it

$$y = \omega^T x \tag{2.5}$$

is maximized.

To break down the derivation of this method, I will first define a convenient distance measure used to optimize the separation between classes, then solve the resulting optimization problem. For clarity, I'll only describe the case of projection of two classes onto one dimension (i.e. using 'line' rather than 'hyperplane' and so on), although the method generalizes easily.

### 2.4.1 A measure of separation for projected Gaussians

If the projected data points for two classes $a$ and $b$ follow distributions $\mathcal{N}_a$ and $\mathcal{N}_b$, which are standard Gaussians, $\mathcal{N}_i(x) = \mathcal{N}(x; \mu_i, \sigma_i^2)$, the joint probability distribution for the distance between the projections will be the convolution

$$P(x) = \int_{-\infty}^{\infty} \mathcal{N}_a(y) \cdot \mathcal{N}_b(x - y) \, dy. \tag{2.6}$$

Computing this for a Gaussian distribution,

$$\mathcal{N}(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \tag{2.7}$$

becomes easier with the convolution theorem, which I'll derive in the following.

Denoting convolution by $*$ and Fourier transforms by

$$\mathcal{F}(f) = \frac{1}{(2\pi)^{n/2}} \int_{\mathbb{R}^n} f(x) \cdot e^{-i\omega x} \, dx, \tag{2.8}$$

the convolution theorem is derived as follows:

$$\mathcal{F}(f * g) = \frac{1}{(2\pi)^{n/2}} \int_{\mathbb{R}^n} \int_{\mathbb{R}^n} f(y) \cdot g(x - y) \, dy \, e^{-i\omega x} \, d\omega, \tag{2.9}$$

$$= \frac{1}{(2\pi)^{n/2}} \int_{\mathbb{R}^n} f(y) \int_{\mathbb{R}^n} g(x - y) e^{-i\omega x} \, dy \, d\omega, \tag{2.10}$$

$$= \frac{1}{(2\pi)^{n/2}} \int_{\mathbb{R}^n} f(y) \int_{\mathbb{R}^n} g(z) e^{-i\omega(z+y)} \, dz \, d\omega, \tag{2.11}$$

$$= \frac{1}{(2\pi)^{n/2}} \int_{\mathbb{R}^n} f(y) e^{-i\omega y} \int_{\mathbb{R}^n} g(z) e^{-i\omega z} \, dz \, d\omega, \tag{2.12}$$

$$\boxed{\mathcal{F}(f * g) = (2\pi)^{n/2} \mathcal{F}(f) \cdot \mathcal{F}(g),} \tag{2.13}$$

where the factor in front of the usual form of the theorem $\mathcal{F}(f * g) = \mathcal{F}(f) \cdot \mathcal{F}(g)$ stems from the convention of using angular frequency in Fourier transforms, as in (2.8), rather than

$$\mathcal{F}(f) = \int_{\mathbb{R}^n} f(x) \cdot e^{-2\pi i v x} \, dx. \tag{2.14}$$

Using this, the convolution of two Gaussians can be calculated as

$$\mathcal{N}_a * \mathcal{N}_b = (2\pi)^{n/2} \mathcal{F}^{-1} \left( \mathcal{F}(\mathcal{N}_a) \cdot \mathcal{F}(\mathcal{N}_b) \right). \tag{2.15}$$

The required Fourier transform can be massaged into a nicer form by displacing the coordinate system and cancelling out terms with odd parity:

$$
\begin{aligned}
\mathcal{F}(\mathcal{N}(x)) &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \cdot e^{-i\omega x}\, \mathrm{d}x, \\
&= \frac{1}{2\pi\sigma} \int_{-\infty}^{\infty} e^{-\frac{x^2}{2\sigma^2}} e^{-i\omega(x+\mu)}\, \mathrm{d}x, \\
&= \frac{1}{2\pi\sigma} e^{-i\omega\mu} \int_{-\infty}^{\infty} e^{-\frac{x^2}{2\sigma^2}} (\cos(\omega x) + i\sin(\omega x))\, \mathrm{d}x, \\
&= \underbrace{\frac{1}{2\pi\sigma} e^{-i\omega\mu}}_{a} \underbrace{\int_{-\infty}^{\infty} e^{-\frac{x^2}{2\sigma^2}} \cos(\omega x)\, \mathrm{d}x}_{I(\omega)}.
\end{aligned}
\tag{2.16}
$$

Noting that $I(\omega)$ reduces to an ordinary Gaussian integral at $\omega = 0$ so $I(0) = \sqrt{2\pi}\sigma$, this can be solved with a cute application of Feynman's trick:

$$
\begin{aligned}
\frac{\partial I}{\partial \omega} &= -\int_{-\infty}^{\infty} x e^{-\frac{x^2}{2\sigma^2}} \sin(\omega x)\, \mathrm{d}x, \\
&= \int_{-\infty}^{\infty} \sigma^2 \frac{\partial}{\partial x}\left(e^{-\frac{x^2}{2\sigma^2}}\right) \sin(\omega x)\, \mathrm{d}x, \\
&= \sigma^2 e^{-\frac{x^2}{2\sigma^2}} \sin(\omega x)\Big|_{-\infty}^{\infty} - \omega \int_{-\infty}^{\infty} \sigma^2 e^{-\frac{x^2}{2\sigma^2}} \cos(\omega x)\, \mathrm{d}x, \\
&= -\omega\sigma^2 I(\omega) \Leftrightarrow \\
I(\omega) &= C e^{-\sigma^2\omega^2/2}, \\
I(0) = C &= \sqrt{2\pi}\sigma, \\
I(\omega) &= \sqrt{2\pi}\sigma e^{-\sigma^2\omega^2/2}.
\end{aligned}
$$

Plugging this into (2.16) gives the result

$$
\mathcal{F}(\mathcal{N}) = \frac{1}{\sqrt{2\pi}} e^{-i\omega\mu} e^{-\sigma^2\omega^2/2}.
\tag{2.17}
$$

This can be used in conjunction with (2.13) to obtain

$$
\begin{aligned}
\mathcal{F}(\mathcal{N}_a * \mathcal{N}_b) &= \sqrt{2\pi} \frac{1}{\sqrt{2\pi}} e^{-i\omega\mu_a} e^{-\sigma_a^2\omega^2/2} \cdot \frac{1}{\sqrt{2\pi}} e^{-i\omega\mu_b} e^{-\sigma_b^2\omega^2/2}, \tag{2.18} \\
&= \frac{1}{\sqrt{2\pi}} e^{-i\omega(\mu_a-\mu_b)} e^{-(\sigma_a^2+\sigma_b^2)\omega^2/2}, \tag{2.19}
\end{aligned}
$$

which is recognized as the transform of another Gaussian describing the separation with $\mu_s = \mu_a - \mu_b$ and $\sigma_s^2 = \sigma_a^2 + \sigma_b^2$, so taking the inverse Fourier transformation gives the convolution

$$\mathcal{N}_a * \mathcal{N}_b = \frac{1}{\sqrt{2\pi}\sigma_s} e^{-\frac{(x-\mu_s)^2}{2\sigma_s^2}} . \tag{2.20}$$

Hence, a reasonable measure of the separation of two projected distributions is

$$d = \frac{(\mu_a - \mu_b)^2}{\sigma_a^2 + \sigma_b^2}. \tag{2.21}$$

### 2.4.2  Optimizing separation

To maximize the separation, the numerator and denominator, respectively, of (2.21) can be rewritten in terms of $w$ in the following way (using $\widetilde{\mu}_i$ to denote projected means) and simplified by introducing scattering matrices:

$$(\widetilde{\mu}_a - \widetilde{\mu}_b)^2 = \left(w^T (\mu_a - \mu_b)\right)^2, \tag{2.22}$$

$$= w^T (\mu_a - \mu_b)(\mu_a - \mu_b)^T w, \tag{2.23}$$

$$= w^T S_B w, \tag{2.24}$$

and

$$\widetilde{\sigma}_i^2 = \sum_{y \in i} \frac{1}{N} (y - \widetilde{\mu}_i)^2, \tag{2.25}$$

$$= w^T \sum_{y \in i} (x - \mu_i)(x - \mu_i)^T w, \tag{2.26}$$

$$= w^T S_i w, \tag{2.27}$$

$$\widetilde{\sigma}_a^2 + \widetilde{\sigma}_b^2 = w^T S_W w, \tag{2.28}$$

having introduced the between-class and within-class scatter matrices $S_B$ and $S_W$ by

$$S_B = (\mu_a - \mu_b)(\mu_a - \mu_b)^T, \tag{2.29}$$

$$S_i = \sum_{y \in i} (x - \mu_i)(x - \mu_i)^T, \tag{2.30}$$

$$S_W = S_a + S_b. \tag{2.31}$$

Hence, the objective is to solve

$$\frac{\mathrm{d}}{\mathrm{d}w}J(w) = \frac{\mathrm{d}}{\mathrm{d}w}\left(\frac{w^T S_B w}{w^T S_W w}\right) = 0, \tag{2.32}$$

$$\frac{\frac{\mathrm{d}[w^T S_B w]}{\mathrm{d}w}w^T S_W w - w^T S_B w\frac{\mathrm{d}[w^T S_W w]}{\mathrm{d}w}}{(w^T S_W w)^2} = 0, \tag{2.33}$$

$$2S_B w \cdot w^T S_W w - w^T S_B w \cdot 2S_W w = 0, \tag{2.34}$$

$$S_B w - \frac{w^T S_B w \cdot S_W w}{w^T S_W w} = 0, \tag{2.35}$$

$$S_B w - S_W w J(w) = 0, \tag{2.36}$$

$$S_B w = S_W w J(w), \tag{2.37}$$

$$S_W^{-1} S_B w = J(w)w. \tag{2.38}$$

The optimal projection vector $w^*$ which satisfies this is

$$w^* = S_W^{-1}(\mu_a - \mu_b). \tag{2.39}$$

> Vær lige sikker på at du forstår det her.

Figure 2.13 shows a visualization of this that I generated by drawing $(x, y)$ points from two random distributions to simulate two distinct classes of points. If the distributions are independent and Gaussian, the projections will also form Gaussian distributions, and the probability of a new point belonging to e.g. class $a$ given its coordinates $d$ can be estimated using Bayesian probability

$$P(a|d) = \frac{P(d|a)P(a)}{P(d|a)p(a) + P(d|b)P(b)}, \tag{2.40}$$

where $P(a)$ and $P(b)$ are simply the prior probabilities for encountering the respective classes, and the conditional probabilities, e.g. $P(d|a)$ are simply given by the value of the projected Gaussian $\mathcal{N}(x'; \widetilde{\mu}_a, \widetilde{\sigma}_a)$ at the projected coordinate $x'$. In practise, even when the points are not independent or Gaussian, so that (2.40) is not a precise estimate of the probability of the point representing a given class, the class with the highest posteriori according to (2.40) still often turns out to be a good guess.

This method accurately predicted the gender of 79.8% of the participants, which is not particularly impressive as 77.3% of participants were male, so a classifier that assumes that every participant is male would have a comparable success rate. An immediate source of concern is the assumption of linearity: It is possible that the data is ordered in such a way that it is possible to separate data points fairly well based on

Figure 2.13: Two collections of points drawn from independent Gaussian distributions, representing class a and class b. If the points are projected onto the straight line, which is given by (2.39), the separation between the peaks representing the two classes is maximized.

gender or some psychological trait, just not using a linear classifier. As an extreme example of this, figure 2.14 shows a situation where the points representing one class are grouped together in an 'island' in the middle, isolating them from points representing the remaining class. While it is clear that there's a pattern here, a linear classifier fails to predict classes more precisely than their ratio. Support Vector Machines, or SVMs are another linear classification technique which can be generalized to detect patterns like that in figure 2.14. This is described in section 3.1

Figure 2.14: An example of data points representing class a are clearly discernible from those of class b, yet a linear Fisher classifier fails to predict the classes more precisely than the ratio of b to a.

# CHAPTER 3

# PSYCHOLOGICAL PROFILING & MACHINE LEARNING

ACHINE learning is currently a strong candidate for prediction of psychological profiles from phone data. This chapter describes the application of the quantitative data described in setion 2.2 and various machine learning schemes, starting with support vector machines (SVMs).

> 1000 kilder!!!

> Uddyb når der er flere modeller.

## 3.1 Support Vector Machines

The purpose of this section is to introduce SVMs and attempt to apply them to the data obtained in 2.2. The introduction is mainly based on introductory texts by Marti Hearst [10] and Christopher Burges [4]. SVMs in their simplest form (*simplest* meaning using a linear kernel, which I'll explain shortly) can be thought of as a slight variation on the linear classifier described in section 2.4. However, where LDA finds a line such that the distribution of the points representing various classes projected onto the line is maximized, the aim of SVMs is to establish the hyperplane that represents the best possible slicing of the feature space into regions containing only points corresponding to the different classes. A simple example of this is shown in figure 3.1. Using labels ±1 to denote classes, the problem may be stated as trying to guess the mapping from an N-dimensional data space to classes $f : \mathbb{R}^N \rightarrow \{\pm 1\}$ based on a set of training data in $\mathbb{R}^N \otimes \{\pm 1\}$. I'll describe separately the properties of

Figure 3.1: The same points as those shown in figure 2.13, except points in class a and class b are now pictured along with their maximally separating hyperplane.

this maximally separating hyperplane, how it is obtained, and how the method is generalized to non-linear classification problems as the 'island' illustrated in figure 2.14.

The well-known equation for a plane is obtained by requiring that its normal vector $\mathbf{w}$ be orthogonal to the vector from some point in the plane $\mathbf{p}$ to any point $\mathbf{x}$ contained in it:

$$\mathbf{w} \cdot (\mathbf{x} - \mathbf{p}) = 0. \tag{3.1}$$

The left hand side of (3.1) gives zero for points in the plane and positive or negative values when the point is displaced in the same or opposite direction as the normal vector, respectively. Hence, $\text{sign}\,(\mathbf{w} \cdot (\mathbf{x} - \mathbf{p}))$ may be taken as the decision function. It is clear from (3.1) that the normal vector may be scaled without changing the actual plane (of course the decision function is inverted if a negative value is chosen), so $\mathbf{w}$ is usually rescaled such that

$$\mathbf{w} \cdot (\mathbf{x} - \mathbf{p}) = \mathbf{w} \cdot \mathbf{x} + b = \pm 1, \tag{3.2}$$

for the points that are closest to the separating plane. Those points located on the margin are encircled in figure 3.1. In general then, the meaning of

the sign and magnitude of

$$\mathbf{w} \cdot \mathbf{x} + b \qquad (3.3)$$

will be the predicted class and a measure of prediction confidence, respectively, for new data points. Finally, note that $\mathbf{w}$ can be expanded in terms of the data points that are on the margin in figure 3.1 as

$$\mathbf{w} = \sum_i v_i \mathbf{x}_i, \qquad (3.4)$$

these $\mathbf{x}_i$, the position vectors of the margin points in data space, are the 'support vectors' that lend their name to the method.

### 3.1.1 Obtaining the Maximally Separating Hyperplane

Assuming first that it is possible to slice the data space into two regions that contain only points corresponding to one class each, and that the plane's normal vector has already been rescaled according to (3.2), the following inequalities hold:

$$\mathbf{x_i} \cdot \mathbf{w} + b \geq 1, y_i = +1,$$
$$\mathbf{x_i} \cdot \mathbf{w} + b \leq -1, y_i = -1. \qquad (3.5)$$

Multiplying by $y_i$, both simply become

$$y_i \left( \mathbf{x_i} \cdot \mathbf{w} + b \right) - 1 \geq 0. \qquad (3.6)$$

The separation between the two margins shown with dashed lines in figure 3.1 is 2/|$\mathbf{w}$|, so the Lagrangian

ikke 100% tryg her...

$$L = \frac{1}{2}|\mathbf{w}|^2 - \sum_i \alpha_i y_i \left( \mathbf{x}_i \cdot \mathbf{w} + b - 1 \right), \qquad (3.7)$$

must be minimized with the constraints

$$\alpha_i \geq 0, \qquad (3.8)$$

$$\frac{\partial L}{\partial \alpha_i} = 0. \qquad (3.9)$$

A result from convex optimization theory known as Wolfe Duality[17] states that one may instead maximize the above Lagrangian subject to

$$\nabla_w L = \frac{\partial L}{\partial b} = 0, \qquad (3.10)$$

which gives conditions

$$\mathbf{w} = \sum_j \alpha_j y_j \mathbf{x}_j, \tag{3.11}$$

$$\sum_j \alpha_j y_j = 0. \tag{3.12}$$

These can be plugged back into (3.7) to obtain

$$L_D = \frac{1}{2} \sum_i \sum_j \alpha_i y_i \alpha_j y_j \mathbf{x}_i \cdot \mathbf{x}_j - \sum_i \alpha_i y_i \left( \mathbf{x}_i \cdot \sum_j \alpha_j y_j \mathbf{x}_j + b \right) + \sum_i \alpha_i, \tag{3.13}$$

$$L_D = -\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \mathbf{x}_j + \sum_i \alpha_i. \tag{3.14}$$

A problem with this is that eqs 3.5 can only be satisfied in the completely separable case, although it is easy to imagine an example in which a classifier performs well but not flawlessly on the training set. For instance, if two points, one from each class, in figure 3.1 were permuted, the classifier shown in the plot would still do a very good job, but eqs. 3.5 would not be satisfiable, causing the method to fail. This is remedied by introducing slack variables

$$\mathbf{x_i} \cdot \mathbf{w} + b \geq 1 - \xi_i, \quad y_i = +1, \tag{3.15}$$

$$\mathbf{x_i} \cdot \mathbf{w} + b \leq -(1 - \xi_i), \quad y_i = -1, \tag{3.16}$$

$$\xi_i \geq 0, \tag{3.17}$$

The above can be rewritten exactly as previously, except another set of non-negative Lagrange multipliers are added to (3.7) to ensure positivity of the $\xi_i$. Also a cost term, usually $C \cdot \sum_i \xi_i$, is added to the Lagrangian to keep separability high:

$$L = \frac{1}{2} |\mathbf{w}|^2 + C \cdot \sum_i \xi_i - \sum_i \alpha_i y_i \, (\mathbf{x}_i \cdot \mathbf{w} + b - 1 + \xi_i) - \sum_i \mu_i \xi_i. \tag{3.18}$$

This results in the same dual Lagrangian $L_D$ as before, but with an upper bound on the $\alpha_i$:

$$0 \leq \alpha_i \leq C. \tag{3.19}$$

The values of the slack variables $\xi_i$ and the cost term $C$ are typically decided by performing a 'grid search' in which the performance is evaluated for each possibly combination of the parameters and the optimal combination

Figure 3.2: Heat map of the result of a grid search over the parameter space of the variables $\xi$ and $C$. The colors signify the mean squared error of a support vector regression problem, which is closely related to the classification problem described in the present section. Note how overfitting damages performance for high values of $C$.

used in the final classifier. Figure 3.2 shows a heat map of the results of a grid search over the parameters in question. The main point to be emphasized here is that the training data $\mathbf{x}_i$ only enter into the dual Lagrangian of (3.14) as inner products. This is essential when extending the SVM model to nonlinear cases, which is the subject of the following section.

### 3.1.2 Generalizing to the non-linear case

The fact that the data $\mathbf{x}_i$ only occur as inner products in (3.14) makes one way of generalizing to non-linearly separable datasets straightforward: Referring back to figure 2.14, one might imagine bending the plane containing the data points by curling the edges outwards in a third dimension after which a two-dimensional plane could separate the points very well. In general, this means applying some mapping

$$\Phi : \mathbb{R}^l \to \mathbb{R}^h, \quad h > l, \tag{3.20}$$

to the $\mathbf{x}_i$ ($l$ and $h$ are for low and high, respectively).  For example, one could look for a mapping such that the new inner product becomes

$$\Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j) = \left(\mathbf{x_i} \cdot \mathbf{x_j}\right)^2 . \tag{3.21}$$

I'll describe the components of each vector separately, so I'm going to change notation to let the subscripts denote coordinates and using $\mathbf{x}$ and $\mathbf{y}$ as two arbitrary feature vectors, where the latter shouldn't be confused with the class label used earlier.  As an example, in two dimensions the above becomes

$$(\mathbf{x} \cdot \mathbf{y})^2 = \left(\sum_{i=1}^{2} x_i y_i\right)^2 = x_1^2 y_1^2 + 2x_1 y_1 x_2 y_2 + x_2^2 y_2^2, \tag{3.22}$$

meaning that one possibility for $\Phi$ is

$$\Phi : \mathbf{x} \mapsto \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1 x_2 \\ x_2^2 \end{pmatrix} \tag{3.23}$$

This can be generalized to $d$-dimensional feature vectors and to taking the $n$'th power rather than the square using the multinomial theorem:

$$\left(\sum_{i=1}^{d} x_i\right)^n = \sum_{\sum_{i=1}^{d} k_i = n} \frac{n!}{\prod_{l=1}^{d} k_l!} \prod_{j=1}^{d} x_j^{k_j}, \tag{3.24}$$

where the subscript $\sum_{i=1}^{d} k_i = n$ simply means that the sum goes over any combination of $d$ non-negative integers $k_i$ that sum to $n$. I wish to rewrite this slightly for two reasons: to simplify the notation in order to make a later proof more manageable, and to help quantify how quickly the number of dimensions in the output space grows to motivate a trick to avoid these explicit mappings.

   As stated, the sum on the RHS of (3.24) runs over all combinations of $d$ integers which sum to $n$. This can be simplified by introducing a function $K$, which simply maps

Dobbeltjek om notation er fornuftig.

$$K : n, d \mapsto \left\{ \{k\} \in \mathbb{N}^d \,\middle|\, \sum_{i=1}^{d} k_i = n \right\}, \tag{3.25}$$

and denoting each of those collections $\{k\}_i$ so each of the coefficients in (3.24) can be written

$$\frac{n!}{\prod_{i=1}^{d} k_i!} = C_{\{k\}}. \tag{3.26}$$

Then, (3.24) becomes

$$\left(\sum_{i=1}^{d} x_i\right)^n = \sum_{K(n,d)} C_{\{k\}} \prod_{j=1}^{d} x_j^{k_j} \tag{3.27}$$

To show how quickly the dimensions of the required embedding space grows, note that the dimension is equal to the number of terms in the sum above, i.e.

$$\dim(\mathbb{R}^h) = |K(n,d)| = \left|\left\{ \{k\} \in \mathbb{N}^d \,\middle|\, \sum_{i=1}^{d} k_i = n \right\}\right|, \tag{3.28}$$

which can be computed using a nice trick known from enumerative combinatorics.

Consider the case where $n = 5$ and $d = 3$. $K(5,3)$ then contains all sets of 3 integers summing to 5, such as $1, 3, 1$ or $0, 1, 4$. Each of these can be uniquely visualized as 5 unit values distributed into 3 partitions in the following fashion:

$$\circ \mid \circ \ \circ \ \circ \mid \circ,$$
$$\mid \circ \mid \circ \ \circ \ \circ \ \circ,$$

and so on. It should be clear that you need $n$ $\circ$-symbols and $d - 1$ $\mid$ separators. The number of possible such combinations, and hence the dimensionality of the embedding space, is then

$$\binom{n + d - 1}{n} = \frac{(n + d - 1)!}{n!(d - 1)!}. \tag{3.29}$$

This number quickly grows to be computationally infeasible, which motivates one to look for a way to compute the inner product in the embedded space without performing the explicit mapping itself. This is the point of the so-called 'kernel trick', which I'll introduce in the following.

The idea of the kernel trick is that since only the inner products between feature vectors in the embedded space are required, one might as well look for some function $K$ of the original feature vectors which gives the same scalar as the inner product in the embedded space, i.e.

$$K(\mathbf{x}, \mathbf{y}) = \Phi(\mathbf{x}) \cdot \Phi(\mathbf{y}). \tag{3.30}$$

In the polynomial case treated above, the correspondence between the kernel function $K(\mathbf{x}, \mathbf{y})$ and the explicit mapping $\Phi$ is straightforward:

$$K(\mathbf{x}, \mathbf{y}) = (\mathbf{x} \cdot \mathbf{y})^n, \tag{3.31}$$

$$\Phi(\mathbf{x}) = \sum_{K(n,d)} \sqrt{C_{\{k\}}} \prod_{j=1}^{d} x_j^{k_j}, \tag{3.32}$$

so that (3.30) is true by the multinomial theorem and the above considerations. However, situations arise in which the explicit mapping $\Phi$ isn't directly obtainable, and the correspondence of the kernel function to inner products in higher dimensional spaces is harder to demonstrate. This is the subject of the following section.

### Radial Basis Functions

One commonly used kernel function is the RBF, or radial basis function, kernel:

$$K(\mathbf{x}, \mathbf{y}) = e^{|\mathbf{x}-\mathbf{y}|^2/2\sigma}. \tag{3.33}$$

Burges [4] shows that the polynomial kernel is valid, so I'll show how the argument extends to the RBF kernel in the following.

find en eller anden kilde

Mercer's condition states that for a kernel function $K(\mathbf{x}, \mathbf{y})$, there exists a corresponding Hilbert space $\mathcal{H}$ and a mapping $\Phi$ as specified earlier, iff any $L^2$-normalizable function $g(\mathbf{x})$ satisfies

$$\int K(\mathbf{x}, \mathbf{y}) g(\mathbf{x}) g(\mathbf{y}) \, d\mathbf{x} \, d\mathbf{y} \geq 0. \tag{3.34}$$

This can be shown be rewriting (3.33) as

$$K(\mathbf{x}, \mathbf{y}) = e^{(\mathbf{x}-\mathbf{y})\cdot(\mathbf{x}-\mathbf{y})/2\sigma} = e^{|\mathbf{x}|^2/2\sigma} e^{|\mathbf{y}|^2/2\sigma} e^{-\mathbf{x}\cdot\mathbf{y}/\sigma}, \tag{3.35}$$

and expanding the last term in $(\mathbf{x} \cdot \mathbf{y})$ as

$$e^{-\mathbf{x}\cdot\mathbf{y}/\sigma} = \sum_{i=0}^{\infty} \frac{(-1)^i}{i!\sigma^i} (\mathbf{x} \cdot \mathbf{y})^i, \tag{3.36}$$

but using (3.27) on the dot product gives

$$(\mathbf{x} \cdot \mathbf{y})^i = \left( \sum_{j=1}^{d} x_j y_j \right)^i = \sum_{K(i,d)} C_{\{k\}} \prod_{j=1}^{d} x_j^{k_j} y_j^{k_j} \tag{3.37}$$

so the Taylor expansion becomes

$$\mathrm{e}^{-\mathbf{x}\cdot\mathbf{y}/\sigma} = \sum_{i=0}^{\infty} \sum_{K(i,d)} \frac{(-1)^i}{i!\sigma^i} C_{\{k\}} \prod_{j=1}^{d} x_j^{k_j} y_j^{k_j}, \tag{3.38}$$

which can be plugged back into (3.35) to yield

$$K(\mathbf{x}, \mathbf{y}) = \sum_{i=0}^{\infty} \sum_{K(i,d)} \frac{(-1)^i}{i!\sigma^i} C_{\{k\}} \mathrm{e}^{|\mathbf{x}|^2/2\sigma} \mathrm{e}^{|\mathbf{y}|^2/2\sigma} \prod_{j=1}^{d} x_j^{k_j} y_j^{k_j}. \tag{3.39}$$

The underlying reason for these algebraic shenanigans is that (3.39) is clearly separable so that the integral in (3.34) from from Mercer's condition becomes

$$\int K(\mathbf{x}, \mathbf{y}) g(\mathbf{x}) g(\mathbf{y}) \,\mathrm{d}\mathbf{x} \,\mathrm{d}\mathbf{y} \tag{3.40}$$

$$= \sum_{i=0}^{\infty} \sum_{K(i,d)} \frac{(-1)^i}{i!\sigma^i} C_{\{k\}} \int_{\mathbb{R}^{2d}} \mathrm{e}^{|\mathbf{x}|^2/2\sigma} \mathrm{e}^{|\mathbf{y}|^2/2\sigma} \prod_{j=1}^{d} x_j^{k_j} y_j^{k_j} g(\mathbf{x}) g(\mathbf{y}) \,\mathrm{d}\mathbf{x} \,\mathrm{d}\mathbf{y} \tag{3.41}$$

$$= \sum_{i=0}^{\infty} \sum_{K(i,d)} \frac{(-1)^i}{i!\sigma^i} C_{\{k\}} \left( \int_{\mathbb{R}^{d}} \mathrm{e}^{|\mathbf{x}|^2/2\sigma} \prod_{j=1}^{d} x_j^{k_j} g(\mathbf{x}) \,\mathrm{d}\mathbf{x} \right) \cdot \left( \int_{\mathbb{R}^{d}} \mathrm{e}^{|\mathbf{y}|^2/2\sigma} \prod_{j=1}^{d} y_j^{k_j} g(\mathbf{y}) \,\mathrm{d}\mathbf{y} \right) \tag{3.42}$$

$$= \sum_{i=0}^{\infty} \sum_{K(i,d)} \frac{(-1)^i}{i!\sigma^i} C_{\{k\}} \left( \int_{\mathbb{R}^{d}} \mathrm{e}^{|\mathbf{x}|^2/2\sigma} \prod_{j=1}^{d} x_j^{k_j} g(\mathbf{x}) \,\mathrm{d}\mathbf{x} \right)^2 \tag{3.43}$$

$$\geq 0. \tag{3.44}$$

Hence, radial basis functions satisfy Mercer's condition and the kernel described above can be plugged into the dual Lagrangian from (3.14) to obtain

$$L_D = -\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathrm{e}^{|\mathbf{x_i}-\mathbf{x_j}|^2/2\sigma} + \sum_{i} \alpha_i, \tag{3.45}$$

which must be maximized subject to the same constraints as earlier. The concrete optimization procedure is complicated and already implemented in most machine learning libraries, so I choose not to go into details with that, but instead to demonstrate the effectiveness of the RBF kernel approach on the non-linear-separable points that were generated earlier. figure 3.3 shows the points again, along with the *decision frontier* i.e. the curve which separates regions in which points are classified into separate

Figure 3.3: The 'island' scenario of figure 2.14 revisited. The points representing class a and class b have been mapped to a higher-dimensional space in which it is possible to construct a separating hyperplane whose decision frontier is also shown.

classes. The danger of overfitting should be clear from figure 3.3. If the cost of misclassification $C$ and the sharpness of the RBFs, usually denoted by $\gamma = 2/\sigma$ are set sufficiently high, the algorithm will simply end up with a tiny decision boundary around every training point of class a, resulting in flawless classification on the training set, but utter failure on new data. The typical way of evaluating this is to perform k-fold validation, meaning that the available data is $k$ equal parts and the SVM is consecutively trained on $k - 1$ parts and tested on the remaining. A variant of this, which my code uses, is stratified k-fold validation, which only differs in that the data is partitioned so as to keep the ratio between the different classes in each parts as close to equal as possible.

The $\gamma$ parameter is often fixed by performing a grid search similar to that discussed earlier. Figure 3.4 shows the resulting heat map from a grid search.

Figure 3.4: Result of a grid search for the optimal combination of values for the cost parameter $C$ and the sharpness $\gamma$ of the Gaussian kernel function used.

### 3.1.3 Statistical subtleties

An important note should be made here about some often neglected subtleties relating to uncertainties. Physicists often deal with measurements that can assumed to be independently drawn from a normal distribution $\mathcal{N}(x_i; \mu, \sigma^2)$ due to the central limit theorem. With a large number of measurements $n$, the standard deviation of a sample

$$\sigma^2 = \frac{1}{N} \sum_i^N (x_i - \mu)^2 , \tag{3.46}$$

converges as $N \to \infty$ to the maximum likelihood, minimum variance unbiased estimator for the true variance of the underlying distribution with unknown mean

$$\hat{\sigma}^2 = \frac{1}{N-1} \sum_i^N (x_i - \mu)^2 . \tag{3.47}$$

The standard deviation $\sigma$ and the width of the underlying gaussian $\hat{\sigma}^2$ can then often be used interchangeably. This tempts some people into the

questionable habit of always assuming that the sample standard deviance can be used as the 68% confidence interval of their results.

When using a K-fold validation scheme, the performance scores for the various folds cannot be assumed to be independently drawn from an underlying distribution, as the test set of one fold is used in the training sets of the remaining folds. In fact, it has been shown [1] that there is no unbiased estimator for the variance of the performance estimated using K-fold validation. However, as K-fold validation is more effective than keeping the test, and training data separate, which can be shown using Jensen's inequality along with some basic properties of expectation values [2], I'll mostly use K-fold regardless. As the standard deviation still provides a qualitative measure of the consistency of the model's performance, I'll still use the sample STD in a usual fashion, such as error bars, unless otherwise is specified, but the reader should keep in mind that these do not indicate precise uncertainties whenever K-fold validation has been involved.

## 3.2   Decision Trees & Random Forests

Another popular machine learning scheme is that of random forests, which consist of an ensemble of decision trees. A decision tree is a very intuitive method for classification problems which can be visualized as a kind of flow chart in the following fashion. As usual, the problem consists of a set of feature vectors $\mathbf{x}_i$ and a set of corresponding class labels $y_i$. A decision tree then resembles a flowchart starting at the root of the tree, at each node splitting into branches and finally branching into leaves at which all class labels should be identical. At each node, part of the feature vector is used to split the dataset into parts. This resembles the 'twenty questions' game, in which one participant thinks of a famous person and another attempts to guess who it is by asking a series of yes/no-questions, each one splitting the set of candidates in two parts. In this riddle game and in decision tree learning, there are good and bad questions (asking whether the person was born on March 14th, 1879 is a very bad first question, for instance). There are several ways of quantifying how 'good' a yes/no-question, corresponding to a partitioning of the dataset, is.

On metric for this is the Gini Impurity Index $I_G$, which is computed by summing over each class label:

$$I_G = \sum_i f_i (1 - f_i) = 1 - \sum_i f_i^2, \tag{3.48}$$

where $f_i$ denotes the fraction of the set the consists of class $y_i$. Using this as a metric, the best partitioning is the one which results in the largest drop in total Gini impurity following a branching. Another metric is the information gain measured by comparing the entropy before a split with a weighted average of the entropy in the groups resulting from the split. Denoting the fractions of the various classes in the parent group, i.e. before splitting, by $f_i$ and the two child groups by $a_i$ and $b_i$, the information gain is

$$I_E = -\sum_i f_i \log_2 f_i + \frac{n_a}{N} \sum_i a_i \log_2 a_i + \frac{n_b}{N} \sum_i b_i \log_2 b_i. \qquad (3.49)$$

However, if too many such nodes are added to a decision tree, over-fitting, i.e. extreme accuracies on training data but poor performance on new data, becomes a problem. This can be remedied by instead predicting with a majority vote, or averaging in the case of regression problems, from an ensemble of randomized decision trees called a random forest. The main merits of random forests are their accuracy and ease of use, and their applications as auxiliary methods in other machine learning schemes, which I'll elaborate on shortly.

The individual trees in a random forest are grown using a randomly selected subset of the training data for each tree. The data used to construct a given tree is referred to as 'in bag', whereas the remaining training data is referred to as 'out of bag' (OOB) for the given tree. At each node, a set number of features is randomly selected and the best possible branching, cf. the above considerations, is determined. The only parameters that must be tweaked manually are the number of trees in the forest, number of features to include in each branching, and the maximum tree depth. While other variables such as the metric for determining branching quality as described above, may be customized, those aren't essential to achieve a decent predictor, which is robust in regard to both overfitting and irrelevant parameters.[3]

There doesn't seem to be a single universally accepted way of adjusting these parameters, so I chose a somewhat pragmatic approach of simply checking how well various choices for each parameter performed on a randomly selected trait. For instance, figure 3.5 shows how well a random forest predicted the tertiles of participants' extroversion as a function of the fraction of available features each tree was allowed to include in each branching. This was done using a large number of trees ($n = 1000$) and using each of the two metrics described earlier. The number of features used pr split doesn't seem to have any significant effect on performance, and as the entropy metric seems to perform as well or slightly better than Gini impurity, I decided to stick to that. A similar plot of the performance

Figure 3.5: Performance of a random forest with 1000 decision trees using various fractions of the available features in each branching using both the entropy and the Gini impurity metric to determine the optimal branching. The number of features seems not to play a major role, and the entropy metric seems to perform slightly better in general.

of various numbers of decision trees in the forest is shown in figure 3.6. The performance seems to stagnate around 100 trees, and remain constant after that, so I usually used at least 500 trees to make sure to get the optimal performance, as runtime wasn't an issue.

The robustness to irrelevant features and overfitting described earlier also plays a role in the application of random forests in conjunction with other schemes. SVMs as described in section 3.1 can be sensitive to irrelevant data[14]. There exist off-the-shelf methods, such as recursive feature elimination (RFE)[9], for use with linear SVMs, but to my knowledge, there is no 'standard' way to eliminate irrelevant features when using a non-linear kernel. However, it is possible to use a random forest approach to obtain the relative importance of the various features and then use only the most important ones in another machine learning scheme which is less tolerant to the inclusion of irrelevant data. The relative importance of feature $j$ can be estimated by first constructing a random forest and evaluating its performance $s$, then randomly permuting the values of feature $j$ across the training sample and measure the damage it does to

Figure 3.6: Example of a random forest performance versus number of decision trees. Performance seems to increase steadily until about 100 trees, then stagnate.

the performance of the forest by comparing with the permuted score $s_p$. The ratio of the mean to standard deviation of those differences:

$$w_j = \frac{\langle s - s_p \rangle}{\text{std}(s - s_p)} \tag{3.50}$$

Random forests also provide a natural measure of similarity between data points. Given data points $i$ and $j$ these can be plugged into all their OOB decision trees, or a random subset thereof, and the fraction of the attempts in which both end up at the same leaf can be taken as a measure of similarity. This can be used to generate a proximity matrix for the data points, and it can be used as a metric for determining the nearest neighbours of a new point in conjunction with a simple nearest neighbour classifier.

## 3.3 Nearest Neighbour-classifiers

Skriv en masse om den smart random forest NN-model. Do iiit!

Figure 3.7: Comparison of performance of models using random forest and support vector regression with a baseline model which always predicts the mean of the training sample. The y axis shows the mean error of each model and the error bars show the 95-percentile around the median scores obtained by running on 1000 bootstrap samples.

## 3.4   Results

HARJ

# Appendix

# A.1   Source Code for Explicit Semantic Analysis

This section contains code pertaining to part I of the thesis.

## A.1.1   Parser

```
1   # −∗− coding: utf−8 −∗−
2   '''Parses a full Wikipedia XML-dump and saves to files containing
3   a maximum of 1000 articles.
4   In the end, each file is saved as a JSON file containing entries like:
5   {{
6       'concept':
7       {
8       'text': <article contents>,
9       'links_in' : Set of links TO the article in question,
10      'links_out' : Set of links FROM the article in question,
11      }
12  }
13  Although links_in is added by the generate_indices script.
14  Also saved are dicts for keeping track of word and concept indices when
15  building a large sparse matrix for the semantic interpreter.
16  The file structure is like {'word blah' : index blah}'''
17
18  import re
19  import xml.sax as SAX
20  import wikicleaner
21  import os
22  import glob
23  import shared
24  import sys
25
26  DEFAULT_FILENAME = 'medium_wiki.xml'
27
28  def canonize_title(title):
29      # remove leading whitespace and underscores
30      title = title.strip(' _')
31      # replace sequencesences of whitespace and underscore chars with a single space
32      title = re.compile(r'[\s_]+').sub(' ', title)
33      #remove forbidden characters
34      title = re.sub('[?/\\\*\"\']','',title)
35      return title.title()
36
37  #Import shared parameters
38  from shared import extensions, temp_dir
39
40  #Cleanup
41  for ext in extensions.values():
42      for f in glob.glob(temp_dir + '*'+ext):
43          os.remove(f)
44
45  def filename_generator(folder):
46      '''Generator for output filenames'''
47      if not os.path.exists(folder):
48          os.makedirs(folder)
49      count = 0
50      while True:
51          filename = folder+"content"+str(count)
```

```python
            count += 1
            yield filename

make_filename = filename_generator(temp_dir)

#Format {right title : redirected title }, e.g. {because : ([cuz, cus])}
redirects = {}

#Minimum number of links/words required to keep an article.
from shared import min_links_out, min_words

#Open log file for writing and import logging function
logfile = open(os.path.basename(__file__)+'.log', 'w')
log = shared.logmaker(logfile)

class WikiHandler(SAX.ContentHandler):
    '''ContentHandler class to process XML and deal with the WikiText.
    It works basically like this:
    It traverses the XML file, keeping track of the type of data being read and
    adding any text to its input buffer. When event handlers register a page
    end, the page content is processed, the processed content is placed in the
    output buffer, and the input buffer is flushed.
    Whenever a set number of articles have been processed, the output buffer is
    written to a file. The point of this approach is to
    limit memory consumption.'''

    def __init__(self):
        SAX.ContentHandler.__init__(self)
        self.current_data = None
        self.title = ''
        self.input_buffer = []
        self.output_buffer = {}
        self.article_counter = 0
        self.links = []
        self.categories = []
        self.redirect = None
        self.verbose = False
        #Harvest unique words here
        self.words = set([])
        #keeps track of ingoing article links. format {to : set([from])}
        self.linkhash = {}

    def flush_input_buffer(self):
        '''Deletes info on the currently processed article.
         This is called when a page end event is registered.'''
        self.input_buffer = []
        self.current_data = None
        self.title = ''
        self.links = []
        self.categories = []
        self.redirect = None

    def flush_output_buffer(self):
        '''Flushes data gathered so far to a file and resets.'''
        self.output_buffer = {}
        self.words = set([])
        self.linkhash = {}

    def startElement(self, tag, attrs):
        '''Eventhandler for element start - keeps track of current datatype.'''
        self.current_data = tag
        #Informs the parser of the redirect destination of the article
```

```python
114          if tag == "redirect":
115              self.redirect = attrs['title']
116              return None
117
118      def endElement(self, name):
119          '''Eventhandler for element end. This causes the parser to process
120           its input buffer when a pageend is encountered.'''
121          #Process content after each page
122          if name == 'page':
123              self.process()
124          #Write remaining data at EOF.
125          elif name == 'mediawiki':
126              self.writeout()
127
128      def characters(self, content):
129          '''Character event handler. This simply passes any raw text from an
130           article field to the input buffer and updates title info.'''
131          if self.current_data == 'text':
132              self.input_buffer.append(content)
133          elif self.current_data == 'title' and not content.isspace():
134              self.title = content
135
136      def process(self):
137          '''Process input buffer contents. This converts wikilanguage to
138           plaintext, registers link information and checks if content has
139           sufficient words and outgoing links (ingoing links can't be checked
140           until the full XML file is processed).'''
141
142          #Ignore everything else if article redirects
143          if self.redirect:
144              self.flush_input_buffer()
145              return None
146              global redirects
147              try:
148                  redirects[self.title].add(self.redirect)
149              except KeyError:
150                  redirects[self.title] = set([self.redirect])
151              self.flush_input_buffer()
152              return None
153
154          #Redirects handled – commence processing
155          print "processing: "+self.title.encode('utf8')
156          #Combine buffer content to a single string
157          text = ''.join(self.input_buffer).lower()
158
159          #Find and process link information
160          link_regexp = re.compile(r'\[\[(.*?)\]')
161          links = re.findall(link_regexp, text) #grap stuff like [[<something>]]
162          #Add links to the parsers link hash
163          for link in links:
164              #Check if link matches a namespace, e.g. 'file:something.png'
165              if any([ns+':' in link for ns in wikicleaner.namespaces]):
166                  continue #Proceed to next link
167              #Namespaces done, so remove any colons:
168              link = link.replace(':', '')
169              if not link:
170                  continue #Some noob could've written an empty link...
171              #remove chapter designations/displaytext – keep article title
172              raw = re.match(r'([^\|\#]*)', link).group(0)
173              title = canonize_title(raw)
174              #note down that current article has outgoing link to 'title'
175              self.links.append(title)
```

```python
176              #also note that 'title' has incoming link from here
177              try:
178                  self.linkhash[title].add(self.title) #maps target−>sources
179              except KeyError:
180                  self.linkhash[title] = set([self.title])
181
182          #Disregard current article if it contains too few links
183          if len(self.links) < min_links_out:
184              self.flush_input_buffer()
185              return None
186
187          #Cleanup text
188          text = wikicleaner.clean(text)
189          article_words = text.split()
190
191          #Disregard article if it contains too few words
192          if len(article_words) < min_words:
193              self.flush_input_buffer()
194              return None
195
196          #Update global list of unique words
197          self.words.update(set(article_words))
198
199          #Add content to output buffer
200          output = {
201              'text': text,
202              #Don't use category info for now
203              #'categories' : self.categories,
204              'links_out': self.links
205          }
206          self.output_buffer[self.title] = output
207          self.article_counter += 1
208
209          #Flush output buffer to file
210          if self.article_counter%1000 == 0:
211              self.writeout()
212
213          #Done, flushing buffer
214          self.flush_input_buffer()
215          return None
216
217      def writeout(self):
218          '''Writes output buffer contents to file'''
219          #Generate filename and write to file
220          filename = make_filename.next()
221          #Write article contents to file
222          with open(filename+extensions['content'], 'w') as f:
223              shared.dump(self.output_buffer, f)
224
225          #Store wordlist as files
226          with open(filename+extensions['words'], 'w') as f:
227              shared.dump(self.words, f)
228
229          #Store linkhash in files
230          with open(filename+extensions['links'], 'w') as f:
231              shared.dump(self.linkhash, f)
232
233          if self.verbose:
234              log("wrote "+filename)
235
236          #Empty output buffer
237          self.flush_output_buffer()
```

```
238          return None
239
240  if __name__ == "__main__":
241      if len(sys.argv) == 2:
242          file_to_parse = sys.argv[1]
243      else:
244          file_to_parse = DEFAULT_FILENAME
245
246      #Create and configure content handler
247      test = WikiHandler()
248      test.verbose = True
249
250      #Create a parser and set handler
251      ATST = SAX.make_parser()
252      ATST.setContentHandler(test)
253
254      #Let the parser walk the file
255      log("Parsing started...")
256      ATST.parse(file_to_parse)
257      log("...Parsing done!")
258
259      #Attempt to send notification that job is done
260      if shared.notify:
261          try:
262              shared.pushme(sys.argv[0]+' completed.')
263          except:
264              log("Job's done. Push failed.")
265
266      logfile.close()
```

## A.1.2 Index Generator

```python
# −*− coding: utf−8 −*−
'''This finishes preproccessing of the output from the XML parser.
This script reads in link data and removes from the content files those
concepts that have too few incoming links. Information on incoming links
is saved to each content file.
Finally, index maps for words and approved concepts are generated and saved.'''

from __future__ import division
import glob
import gc
import shared
import os
import sys

logfile = open(os.path.basename(__file__)+'.log','w')
log = shared.logmaker(logfile)

#Import shared parameters
from shared import extensions, temp_dir, min_links_in, matrix_dir

def listchopper(l):
    '''Generator to chop lists into chunks of a predefined length'''
    n = shared.link_chunk_size
    ind = 0
    while ind < len(l):
        yield l[ind:ind+n]
        ind += n

def main():
    #Import shared parameters and verify output dir exists
    if not os.path.exists(temp_dir):
        raise IOError

#===============================================================================
#     Read in link data and update content files accordingly
#===============================================================================

    #Get list of files containing link info and chop it up
    linkfiles = glob.glob(temp_dir + '*'+extensions['links'])
    linkchunks = listchopper(linkfiles)

    linkfiles_read = 0
    for linkchunk in linkchunks:
        #Hash mapping each article to a set of articles linking to it
        linkhash = {}

        for filename in linkchunk:
            with open(filename, 'r') as f:
                newstuff = shared.load(f)
            #Add link info to linkhash
            for target, sources in newstuff.iteritems():
                try:
                    linkhash[target].update(set(sources))
                except KeyError:
                    linkhash[target] = set(sources)

            #Log status
```

```
58              linkfiles_read += 1
59              log("Read " + filename + " - " +
60                  str(100*linkfiles_read/len(linkfiles))[:4] + " % of link data.")
61
62          log("Chunk finished - updating content files")
63          #Update concept with newly read link data
64          contentfiles = glob.glob(temp_dir + '*'+extensions['content'])
65          contentfiles_read = 0
66          for filename in contentfiles:
67              #Read file. Content is like {' article  title ' : {'text' : blah}}
68              with open(filename, 'r') as f:
69                  content = shared.load(f)
70
71              #Search linkhash for links going TO concept
72              for concept in content.keys():
73                  try:
74                      sources = linkhash[concept]
75                  except KeyError:
76                      sources = set([])   #Missing key => zero incoming links
77
78                  #Update link info for concept
79                  try:
80                      content[concept]['links_in'] = set(content[concept]['links_in'])
81                      content[concept]['links_in'].update(sources)
82                  except KeyError:
83                      content[concept]['links_in'] = sources
84
85              #Save updated content
86              with open(filename, 'w') as f:
87                  shared.dump(content, f)
88
89              contentfiles_read += 1
90              if contentfiles_read % 100 == 0:
91                  log("Fixed " + str(100*contentfiles_read/len(contentfiles))[:4]
92                      + "% of content files")
93          pass  #Proceed to next link chunk
94
95  #==============================================================================
96  #       Finished link processing
97  #       Remove unworthy concepts and combine concept/word lists.
98  #==============================================================================
99
100     #What, you think memory grows on trees?
101     del linkhash
102     gc.collect()
103
104     #Set of all approved concepts
105     concept_list = set([])
106
107     #Purge inferior concepts (with insufficient incoming links)
108     for filename in contentfiles:
109         #Read in content file
110         with open(filename, 'r') as f:
111             content = shared.load(f)
112
113         for concept in content.keys():
114             entry = content[concept]
115             if 'links_in' in entry and len(entry['links_in']) >= min_links_in:
116                 concept_list.add(concept)
117             else:
118                 del content[concept]
119
```

```
120        with open(filename, 'w') as f:
121            shared.dump(content, f)
122
123    log("Links done - saving index files")
124
125    #Make sure output dir exists
126    if not os.path.exists(matrix_dir):
127        os.makedirs(matrix_dir)
128
129    #Generate and save a concept index map. Structure: {concept : index}
130    concept_indices = {n: m for m,n in enumerate(concept_list)}
131    with open(matrix_dir+'concept2index.ind', 'w') as f:
132        shared.dump(concept_indices, f)
133
134    #Read in all wordlists and combine them.
135    words = set([])
136    for filename in glob.glob(temp_dir + '*'+extensions['words']):
137        with open(filename, 'r') as f:
138            words.update(shared.load(f))
139
140    #Generate and save a word index map. Structure: {word : index}
141    word_indices = {n: m for m,n in enumerate(words)}
142    with open(matrix_dir+'word2index.ind', 'w') as f:
143        shared.dump(word_indices, f)
144
145    log("Wrapping up.")
146    #Attempt to notify that job is done
147    if shared.notify:
148        try:
149            shared.pushme(sys.argv[0]+' completed.')
150        except:
151            log("Job's done. Push failed.")
152
153    logfile.close()
154
155 if __name__ == '__main__':
156     main()
```

### A.1.3   Matrix Builder

```
1    # −*− coding: utf−8 −*−
2    '''Builds a huge sparse matrix of Term frequency/Inverse Document Frequency
3    (TFIDF) of the previously extracted words and concepts.
4    First a matrix containing simply the number of occurrences of word i in the
5    article corresponding to concept j is build (in DOK format as that is faster
6    for iterative construction), then the matrix is converted to sparse row format
7    (CSR), TFIDF values are computed, each row is normalized and finally pruned.'''
8
9    from __future__ import division
10   import scipy.sparse as sps
11   import numpy as np
12   from collections import Counter
13   import glob
14   import shared
15   import sys
16   import os
17
18   def percentof(small, large):
19       return str(100*small/large) + "%"
20
21   logfile = open(os.path.basename(__file__)+'.log', 'w')
22   log = shared.logmaker(logfile)
23
24   #import shared parameters
25   from shared import (extensions, matrix_dir, prune, temp_dir, column_chunk_size,
26                       row_chunk_size, datatype)
27
28   def main():
29       #Cleanup
30       for f in glob.glob(matrix_dir + '/*'+extensions['matrix']):
31           os.remove(f)
32
33       #Set pruning parameters
34       window_size = shared.window_size
35       cutoff = shared.cutoff
36
37       #Read in dicts mapping words and concepts to their respective indices
38       log("Reading in word/index data")
39       word2index = shared.load(open(matrix_dir+'word2index.ind', 'r'))
40       concept2index = shared.load(open(matrix_dir+'concept2index.ind', 'r'))
41       log("...Done!")
42
43   #===============================================================================
44   #       Construct count matrix in small chunks
45   #===============================================================================
46
47       #Count words and concepts
48       n_words = len(word2index)
49       n_concepts = len(concept2index)
50
51       #Determine matrix dimensions
52       matrix_shape = (n_words, n_concepts)
53
54       #Allocate sparse matrix. Dict−of−keys should be faster for iterative
55       #construction. Convert to csr for fast row operations later.
56       mtx = sps.dok_matrix(matrix_shape, dtype = datatype)
57
```

```python
58    def matrix_chopper(matrix, dim):
59        '''Generator to split a huge matrix into small submatrices, which can
60         then be stored in individual files.
61         This is handy both when constructing the matrix (building the whole
62         matrix without saving to files in the process takes about 50 gigs RAM),
63         and when applying it, as this allows one to load only the submatrix
64         relevant to a given word.'''
65        ind = 0
66        counter = 0
67        rows = matrix.get_shape()[0]
68        while ind < rows:
69            end = min(ind+dim, rows)
70            #Return pair of submatrix number and the submatrix itself
71            yield counter, sps.vstack([matrix.getrow(i)\
72                                        for i in xrange(ind, end)], format = 'csr')
73            counter += 1
74            ind += dim
75
76    def writeout():
77        '''Saves the matrix as small submatrices in separate files.'''
78        for n, submatrix in matrix_chopper(mtx, row_chunk_size):
79            filename = matrix_dir+str(n)+extensions['matrix']
80            #Update submatrix if it's already partially calculated
81            log("Writing out chunk %s" % n)
82            try:
83                with open(filename, 'r') as f:
84                    submatrix = submatrix + shared.mload(f)
85                #
86            except IOError:
87                pass #File doesn't exist yet, so no need to change mtx
88
89            #Dump the submatrix to file
90            with open(filename, 'w') as f:
91                shared.mdump(submatrix, f)
92        return None
93
94    log("Constructing matrix.")
95    filelist = glob.glob(temp_dir + '*'+extensions['content'])
96    files_read = 0
97    for filename in filelist:
98        with open(filename, 'r') as f:
99            content = shared.load(f)
100
101        #Loop over concepts (columns) as so we don't waste time with rare words
102        for concept, entry, in content.iteritems():
103            #This is the column index (concept w. index j)
104            j = concept2index[concept]
105
106            #Convert concept 'countmap' like so: {word : n}
107            wordmap = Counter(entry['text'].split()).iteritems()
108
109            #Add them all to the matrix
110            for word, count in wordmap:
111                #Find row index of the current word
112                i = word2index[word]
113
114                #Add the number of times word i occurs in concept j to the matrix
115                mtx[i,j] = count
116            #
117        #Update file count
118        files_read += 1
119        log("Processed content file no. %s of %s - %s"
```

```
120                 % (files_read, len(filelist)−1, percentof(files_read, len(filelist))))
121
122          if files_read % column_chunk_size == 0:
123              mtx = mtx.tocsr()
124              writeout()
125              mtx = sps.dok_matrix(matrix_shape)
126          #
127
128      #Convert matrix to CSR format and write to files.
129      mtx = mtx.tocsr()
130      writeout()
131
132  #==============================================================================
133  # Count matrix/matrices constructed − computing TF−IDF
134  #==============================================================================
135
136      log("Done - computing TF-IDF")
137
138      #Grap list of matrix files (containing the submatrices from before)
139      matrixfiles = glob.glob(matrix_dir + "*" + extensions['matrix'])
140      words_processed = 0  #for logging purposes
141
142      for filename in matrixfiles:
143          with open(filename, 'r') as f:
144              mtx = shared.mload(f)
145
146          #Number of words in a submatrix
147          n_rows = mtx.get_shape()[0]
148
149          for w in xrange(n_rows):
150              #Grap non−zero elements from the row corresonding to word w
151              row = mtx.data[mtx.indptr[w] : mtx.indptr[w+1]]
152              if len(row) == 0:
153                  continue
154
155              #Make a vectorized function to convert a full row to TF−IDF
156              f = np.vectorize(lambda m_ij: (1+np.log(m_ij))*
157                                  np.log(n_concepts/len(row)))
158
159              #Map all elements to TF−IDF and update matrix
160              row = f(row)
161
162              #Normalize the row
163              assert row.dtype.kind == 'f'  #Non floats round to zero w/o warning
164              normfact = 1.0/np.linalg.norm(row)
165              row *= normfact
166
167              #Start inverted index pruning
168              if prune:
169                  #Number of documents containing w
170                  n_docs = len(row)
171
172                  #Don't prune if the windows exceeds the array bounds (duh)
173                  if window_size < n_docs:
174
175                      #Obtain list of indices such that row[index] is sorted
176                      indices = np.argsort(row)[::−1]
177
178                      #Generate a sorted row
179                      sorted_row = [row[index] for index in indices]
180
181                      #Go through sorted row and truncate when pruning condition is met
```

```python
182                         for i in xrange(n_docs-window_size):
183                             if sorted_row[i+window_size] >= cutoff*sorted_row[i]:
184                                 #Truncate, i.e. set the remaining entries to zero
185                                 sorted_row[i:] = [0]*(n_docs-i)
186                                 break
187                             else:
188                                 pass
189
190                         #Unsort to original positions
191                         for i in xrange(n_docs):
192                             row[indices[i]] = sorted_row[i]
193
194                 #Update matrix
195                 mtx.data[mtx.indptr[w] : mtx.indptr[w+1]] = row
196
197                 #Log it
198                 words_processed += 1
199                 if words_processed % 10**3 == 0:
200                     log("Processing word %s of %s - %s" %
201                         (words_processed, n_words,
202                          percentof(words_processed, n_words)))
203
204         #Keep it sparse - no need to store zeroes
205         mtx.eliminate_zeros()
206         with open(filename, 'w') as f:
207             shared.mdump(mtx, f)
208
209     log("Done!")
210
211     #Notify that the job is done
212     if shared.notify:
213         try:
214             shared.pushme(sys.argv[0]+' completed.')
215         except:
216             log("Job's done. Push failed.")
217
218     logfile.close()
219     return None
220
221 if __name__ == '__main__':
222     main()
```

### A.1.4 Library for Computational Linguistics

```python
# −*− coding: utf−8 −*−
'''Small module for computational linguistics applied to Twitter.
The main classes are a TweetHarvester, which gathers data from Twitters' API,
and a SemanticAnalyser, which relies on the previously constructed TFIDF
matrices.'''

from __future__ import division
from scipy import sparse as sps
from collections import Counter
from numpy.linalg import norm
import re
import shared
import tweepy
from datetime import date
import json
import time
import sys
import codecs
from pprint import pprint
sys.stdout = codecs.getwriter('utf8')(sys.stdout)
sys.stderr = codecs.getwriter('utf8')(sys.stderr)

#==============================================================================
# This stuff defines a twitter 'harvester' for downloading Tweets
#==============================================================================

#Import credentials for accessing Twitter API
from supersecretstuff import consumer_key, consumer_secret, access_token, access_token_secret
auth = tweepy.OAuthHandler(consumer_key, consumer_secret)
auth.set_access_token(access_token, access_token_secret)

class listener(tweepy.StreamListener):
    '''Listener class to access Twitter stream.'''
    #What to do with a tweet (override later )
    def process(self, content):
        print content
        return None

    def on_status(self, status):
        self.process(status)
        return True

    def on_error(self, status):
        print status

# Exception to be raised  when the Twitter API messes up. Happens occasionally.
class IncompleteRead(Exception):
    pass

class TweetHarvester(object):
    '''Simple class to handle tweet harvest.
     Harvest can be performed actively or passively, i.e. using the 'mine'
     method to gather a fixed number of tweets or using the 'listen' method
     to stream tweets matching a given search term.
     Harvested tweets are sent to the process method which by default simply
     stores them inside the object.'''
```

```python
58    def __init__(self, max_tweets=-1, verbose = False, tweets_pr_file = 10**5):
59        #Set parameters
60        self.max_tweets = max_tweets #-1 for unlimited stream
61        self.verbose = verbose
62        self.tweets_pr_file = tweets_pr_file
63
64        #Internal parameters to keep track of harvest status
65        self.files_saved = 0
66        self.harvested_tweets = []
67        self.current_filenumber = 0
68        self.current_date = date.today()
69
70    def filename_maker(self):
71        #Update counter and date if neccessary
72        if not self.current_date == date.today():
73            self.current_date = date.today()
74            self.current_filenumber = 0
75        else:
76            pass #Date hasn't changed. Proceed.
77        filename = str(self.current_date) + "-data%s.json" % self.current_filenumber
78        self.current_filenumber += 1
79        return filename
80
81    #Simple logging function
82    def log(self, text):
83        string = text+" at "+time.asctime()+"\n"
84        if self.verbose:
85            print string
86        with open('tweetlog.log', 'a') as logfile:
87            logfile.write(string)
88        #Must return true so I can log errors without breaking the stream.
89        return True
90
91    def listen(self, search_term):
92        #Make a listener
93        listener = tweepy.StreamListener()
94        #Override relevant methods.
95        listener.on_status = self.process
96        listener.on_error = lambda status_code: self.log("Error: "+status_code)
97        listener.on_timeout = lambda: self.log("Timeout.")
98
99        twitterStream = tweepy.Stream(auth, listener)
100       twitterStream.filter(track=search_term)
101
102   def mine(self, search_term, n = None):
103       '''Mine a predefined number of tweets using input search word'''
104       if n == None:
105           n = self.max_tweets
106
107       api = tweepy.API(auth)
108       tweets = tweepy.Cursor(api.search, q=search_term).items(n)
109       for tweet in tweets:
110           self.process(tweet)
111
112   def process(self, tweet):
113       self.harvested_tweets.append(tweet)
114       if self.verbose:
115           print "Holding %s tweets." % len(self.harvested_tweets)
116
117       #Write to file if buffer is full
118       if len(self.harvested_tweets) == self.tweets_pr_file:
119           self.writeout()
```

```
120
121             #Check if limit  has been reached (returning false  cuts  off  listener )
122             return not (len(self.harvested_tweets) == self.max_tweets)
123
124     def writeout(self):
125         filename = self.filename_maker()
126         with open(filename,'w') as outfile:
127             outfile.writelines([json.dumps(t._json)+"\n"
128                             for t in self.harvested_tweets])
129
130         self.harvested_tweets = []
131         self.files_saved += 1
132         #Log event
133         s = "Saved %s files" % self.files_saved
134         self.log(s)
135
136
137 #===============================================================================
138 # Defines stuff  to analyse text  using an already constructed interpretation
139 # matrix.
140 #===============================================================================
141
142 from shared import matrix_dir, row_chunk_size, extensions
143
144 class SemanticAnalyser(object):
145     '''Analyser class using Explicit Semantic Analysis (ESA) to process
146      text fragments. It can compute semantic (pseudo) distance and similarity,
147      as well'''
148     def __init__(self, matrix_filename = 'matrix.mtx', display_concepts = 20):
149         #Number of top concepts to display
150         self.display_concepts = display_concepts
151
152         #Hashes for word and concept indices
153         with open(matrix_dir+'word2index.ind', 'r') as f:
154             self.word2index = shared.load(f)
155         with open(matrix_dir+'concept2index.ind', 'r') as f:
156             self.concept2index = shared.load(f)
157         self.index2concept = {i : c for c, i in self.concept2index.iteritems()}
158
159         #Count number of words and concepts
160         self.n_words = len(self.word2index)
161         self.n_concepts = len(self.concept2index)
162
163     def clean(self, text):
164         text = re.sub('[^\w\s\d\'\-]','', text)
165         text = text.lower()
166
167         return text
168
169     def interpretation_vector(self, text):
170         '''Converts a text fragment string into a row vector where the i'th
171          entry corresponds to the total TF-IDF score of the text fragment
172          for concept i'''
173
174         #Remove mess (quotes, parentheses etc) from text
175         text = self.clean(text)
176
177         #Convert string to hash like  {'word' : no. of occurrences}
178         countmap = Counter(text.split()).iteritems()
179
180         #Interpretation  vector to  be returned
181         result = sps.csr_matrix((1, self.n_concepts), dtype = float)
```

```python
182
183             #Add word count in the correct position of the vector
184             for word, count in countmap:
185                 try:
186                     ind = self.word2index[word]
187                     #Which file to look in
188                     file_number = int(ind/row_chunk_size)
189                     filename = matrix_dir+str(file_number)+extensions['matrix']
190
191                     #And which row to extract
192                     row_number = ind % row_chunk_size
193
194                     #Do it! Do it naw!
195                     with open(filename, 'r') as f:
196                         temp = shared.mload(f)
197                         result = result + count*temp.getrow(row_number)
198                 except KeyError:
199                     pass    #No data on this word -> discard
200
201             #Done. Return row vector as a 1x#concepts CSR matrix
202             return result
203
204     def interpret_text(self, text):
205         '''Attempts to guess the core concepts of the given text fragment'''
206         #Compute the interpretation vector for the text fragment
207         vec = self.interpretation_vector(text)
208
209         #Magic, don't touch
210         top_n = vec.data.argsort()[:len(vec.data)-1-self.display_concepts:-1]
211
212         #List top scoring concepts and their TD-IDF
213         concepts = [self.index2concept[vec.indices[i]] for i in top_n]
214         return concepts
215 #       scores = [vec.data[i] for i in top_n]
216 #       #Return as dict {concept : score}
217 #       return dict(zip(concepts, scores))
218
219     def interpret_file(self, filename):
220         with open(filename, 'r') as f:
221             data = self.clean(f.read())
222         return self.interpret_text(data)
223
224     def interpret_input(self):
225         text = raw_input("Enter text fragment: ")
226         topics = self.interpret_text(text)
227         print "Based on your input, the most probable topics of your text are:"
228         print topics[:self.display_concepts]
229
230     def compare_texts(self, text1, text2):
231         '''Determines cosine similarity between input texts.
232          Returns float in [0,1]'''
233
234         #Determine intepretation vectors
235         v1 = self.interpretation_vector(text1)
236         v2 = self.interpretation_vector(text2)
237
238         #Compute their inner product and make sure it's a scalar
239         dot = v1.dot(v2.transpose())
240         assert dot.shape == (1,1)
241
242         if dot.data:
243             scal = dot.data[0]
```

```
244              else:
245                  scal = 0      #Empty sparse matrix means zero
246
247              #Normalize and return
248              sim = scal/(norm(v1.data)*norm(v2.data))
249              return sim
250
251      def cosine_distance(self, text1, text2):
252          return 1-self.compare_texts(text1, text2)
253
254  if __name__ == '__main__':
255      th = TweetHarvester(verbose=True, max_tweets=10)
256      th.mine('carlsberg', n=10)
257      temp = [t._json for t in th.harvested_tweets if t._json['lang'] == 'en']
258      js = temp[4]
259      with open('tweet_example.json', 'w') as f:
260          pprint(js, stream=f)
261
262  #    if len(sys.argv) > 1:
263  #        fn = sys.argv[1]
264  #    else:
265  #        fn = 'interpret_me.txt'
266  #        with open(fn, 'r') as f:
267  #            data = f.read()
268  #        #
269  #    data = sa.clean(data)
270  #    guesses = sa.interpret_text(data)
271  #
272  #    if len(sys.argv) > 2:
273  #        output_filename = sys.argv[2]
274  #    else:
275  #        output_filename = 'guesses.txt'
276  #    with open(output_filename, 'w') as f:
277  #        for line in guesses:
278  #            f.write(line.encode('utf8'))
279  #            f.write('\n')
```

## A.1.5   Wikicleaner

```
1  # -*- coding: utf-8 -*-
2  import re
3  from htmlentitydefs import name2codepoint
4
5  namespaces = set(['help', 'file talk', 'module', 'topic', 'mediawiki',
6  'wikipedia talk', 'file', 'user talk', 'special', 'category talk', 'category',
7  'media', 'wikipedia', 'book', 'draft', 'book talk', 'template', 'help talk',
8  'timedtext', 'mediawiki talk', 'portal talk', 'portal', 'user', 'module talk',
9  'template talk', 'education program talk', 'education program',
10 'timedtext talk', 'draft talk', 'talk'])
11
12 def dropNested(text, openDelim, closeDelim):
13     '''Helper function to match nested expressions which may cause problems
14      example: {{something something {{something else}} and something third}}
15      cannot be easily matched with a regexp to remove all occurrences.
16      Copied from the WikiExtractor project.'''
17     openRE = re.compile(openDelim)
18     closeRE = re.compile(closeDelim)
```

```python
19          # partition text in separate blocks { } { }
20          matches = []                # pairs (s, e) for each partition
21          nest = 0                    # nesting level
22          start = openRE.search(text, 0)
23          if not start:
24              return text
25          end = closeRE.search(text, start.end())
26          next = start
27          while end:
28              next = openRE.search(text, next.end())
29              if not next:            # termination
30                  while nest:         # close all pending
31                      nest -=1
32                      end0 = closeRE.search(text, end.end())
33                      if end0:
34                          end = end0
35                      else:
36                          break
37                  matches.append((start.start(), end.end()))
38                  break
39              while end.end() < next.start():
40                  # { } {
41                  if nest:
42                      nest -= 1
43                      # try closing more
44                      last = end.end()
45                      end = closeRE.search(text, end.end())
46                      if not end:     # unbalanced
47                          if matches:
48                              span = (matches[0][0], last)
49                          else:
50                              span = (start.start(), last)
51                          matches = [span]
52                          break
53                  else:
54                      matches.append((start.start(), end.end()))
55                      # advance start, find next close
56                      start = next
57                      end = closeRE.search(text, next.end())
58                      break           # { }
59              if next != start:
60                  # { { }
61                  nest += 1
62          # collect text outside partitions
63          res = ''
64          start = 0
65          for s, e in matches:
66              res += text[start:s]
67              start = e
68          res += text[start:]
69          return res
70
71  def unescape(text):
72      '''Removes HTML or XML character references and entities
73       from a text string.
74       @return nice text'''
75      def fixup(m):
76          text = m.group(0)
77          code = m.group(1)
78          return text
79          try:
80              if text[1] == "#": # character reference
```

```python
                    if text[2] == "x":
                        return unichr(int(code[1:], 16))
                    else:
                        return unichr(int(code))
                else:                    # named entity
                    return unichr(name2codepoint[code])
            except UnicodeDecodeError:
                return text # leave as is

        return re.sub("&#?(\w+);", fixup, text)

def drop_spans(matches, text):
    """Drop from text the blocks identified in matches"""
    matches.sort()
    res = ''
    start = 0
    for s, e in matches:
        res += text[start:s]
        start = e
    res += text[start:]
    return res

###Compile regexps for text cleanup:
#Construct patterns for elements to be discarded:
discard_elements = set([
        'gallery', 'timeline', 'noinclude', 'pre',
        'table', 'tr', 'td', 'th', 'caption',
        'form', 'input', 'select', 'option', 'textarea',
        'ul', 'li', 'ol', 'dl', 'dt', 'dd', 'menu', 'dir',
        'ref', 'references', 'img', 'imagemap', 'source'
        ])
discard_element_patterns = []
for tag in discard_elements:
    pattern = re.compile(r'<\s*%s\b[^>]*>.*?<\s*/\s*%s>' % (tag, tag), re.DOTALL | re.IGNORECASE)
    discard_element_patterns.append(pattern)

#Construct patterns to recognize HTML tags
selfclosing_tags = set([ 'br', 'hr', 'nobr', 'ref', 'references' ])
selfclosing_tag_patterns = []
for tag in selfclosing_tags:
    pattern = re.compile(r'<\s*%s\b[^/]*/\s*>' % tag, re.DOTALL | re.IGNORECASE)
    selfclosing_tag_patterns.append(pattern)

#Construct patterns for tags to be ignored
ignored_tags = set([
        'a', 'b', 'big', 'blockquote', 'center', 'cite', 'div', 'em',
        'font', 'h1', 'h2', 'h3', 'h4', 'hiero', 'i', 'kbd', 'nowiki',
        'p', 'plaintext', 's', 'small', 'span', 'strike', 'strong',
        'sub', 'sup', 'tt', 'u', 'var',
])
ignored_tag_patterns = []
for tag in ignored_tags:
    left = re.compile(r'<\s*%s\b[^]*>' % tag, re.IGNORECASE)
    right = re.compile(r'<\s*/\s*%s>' % tag, re.IGNORECASE)
    ignored_tag_patterns.append((left, right))

#Construct patterns to recognize math and code
placeholder_tags = {'math':'formula', 'code':'codice'}
placeholder_tag_patterns = []
for tag, repl in placeholder_tags.items():
    pattern = re.compile(r'<\s*%s(\s*| [^>]+?)>.*?<\s*/\s*%s\s*>' % (tag, tag), re.DOTALL | re.IGNORECASE)
    placeholder_tag_patterns.append((pattern, repl))
```

```python
143
144    #HTML comments
145    comment = re.compile(r'<!--.*?-->', re.DOTALL)
146
147    #Wikilinks
148    wiki_link = re.compile(r'\[\[([^[]]*?)(?:\|([^[]]*?))?\]\](\w*)')
149    parametrized_link = re.compile(r'\[\[.*?\]\]')
150
151    #External links
152    externalLink = re.compile(r'\[\w+.*? (.*?)\]')
153    externalLinkNoAnchor = re.compile(r'\[\w+[&\]]*\]')
154
155    #Bold/italic  text
156    bold_italic = re.compile(r"'''''([^']*?)'''''")
157    bold = re.compile(r"'''(.*?)'''")
158    italic_quote = re.compile(r"'\"(.*?)\"'")
159    italic = re.compile(r"''([^']*)''")
160    quote_quote = re.compile(r'""(.*?)""')
161
162    #Spaces
163    spaces = re.compile(r' {2,}')
164
165    #Dots
166    dots = re.compile(r'\.{4,}')
167
168    #Sections
169    section = re.compile(r'(==+)\s*(.*?)\s*\1')
170
171    # Match preformatted lines
172    preformatted = re.compile(r'^ .*?$', re.MULTILINE)
173
174    #Wikilinks
175    def make_anchor_tag(match):
176        '''Recognizes links and returns only their anchor. Example:
177         <a href="www.something.org">Link text</a> -> Link text'''
178        link = match.group(1)
179        colon = link.find(':')
180        if colon > 0 and link[:colon] not in namespaces:
181            return ''
182        trail = match.group(3)
183        anchor = match.group(2)
184        if not anchor:
185            if link[:colon] in namespaces:
186                return '' #Don't keep stuff like "category: shellfish "
187            anchor = link
188        anchor += trail
189        return anchor
190
191    def clean(text):
192        '''Outputs an article in plaintext from its format in the raw xml dump.'''
193        # Drop transclusions (template, parser functions)
194        # See: http :// www.mediawiki.org/wiki/Help:Templates
195        text = dropNested(text, r'{{', r'}}')
196        # Drop tables
197        text = dropNested(text, r'{\|', r'\|}')
198
199        # Convert wikilinks links  to  plaintext
200        text = wiki_link.sub(make_anchor_tag, text)
201        # Drop remaining links
202        text = parametrized_link.sub('', text)
203
204        # Handle external links
```

```
205    text = externalLink.sub(r'\1', text)
206    text = externalLinkNoAnchor.sub('', text)
207
208    #Handle text formatting
209    text = bold_italic.sub(r'\1', text)
210    text = bold.sub(r'\1', text)
211    text = italic_quote.sub(r'&quot;\1&quot;', text)
212    text = italic.sub(r'&quot;\1&quot;', text)
213    text = quote_quote.sub(r'\1', text)
214    text = text.replace("'''", '').replace("''", '&quot;')
215
216    ############### Process HTML ##############
217
218    # turn into HTML
219    text = unescape(text)
220
221    # do it again (&amp;nbsp;)
222    text = unescape(text)
223
224    # Collect spans
225
226    matches = []
227    # Drop HTML comments
228    for m in comment.finditer(text):
229            matches.append((m.start(), m.end()))
230
231    # Drop self-closing tags
232    for pattern in selfclosing_tag_patterns:
233        for m in pattern.finditer(text):
234            matches.append((m.start(), m.end()))
235
236    # Drop ignored tags
237    for left, right in ignored_tag_patterns:
238        for m in left.finditer(text):
239            matches.append((m.start(), m.end()))
240        for m in right.finditer(text):
241            matches.append((m.start(), m.end()))
242
243    # Bulk remove all spans
244    text = drop_spans(matches, text)
245
246    # Cannot use dropSpan on these since they may be nested
247    # Drop discarded elements
248    for pattern in discard_element_patterns:
249        text = pattern.sub('', text)
250
251    # Expand placeholders
252    for pattern, placeholder in placeholder_tag_patterns:
253        index = 1
254        for match in pattern.finditer(text):
255            text = text.replace(match.group(), '%s_%d' % (placeholder, index))
256            index += 1
257
258    ##############################################
259
260    # Drop preformatted
261    # This can't be done before since it may remove tags
262    text = preformatted.sub('', text)
263
264    # Cleanup text
265    text = text.replace('\t', ' ')
266    text = spaces.sub(' ', text)
```

```
267    text = dots.sub('...', text)
268    text = re.sub(u' (,:\.\)\]»)', r'\1', text)
269    text = re.sub(u'(\[\(«) ', r'\1', text)
270    text = re.sub(r'\n\W+?\n', '\n', text) # lines with only punctuations
271    text = text.replace(',,', ',').replace(',.', '.')
272
273    #Handle section headers, residua etc.
274    page = []
275    headers = {}
276    empty_section = False
277
278    for line in text.split('\n'):
279
280        if not line:
281            continue
282        # Handle section titles
283        m = section.match(line)
284        if m:
285            title = m.group(2)
286            lev = len(m.group(1))
287            if title and title[-1] not in '!?':
288                title += '.'
289            headers[lev] = title
290            # drop previous headers
291            for i in headers.keys():
292                if i > lev:
293                    del headers[i]
294            empty_section = True
295            continue
296        # Handle page title
297        if line.startswith('++'):
298            title = line[2:-2]
299            if title:
300                if title[-1] not in '!?':
301                    title += '.'
302                page.append(title)
303        # handle lists
304        elif line[0] in '*#:;':
305            continue
306        # Drop residuals of lists
307        elif line[0] in '{|' or line[-1] in '}':
308            continue
309        # Drop irrelevant lines
310        elif (line[0] == '(' and line[-1] == ')') or line.strip('.-') == '':
311            continue
312        elif len(headers):
313            items = headers.items()
314            items.sort()
315            for (i, v) in items:
316                page.append(v)
317            headers.clear()
318            page.append(line)  # first line
319            empty_section = False
320        elif not empty_section:
321            page.append(line)
322
323    text = ''.join(page)
324
325    #Remove quote tags.
326    text = text.replace("&quot;", '')
327
328    #Get rid of parentheses, punctuation and the like
```

```
329    text = re.sub('[^\w\s\d\'\-]','', text)
330    return text
```

## A.2   Social Fabric-related Code

This section contains the code referred to in part II of the thesis.

### A.2.1   Phonetools

```python
# -*- coding: utf-8 -*-
"""
Created on Tue Dec 02 15:12:29 2014

@author: Bjarke
"""

import datetime
from pkg_resources import resource_filename
from ast import literal_eval as LE
import numpy as np
import json
import os

def _global_path(path):
    '''Helper method to ensure that data files belonging to the social_fabric
    module are available both when importing the module and when running
    individual parts from it for testing.
    usage: Always use _global_path(somefile) rather than somefile'''
    if __name__ == '__main__':
        return path
    else:
        return resource_filename('social_fabric', path)

def make_filename(prefix, interval, bin_size, ext=''):
    '''Returns a filename like "call_out-intervalsize-binsize_N", where N is
    an int so files aren't accidentally overwritten.
    I've been known to do that.'''
    stem = "%s-int%s-bin%s" % (prefix, interval, bin_size)
    n = 0
    ext = '.'+ext
    attempt = stem+ext
    while os.path.isfile(attempt):
        n+=1
        attempt = stem+"_"+str(n)+ext
    return attempt

def unix2str(unixtime):
    '''Converts timestamp to datetime object'''
    dt = datetime.datetime.fromtimestamp(unixtime)
    return str(dt)

#Converts bluetooth MAC-addresses to users
with open(_global_path('user_mappings/bt_to_user.txt'), 'r') as f:
    bt2user = LE(f.read())

#Converts phone 'number' code to users
with open(_global_path('user_mappings/phonenumbers.txt'), 'r') as f:
    number2user = LE(f.read())

#Converts IDs from psychological profile to users
```

```python
52   with open(_global_path('user_mappings/user_mapping.txt'), 'r') as f:
53       psych2user = {}
54       #This is in tab-separated values, for some reason.
55       for line in f.read().splitlines():
56           (ID, usercode) = line.split('\t')
57           assert ID.startswith('user_')
58           psych2user[ID] = usercode
59
60   #Converts users to info on their psych profiles
61   with open(_global_path('user_mappings/user2profile.json'), 'r') as f:
62       user2profile = json.load(f)
63
64   def is_valid_call(call_dict):
65       '''Determine whether an entry is 'valid', i.e. make sure user isn't
66        calling/texting themselves, which people apparently do...'''
67       caller = call_dict['user']
68       try:
69           receiver = number2user[call_dict['number']]
70       except KeyError:
71           receiver = None
72       return caller != receiver
73
74   def readncheck(path):
75       '''Reads in all valid call info from the file at path'''
76       try:
77           with open(path,'r') as f:
78               raw = [LE(line) for line in f.readlines()]
79       except IOError:
80           return []   #no file :(
81       #File read. Return proper calls
82       return [call for call in raw if is_valid_call(call)]
83
84   class Binarray(list):
85       '''Custom array type to automatically bin the time around a set center
86        and place elements in each bin.
87        The array can be centered using Binarray.center = <some time>.
88        After centering, timestamps can be placed in bins around the center with
89        Binarray.place_event(<some other time>).'''
90
91       def __init__(self, interval = 3*60**2, bin_size = 3*60, center = None,
92                   initial_values = None):
93           '''
94           Args:
95           ----------------------
96           interval : int
97             Total number of seconds covered by the Binarray.
98
99           bin_size : int
100            Width of each bin measured in seconds. The total interval must be
101            an integer multiplum of the bin size.
102
103          center : int
104            Where to center the Binarray. If the array is centered at time t,
105            any event placed in it will placed in a bin depending on how long
106            before or after t the event occured.
107
108          initial_values : list
109            List of values to start the Binarray with. Default is zeroes.'''
110
111
112          #Make sure interval is an integer multiplum of bins
113          if not interval % bin_size == 0:
```

```python
114             suggest = interval − interval % bin_size
115             error = "Interval isn't an integer multiple of bin size. \
116               Consider changing interval to %s." % suggest
117             raise ValueError(error)
118
119         #Set parameters
120         self.bin_size = bin_size
121         self.interval = interval
122         self.size = 2*int(interval/bin_size)
123         self.centerindex = int(self.size//2)
124         self.center = center
125         #Keep track of how many events missed the bins completely
126         self.misses = 0
127         #Call parent constructor
128         if not initial_values:
129             startlist = [0]*self.size
130         else:
131             if not len(initial_values) == self.size:
132                 msg = '''Array of start value must have length %d. Tried to
133                   instantiate with length of %d.''' % (self.size,
134                                                     len(initial_values))
135                 raise ValueError(msg)
136             startlist = initial_values
137         super(Binarray, self).__init__(startlist)
138
139
140     def place_event(self, position):
141         '''Places one count in the appropriate bin if event falls within
142          <interval> of <center>. Returns True on success.'''
143         if self.center == None:
144             raise TypeError('Center must be set!')
145         delta = position − self.center
146         #Check if event is outside current interval
147         if np.abs(delta) >= self.interval:
148             self.misses += 1
149             return False
150         #Woo, we're in the correct interval
151         index = int(delta//self.bin_size) #relative to middle of array
152         self[self.centerindex + index] += 1
153         return True
154
155     def normalized(self):
156         '''Returns a normalized copy of the array's contents.'''
157         events = sum(self) + self.misses
158         #Use numpy vectorized function for increased speed.
159         f = np.vectorize(lambda x: x*1.0/events if events else 0)
160         return f(self)
161
162     def _todict(self):
163         '''Helper method to allow dumping to JSON format.'''
164         attrs = ['misses', 'interval', 'bin_size']
165         d = {att : self.__getattribute__(att) for att in attrs}
166         d['values'] = list(self)
167         d['type'] = 'binarray'
168         return d
169
170
171 def _dumphelper(obj):
172     '''Evil recursive helper method to convert various nested objects to
173      a JSON-serializeable format.
174      This should only be called by the dump method!'''
175     if isinstance(obj, Binarray):
```

```python
176            d = obj._todict()
177            return _dumphelper(d)
178        elif isinstance(obj, tuple):
179            hurrayimhelping = [_dumphelper(elem) for elem in obj]
180            return {'type' : 'tuple', 'values' : hurrayimhelping}
181        elif isinstance(obj, dict):
182            temp = {'type' : 'dict'}
183            contents = [{'key' : _dumphelper(key), 'value' : _dumphelper(value)}
184                        for key, value in obj.iteritems()]
185            temp['contents'] = contents
186            return temp
187        #Do nothing if obj is an unrecognized type. Let JSON raise errors.
188        else:
189            return obj
190
191    def _hook(obj):
192        '''Evil recursive object hook method to reconstruct various nested
193        objects from a JSON dump.
194        This should only be called by the load method!'''
195        if isinstance(obj, (unicode, str)):
196            try:
197                return _hook(LE(obj))
198            except ValueError: #happens for simple strings that don't need eval
199                return obj
200        elif isinstance(obj, dict):
201            if not 'type' in obj:
202                raise KeyError('Missing type info')
203            if obj['type'] == 'dict':
204                contents = obj['contents']
205                d = {_hook(e['key']) : _hook(e['value']) for e in contents}
206                #Make sure we also catch nested expressions
207                if 'type' in d:
208                    return _hook(d)
209                else:
210                    return d
211            elif obj['type'] == 'binarray':
212                instance = Binarray(initial_values = obj['values'],
213                                    bin_size = obj['bin_size'],
214                                    interval = obj['interval'])
215                instance.misses = obj['misses']
216                for key, val in obj.iteritems():
217                    if key == 'values':
218                        continue
219                    instance.__setattr__(key, val)
220                return instance
221
222            elif obj['type'] == 'tuple':
223                #Hook elements individually, then convert back to tuple
224                restored = [_hook(elem) for elem in obj['values']]
225                return tuple(restored)
226            else:
227                temp = {}
228                for k, v in obj.iteritems():
229                    k = _hook(k)
230                    temp[k] = _hook(v)
231                return temp
232            #
233        #Do nothing if obj is an unrecognized type
234        else:
235            return obj
236
237    def load(file_handle):
```

```python
238        '''Reads in json serialized nested combinations of dicts, binarrays
239         and tuples.'''
240        temp = json.load(file_handle, encoding='utf-8')
241        return _hook(temp)
242
243    def dump(obj, file_handle):
244        '''json serializes nested combinations of dicts, binarrays
245         and tuples.'''
246        json.dump(_dumphelper(obj), file_handle, indent=4, encoding='utf-8')
247
248
249    if __name__=='__main__':
250        from time import time
251        from random import randint
252        #Create Binarray with interval +/- one hour and bin size ten minutes.
253        ba = Binarray(interval = 60*60, bin_size = 10*60)
254        #Center it on the present
255        now = int(time())
256        ba.center = now
257        #Generate some timestamps around the present
258        new_times = [now + randint(-60*60, 60*60) for _ in xrange(100)]
259        for tt in new_times:
260            ba.place_event(tt)
261
262        #Save it
263        with open('filename.sig', 'w') as f:
264            dump(ba, f)
265
266        print ba
```

### A.2.2  Code Communication Dynamics

### Code for extracting Bluetooth signals

```python
# −∗− coding: utf−8 −∗−
from __future__ import division

import os
import glob
import itertools
import random
from ast import literal_eval as LE
from social_fabric.phonetools import readncheck, Binarray, dump, make_filename

#============================================================================
# Parameters
#============================================================================

#Grap all the files we need to read
userfile_path = "userfiles/"  #linux


#Interval of interest and bin size (seconds)
interval = 12*60**2
bin_size = 10*60

#How many other users contribute to background for a given user
number_of_background_samples = 1

#Number of users to analyse. Use None or 0 to include everyone.
max_users = 2

#Number of repeated signals required to be considered social
social_threshold = 2


#============================================================================
# Functions
#============================================================================

def call_type(n):
    if n == 1:
        return 'call_in'
    elif n == 2:
        return 'call_out'
    else:
        return None

def text_type(n):
    if n == 1:
        return 'text_in'
    elif n == 2:
        return 'text_out'
    else:
        return None


def user2social_times(user):
    '''Converts user code to a list of times when the user was social i.e. had
    two or more repeated Bluetooth signals.'''
```

```
56        with open(userfile_path+user+"/bluetooth_log.txt", 'r') as f:
57            raw = [LE(line) for line in f.readlines()]
58    #List of times when user was social
59    social_times = []
60    #Temporary variables
61    current_time = 0
62    previous_time = 0
63    current_users = []
64    previous_users = []
65    for ind in xrange(len(raw)-1):
66        signal = raw[ind]
67        new_time = signal['timestamp']
68        #Check if line represents a new signal and if so, update values
69        if new_time != current_time:
70            #Determine if previous signal was social and append to results
71            overlap = set(previous_users).intersection(set(current_users))
72            if len(overlap) >= social_threshold:
73                social_times.append(previous_time)
74            #Update variables
75            previous_time = current_time
76            previous_users = current_users
77            current_users = []
78            current_time = new_time
79        new_user = signal['name']
80        if new_user=='-1' or not new_user:
81            continue
82        else:
83            current_users.append(new_user)
84        #
85    return social_times
86
87    #==============================================================================
88    # Time to crunch some numbers
89    #==============================================================================
90
91    user_folders = [f for f in glob.glob(userfile_path+"/*")
92                        if os.path.isfile(f+"/call_log.txt")
93                        and os.path.isfile(f+"/sms_log.txt")
94                        and os.path.isfile(f+"/bluetooth_log.txt")]
95    users = [folder.split(userfile_path)[1] for folder in user_folders]
96
97    if not users:
98        raise IOError('Found no users. Check userfile path.')
99
100   if max_users:
101       users = users[:max_users]
102
103   trigger = 'bluetooth'
104   events = ['call_out', 'call_in', 'text_in', 'text_out']
105   pairs = [p for p in itertools.product([trigger], events)]
106
107   activity = {p : Binarray(interval,bin_size) for p in pairs}
108   background = {p : Binarray(interval,bin_size) for p in pairs}
109
110
111   #Read in data
112   call_data = {user : readncheck(userfile_path+user+"/call_log.txt")
113               for user in users}
114   text_data = {user : readncheck(userfile_path+user+"/sms_log.txt")
115               for user in users}
116
117   count = 0
```

```
118
119  for user in users:
120      count += 1
121      print "Analyzing user %s out of %s. Code: %s" % (count, len(users), user)
122
123      #Get user data
124      user_calls = call_data[user]
125      user_texts = text_data[user]
126      user_social_times = user2social_times(user)
127
128      if not user_social_times:
129          continue
130
131      #Get background data
132      others = []
133      other_social_times = []
134      while len(others) < number_of_background_samples:
135          temp = random.choice(users)
136          if not (temp in others or temp == user):
137              newstuff = user2social_times(temp)
138              if not newstuff:
139                  continue
140              others.append(temp)
141              other_social_times += newstuff
142
143      #Determine the interval in  which we have data on the current user
144      first = min(user_social_times)
145      last = max(user_social_times)
146
147      #=================================================================================
148      #      #Establish  activity  signal
149      #=================================================================================
150      for time in user_social_times:
151          for e in events:
152              activity[(trigger, e)].center = time
153
154          for user_call in user_calls:
155              event = call_type(user_call['type'])
156              if not event:
157                  continue
158              time = user_call['timestamp']
159              activity[(trigger, event)].place_event(time)
160
161          for user_text in user_texts:
162              event = text_type(user_text['type'])
163              if not event:
164                  continue
165              time = user_text['timestamp']
166              activity[(trigger, event)].place_event(time)
167
168      #=================================================================================
169      #      #Establish  background  signal
170      #=================================================================================
171      for other_time in other_social_times:
172          #Reposition the relevant  binarrays
173          if not first <= other_time <= last:
174              continue
175          for e in events:
176              background[(trigger, e)].center = other_time
177
178          #Determine call background
179          for user_call in user_calls or []:
```

```
180              event = call_type(user_call['type'])
181              if not event:
182                  continue
183              time = user_call['timestamp']
184              background[(trigger,event)].place_event(time)
185
186          #Determine text background
187          for user_text in user_texts or []:
188              event = text_type(user_text['type'])
189              if not event:
190                  continue
191              time = user_text['timestamp']
192              background[(trigger,event)].place_event(time)
193          #
194      #
195  #
196
197  #==============================================================================
198  # Done. Save signals
199  #==============================================================================
200
201  #Make a filename for the output file
202  filename = make_filename(prefix = trigger, interval=interval,
203                           bin_size=bin_size, ext = 'json')
204  with open(filename, 'w') as f:
205      dump((activity, background), f)
206
207  print "Saved to "+filename
```

## Code for loading and plotting Bluetoot data

```python
1   # −∗− coding: utf−8 −∗−
2   from __future__ import division
3
4   import glob
5   import numpy as np
6   import matplotlib.pyplot as plt
7
8   from social_fabric.phonetools import (load, make_filename)
9
10  #=================================================================================
11  display = True
12  #These are just  default − they're  updated when data is read
13  interval = 12∗60∗∗2
14  bin_size = 10∗60
15  #=================================================================================
16
17
18  #Numpy−compliant method to convert activity and background to relative signal
19  get_signal = np.vectorize(lambda act, back : act ∗ 1.0/back if back else 0)
20
21  def read_data(filename):
22      '''Reads in data to plot. Updates the interval and bin sizes parameters.
23       Returns data to be plotted as a hashmap with the following structure:
24      {trigger : {event : signal}}.'''
25      data = {}
26      with open(filename, 'r') as f:
27          (a,b) = load(f)
28      first = True
29      for key in a.keys():
30          #Update interval and bin info
31          if first:
32              global interval
33              global bin_size
34              interval = a[key].interval
35              bin_size = a[key].bin_size
36              first = False
37
38          (trigger, event) = key
39          act = a[key].normalized()
40          back = b[key].normalized()
41          signal = get_signal(list(act), list(back))
42          if not trigger in data:
43              data[trigger] = {}
44          data[trigger][event] = signal
45      return data
46
47
48  def make_plot(trigger, signals):
49      '''Generate a plot of relative user activity signal.'''
50      legendstuff = []
51      for event, vals in signals.iteritems():
52          t = np.arange(−interval, interval, bin_size)
53          t = map(lambda x: x∗1.0/3600,  t)
54          legendstuff.append(event)
55          s = list(vals)
56          plt.plot(t, s)
57      plt.legend(legendstuff, loc='upper left')
58      plt.xlabel('Time (h)')
59      plt.ylabel('Relative signal')
60      plt.title("Trigger: " + trigger)
61      plt.grid(True)
```

```python
62        saveas = make_filename(trigger, interval, bin_size, ext = 'pdf')
63        plt.savefig(saveas)
64        print "Saved to "+saveas
65        if display:
66            plt.show()
67
68 if __name__ == '__main__':
69     # Read in data
70     filenames = glob.glob('bluetooth-int43200-bin600.json')
71     for filename in filenames:
72         data = read_data(filename)
73         for trigger, signals in data.iteritems():
74             make_plot(trigger, signals)
75             #
```

### A.2.3   Preprocessing

```python
# -*- coding: utf-8 -*-
from __future__ import division

import os
import glob
import numpy as np
from ast import literal_eval as LE
from social_fabric.phonetools import readncheck, user2profile
from social_fabric.secrets import pushme
from social_fabric.smallestenclosingcircle import make_circle
from social_fabric import lloyds
from collections import Counter
import math
from datetime import datetime, timedelta
import pytz
import json
from statsmodels.tsa import ar_model
import multiprocessing


#===============================================================================
# Parameters
#===============================================================================

#Grap all the files we need to read
#userfile_path = "c :\\ userfiles \\"    #@Windows
userfile_path = "/lscr_paper/amoellga/Data/Telefon/userfiles/" #linux

#Filename to save output to.
output_filename = 'data.json'

#Number of users to analyse. Use None or 0 to include everyone.
max_users = None

#Specific user codes to analyze for debugging purposes. Empty means include all
exclusive_users = []
#Den her driller :  '28b76d7b7879d364321f164df5169f'

#Conversion factors from degrees to meters (accurate around Copenhagen)
longitude2meters = 111319
latitude2meters = 110946

#In meters. These improve convergence of the stochastic SEC algorithm.
x_offset = 1389425.2238223257
y_offset = 6181209.0059678229

required = ['bluetooth_log.txt', 'call_log.txt', 'facebook_log.txt',
            'gps_log.txt', 'sms_log.txt']

user_folders = [f for f in glob.glob(userfile_path+"*")
                if all(os.path.isfile(f+"/"+stuff) for stuff in required)]
users = [folder.split(userfile_path)[1] for folder in user_folders]

if not users:
    raise IOError('Found no users. Check userfile path.')

if exclusive_users:
    users = list(set(exclusive_users).intersection(users))
```

```
58
59   if max_users:
60       users = users[:max_users]
61
62   # Number of wallclock hours pr bin when fitting autoregressive series
63   hours_pr_ar_bin = 6
64   assert 24%hours_pr_ar_bin == 0
65
66   # Number of hours pr bin when compute daily rythm entropy
67   hours_pr_daily_rythm_bin = 1
68   assert 24%hours_pr_daily_rythm_bin == 0
69
70   # Time zone information
71   cph_tz = pytz.timezone('Europe/Copenhagen')
72
73   n_jobs = 16  #maximum number of processors to use.
74
75   #Which data kinds to include
76   include_calls = True
77   include_ar = True
78   include_gps = True
79   include_network = False  #Allow geolocation from network data
80   include_bluetooth = True
81   include_facebook = True
82
83   #Whether to include not−a−number values in final output
84   allow_nan = True
85
86   #Threshold values to discard users with insufficient data
87   minimum_number_of_texts = 10
88   minimum_number_of_calls = 5
89   minimum_number_of_gps_points = 100
90   minimum_number_of_facebook_friends = 1
91
92   #N_pings required to be considered social
93   bluetooth_social_threshold = 2
94
95   #Whether to output plots of cluster analysis
96   plot_clusters = False
97
98   #==============================================================================
99   # Define helper methods
100  #==============================================================================
101
102  def get_distance(p, q):
103      return math.sqrt(sum([(p[i]−q[i])**2 for i in xrange(len(p))]))
104
105  def is_sorted(l):
106      return all([l[i+1] >= l[i] for i in xrange(len(l)−1)])
107
108  def get_entropy(event_list):
109      '''Takes a list of contacts from in/ or outgoing call/text events and
110       computes its entropy. event_list must be simply a list of user codes
111       corresponding to events.'''
112      n = len(event_list)
113      counts = Counter(event_list)
114      ent = sum([−v/n * math.log(v/n, 2) for v in counts.values()])
115      return ent
116
117  def next_time(dt, deltahours):
118      '''Accepts a datetime object and returns the next datetime object at which
119       the 'hour' count modulo deltahours is zero.
```

```python
120        For example, deltahours = 6 gives the next time clock time is
121        0, 6, 12 or 18.
122        Sounds simple but is pretty annoying due to daylight saving time and so on,
123        so take care not to mess with this.'''
124     base = datetime(dt.year, dt.month, dt.day, dt.hour)
125     interval = timedelta(hours = deltahours − dt.hour%deltahours)
126     naive_guess = base + interval
127     return cph_tz.localize(naive_guess)
128
129 def next_midnight(dt):
130     '''Takes a datetime object and returns dt object of following midnight'''
131     base = datetime(dt.year, dt.month, dt.day)
132     naive = base + timedelta(days = 1)
133     return cph_tz.localize(naive)
134
135 def epoch2dt(timestamp):
136     '''Converts unix timestamp into a pytz timezone-aware datetimeobject.'''
137     utc_time = datetime.utcfromtimestamp(timestamp)
138     smart_time = cph_tz.fromutc(utc_time)
139     return smart_time
140
141 def sort_dicts_by_key(dictlist, key):
142     '''Takes a list of dicts and returns the same list sorted by its
143      key-entries.'''
144     decorated = [(d[key], d) for d in dictlist]
145     decorated.sort()
146     return [d for (k, d) in decorated]
147
148 def get_autocovar_coefficient(X, lag):
149     '''Returns the autocovariance coefficient for the input series at the
150      input lag.'''
151     mu = np.mean(X)
152     temp = sum((X[i] − mu)*(X[i+lag]−mu) for i in xrange(len(X)−lag))
153     return temp*1.0/len(X)
154
155 def get_autocorrelation_coefficients(series, lags):
156     '''Determines the autocorrelation coefficients of input series at
157      each of the input lags. Uses .r_k = c_k/c_0.
158      Accepts a list of lags or an int in which case it returns lags up to
159      and including the input.'''
160     if isinstance(lags, int):
161         lags = range(lags+1)
162     c0 = get_autocovar_coefficient(series, 0)
163     inv = 1.0/c0
164     return [inv*get_autocovar_coefficient(series, lag) for lag in lags]
165
166 def make_time_series(dts, hours_pr_ar_bin):
167     '''Takes a sorted list of datetime objects and converts to a time series
168      where each entry denotes the number of events in the corresponding bin.'''
169     first_time = next_time(dts[0], hours_pr_ar_bin)
170     for i in xrange(len(dts)):
171         if dts[i] >= first_time:
172             dts = dts[i:]
173             break
174         #
175     last_time = next_time(first_time, hours_pr_ar_bin)
176     time_series = []
177
178     summer = 0
179     for dt in dts:
180         while not first_time <= dt < last_time:
181             time_series.append(summer)
```

```
182              summer = 0
183              first_time = last_time
184              last_time = next_time(last_time, hours_pr_ar_bin)
185          summer += 1
186      return time_series
187
188
189  def timestamps2daily_entropy(timestamps, hours_pr_bin):
190          '''Constructs a histogram of hour-values of the imput timestamps
191           and computes its entropy.'''
192          if not 24%hours_pr_bin == 0:
193              raise ValueError("24 must be divisible by hours_pr_bin.")
194          bins = {}
195          for timestamp in timestamps:
196              hour = epoch2dt(timestamp).hour
197              _bin = int(hour/hours_pr_bin)
198              try:
199                  bins[_bin] += 1
200              except KeyError:
201                  bins[_bin] = 1
202              #
203          total = sum(bins.values())
204          entropy = sum([-v/total*math.log(v/total, 2) for v in bins.values()])
205          return entropy
206
207  #Make sure we have a blank file to write to.
208  open(output_filename, 'w').close()
209  assert os.stat(output_filename).st_size == 0  #Check that it worked.
210
211  '''Main processing method. This is written as a separate method to allow
212  easy multiprocessing.'''
213  def process_user(user, user_counter):
214      msg = "Processing user %s of %s: %s" % (user_counter, len(users), user)
215      print msg
216      if user_counter % 100 == 0:
217          pushme(msg)
218
219      #Dict to hold all data extracted on current user
220      data = {}
221
222      #Try to load user psych profile- Discard user if they're not in file
223      try:
224          profile = user2profile[user]
225      except KeyError:
226          print "No psychological data on user."
227          return None  # No questionaire data. Shouldn't happen.
228
229      #Read in calls/texts
230      calls = readncheck(userfile_path+user+"/call_log.txt")
231
232      if len(calls) < minimum_number_of_calls:
233          print "too few calls"
234          return None
235      texts = readncheck(userfile_path+user+"/sms_log.txt")
236      if len(texts) < minimum_number_of_texts:
237          print "too few texts"
238          return None
239
240      #Extract list of times for calls, texts and combination
241      call_times = sorted([d['timestamp'] for d in calls])
242      text_times = sorted([d['timestamp'] for d in texts])
243      call_text_times = sorted(call_times + text_times)
```

```
244
245        #Get calls  from the  first  three  months
246        tmin = call_times[0]
247        tmax = call_times[0] + 60*60*24*30*3
248        early_calls = [c['number'] for c in calls if tmin<=c['timestamp']<=tmax]
249
250        #Repeat for texts
251        tmin = text_times[0]
252        tmax = text_times[0] + 60*60*24*30*3
253        early_texts = [t['body'] for t in texts if tmin<=t['timestamp']<=tmax]
254
255        #Get number of unique contacts for first  3 months and append to data.
256        uniques = len(set(early_calls + early_texts))
257        data['n_contacts_first_three_months'] = uniques
258
259        #Compute daily entropy for calls and texts
260        data['call_daily_entropy'] = timestamps2daily_entropy(call_times,
261                                               hours_pr_daily_rythm_bin)
262        data['text_daily_entropy'] = timestamps2daily_entropy(text_times,
263                                               hours_pr_daily_rythm_bin)
264
265        #Compute median and std for call durations
266        call_durations = [c['duration'] for c in calls if not c['duration'] == 0]
267        data['call_duration_med'] = np.median(call_durations)
268        data['call_duration_std'] = np.std(call_durations)
269
270  #================================================================================
271  # Crunch time−series info
272  #================================================================================
273      if include_ar:
274          #Grap a sorted  list  of only times of  events caused by user
275          outgoing_stuff = sorted([d['timestamp'] for d in calls
276                              if d['type'] == 2] + [d['timestamp']
277                              for d in texts if d['type'] == 2])
278
279          n_params = int(24*7/hours_pr_ar_bin + 1) #1 week plus 1 extra bin
280
281          #−Fit time series  and extract  parameters
282          try:
283              #Convert into timezone aware datetime objects
284              dts = [epoch2dt(timestamp) for timestamp in outgoing_stuff]
285
286              time_series = make_time_series(dts, hours_pr_ar_bin)
287
288              model = ar_model.AR(time_series)
289              result = model.fit(n_params)
290              #Grab parameters from fitted model
291              params = result.params[1:]
292              while len(params) < n_params:
293                  params.append(float('nan'))
294          except:
295              if not allow_nan:
296                  return None
297              else:
298                  params = [float('nan') for _ in xrange(n_params)]
299
300          #Append AR−coefficients to user data
301          count = 0
302          for par in params:
303              count += 1
304              name = "outgoing_activity_AR_coeff_"+str(count)
305              data[name] = par
```

```python
306
307                 #Get autocorrelation  coefficients  as well
308                 try:
309                     accs = get_autocorrelation_coefficients(time_series, n_params)[:1]
310                 except:
311                     if not allow_nan:
312                         return None
313                     accs = [float('nan') for _ in xrange(n_params + 1)]
314
315                 #Append autocorrelation coefficients to user data
316                 for i in xrange(len(accs)):
317                     name = "outgoing_activity_acc_"+str(i)
318                     data[name] = accs[i]
319
320                 #Repeat with incoming signals. Might be interesting.
321                 incoming_stuff = sorted([d['timestamp'] for d in calls
322                                        if d['type'] == 1] + [d['timestamp']
323                                        for d in texts if d['type'] == 1])
324
325                 try:
326                     # Convert into timezone aware datetime objects
327                     dts = [epoch2dt(timestamp) for timestamp in incoming_stuff]
328
329                     time_series = make_time_series(dts, hours_pr_ar_bin)
330                     model = ar_model.AR(time_series)
331                     result = model.fit(n_params)
332                     params = result.params[1:]
333                     while len(params) < n_params:
334                         params.append(float('nan'))
335                 except:
336                     if not allow_nan:
337                         return None
338                     else:
339                         params = [float('nan') for _ in xrange(n_params)]
340
341                 # Name each of them and append to user data
342                 count = 0
343                 for par in params:
344                     count += 1
345                     name = "incoming_activity_AR_coeff_"+str(count)
346                     data[name] = par
347
348                 #Get autocorrelation  coefficients  as well
349                 try:
350                     accs = get_autocorrelation_coefficients(time_series, n_params)[:1]
351                 except:
352                     if not allow_nan:
353                         return None
354                     accs = [float('nan') for _ in xrange(n_params + 1)]
355
356                 for i in xrange(len(accs)):
357                     name = "incoming_activity_acc_"+str(i)
358                     data[name] = accs[i]
359
360
361     #===============================================================================
362     #     Crunch call/text  info
363     #===============================================================================
364
365         if include_calls:
366             #Add values to temporary data map.
367             d = {'call' : call_times, 'text' : text_times, 'ct' : call_text_times}
```

```python
368        for label, times in d.iteritems():
369            timegaps = [times[i+1] − times[i] for i in xrange(len(times)−1)]
370            timegaps = filter(lambda x: x < 259200, timegaps)  #3 days, tops
371            data[label+'_iet_med'] = np.median(timegaps)
372            data[label+'_iet_std'] = np.std(timegaps)
373
374        #Generate lists of the contact for each text/call event
375        call_numbers = [call['number'] for call in calls]
376        text_numbers = [text['address'] for text in texts]
377        ct_numbers = call_numbers + text_numbers
378
379        #Compute entropy and add to data
380        data['call_entropy'] = get_entropy(call_numbers)
381        data['text_entropy'] = get_entropy(text_numbers)
382        data['ct_entropy'] = get_entropy(ct_numbers)
383
384        #Compute contact list info
385        call_contacts = Counter([c['number'] for c in calls
386                                 if c['type'] == 2]).keys()
387        text_contacts = Counter([t['address'] for t in texts
388                                 if t['type'] == 2]).keys()
389        #Grap number of contacts
390        n_call_contacts = len(call_contacts)
391        n_text_contacts = len(text_contacts)
392        n_ct_contacts = len(set(call_contacts).union(set(text_contacts)))
393
394        #Add to data map
395        data['n_call_contacts'] = n_call_contacts
396        data['n_text_contacts'] = n_text_contacts
397        data['n_ct_contacts'] = n_ct_contacts
398
399        #Compute and add contact/interaction ratio (cir)
400        data['call_cir'] = n_call_contacts/len(calls)
401        data['text_cir'] = n_text_contacts/len(texts)
402        data['ct_cir'] = n_ct_contacts/(len(calls) + len(texts))
403
404        #Add data on number of interactions
405        data['n_calls'] = len(calls)
406        data['n_texts'] = len(texts)
407        data['n_ct'] = len(calls + texts)
408
409        #Determine percentage of calls/texts that were initiated by user.
410        initiated_calls = len([c for c in calls if c['type'] == 2])
411        data['call_percent_initiated'] = initiated_calls/len(calls)
412        initiated_texts = len([t for t in texts if t['type'] == 2])
413        data['call_percent_initiated'] = initiated_texts/len(texts)
414
415        #Determine call response rate.
416        with open(userfile_path+user+"/call_log.txt", 'r') as f:
417            all_calls = [LE(line) for line in f.readlines()]
418        #Make sure the call data is sorted
419        if not is_sorted([c['timestamp'] for c in all_calls]):
420            all_calls = sort_dicts_by_key(all_calls, 'timestamp')
421
422        '''Check for unanswered called that are replied to within an hour.
423         This is performed in the following fashion: iterate through all the
424         calls. If a call is unanswered, add it to "holding" list. If a call
425         from holding matches the current call, it counts as a reply.
426         If the time of the current call is more than hour after a held call,
427         it is discarded'''
428        missed = 0
429        replied = 0
```

```
430        holding = []
431        for call in all_calls:
432            if call['type']==3 or call['type']==1 and call['duration']==0:
433                holding.append(call)
434                missed += 1
435            else:
436                for held_call in holding:
437                    #Drop calls that have been held for too long
438                    if call['timestamp'] − held_call['timestamp'] > 3600:
439                        holding.remove(held_call)
440                    #Check if given call is a resonse
441                    elif (call['type'] == 2
442                          and call['number'] == held_call['number']):
443                        holding.remove(held_call)
444                        replied += 1
445                    #
446                #
447            #
448        data['call_response_rate'] = replied/(missed+replied) if replied else 0
449
450        #Determine text response rate
451        missed = 0
452        replied = 0
453        holding = []
454        response_times = []
455        if not is_sorted([t['timestamp'] for t in texts]):
456            texts = sort_dicts_by_key(texts, 'timestamp')
457
458        for text in texts:
459            #Make sure incoming text is not from a user already held
460            if text['type'] == 1:
461                if not holding or all([text['address'] !=
462                                       t['address'] for t in holding]):
463                    #It's good − append it
464                    holding.append(text)
465                    missed += 1
466                #
467            else:
468                for held_text in holding:
469                    if text['timestamp'] − held_text['timestamp'] > 3600:
470                        holding.remove(held_text)
471                    #Check if text counts as reply
472                    elif (text['type'] == 2
473                    and text['address'] == held_text['address']):
474                        holding.remove(held_text)
475                        replied += 1
476                        dt = text['timestamp'] − held_text['timestamp']
477                        response_times.append(dt)
478                    #
479                #
480            #
481        data['text_response_rate'] = replied/(missed+replied) if replied else 0
482        data['text_latency'] = np.median(response_times)
483
484        #Check percentage of calls taken place in during the night
485        count = 0
486        for call in calls:
487            hour = epoch2dt(call['timestamp']).hour
488            if not (8 <= hour <22):
489                count += 1
490            #
491        data['call_night_activity'] = count/len(calls)
```

```python
492
493          #Compute % of calls/texts outgoing from user. This works because true=1
494          data['call_outgoing'] = sum([c['type'] == 2 for c in calls])/len(calls)
495          data['text_outgoing'] = sum([t['type'] == 2 for t in texts])/len(texts)
496
497  #================================================================================
498  # Crunch location data
499  #================================================================================
500
501      if include_gps:
502          with open(userfile_path+user+"/gps_log.txt", 'r') as f:
503              raw = [LE(line) for line in f.readlines()]
504
505          if not is_sorted([l['timestamp'] for l in raw]):
506              raw = sort_dicts_by_key(raw, 'timestamp')
507
508          #We only want measurements taken at least 500s apart.
509          prev = 0
510          gps_data = []
511          allowed_providers = ['gps', 'network'] if include_network else ['gps']
512          for line in raw:
513              now = line['timestamp']
514              if line['provider'] in allowed_providers and now − prev >= 500:
515                  #Convert coordinates to km and note it down
516                  x = (longitude2meters*line['lon'] − x_offset)*0.001
517                  y = (latitude2meters*line['lat'] − y_offset)*0.001
518                  gps_data.append({'point' : (x,y), 'timestamp' : now,
519                                   'smarttime' : epoch2dt(now)})
520                  prev = now
521              #
522          #ignore user if there aren't enough data
523          if not len(gps_data) >= minimum_number_of_gps_points:
524              return None
525
526          # We want to investigate each day saparately so start at midnight.
527          first_midnight = next_midnight(gps_data[0]['smarttime'])
528          for i in xrange(len(gps_data)):
529              if gps_data[i]['smarttime'] >= first_midnight:
530                  gps_data = gps_data[i:]
531                  break
532
533
534          #Generate list of radii of smallest enclosing circle , SEC, for each day
535          current_points = []
536          prev = gps_data[0]['timestamp']
537          radii = []
538
539          distances = []
540          early_day = gps_data[0]['smarttime']
541          late_day = next_midnight(early_day)
542          for datum in gps_data:
543              now = datum['smarttime']
544              while not early_day <= now < late_day:
545                  if len(current_points) > 2:
546                      crds = [p['point'] for p in current_points]
547                      circle = make_circle(crds)
548                      r = circle[2] if circle else 0
549                      if circle and r > 0:
550                          if r <= 500:
551                              radii.append(r)
552                      distances.append(sum([get_distance(crds[i], crds[i+1])
553                                       for i in xrange(len(crds)−1)]))
```

```
554                     # Reset counters and update bins
555                     current_points = []
556                     early_day = late_day
557                     late_day = next_midnight(late_day)
558                 current_points.append(datum)
559
560
561          data['radius_of_gyration_med'] = np.median(radii)
562          data['radius_of_gyration_std'] = np.std(radii)
563          data['travel_med'] = np.median(distances)
564          data['travel_std'] = np.std(distances)
565
566          #Run Lloyd's algorithm to identify  clusters
567          #Determine which points are stationary − less  movement than 100m
568          stationary_data = []
569          try:
570              for i in xrange(1,len(gps_data)−1):
571                  a,b,c = tuple([gps_data[ind]['point'] for ind in
572                                 [i−1, i, i+1]])
573                  if (get_distance(a, b) < 0.1 and get_distance(b, c) < 0.1):
574                      stationary_data.append(gps_data[i])
575                  #
576              initial_clusters = 50
577              threshold_percent = 0.05
578              points = [elem['point'] for elem in stationary_data]
579
580              minimum_points = int(threshold_percent*len(points))
581
582
583              clusters_scatter = lloyds.lloyds(points, initial_clusters, runs=3,
584                                 init='scatter')
585
586              clusters_sample = lloyds.lloyds(points, initial_clusters, runs=3,
587                                 init='sample')
588
589              #Determine most succesful method
590              locs = lambda c: [p for p in c.values() if len(p)>=minimum_points]
591              if len(locs(clusters_scatter)) >= len(locs(clusters_sample)):
592                  method = 'scatter'
593                  best = clusters_scatter
594              else:
595                  method = 'sample'
596                  best = clusters_sample
597
598              # Number of places which survive cutoff (true = 1,  so  just  sum)
599              n_places = sum([len(pl) >= minimum_points for pl in best.values()])
600
601              #Output plots of clustering
602              if plot_clusters:
603                  if not os.path.isdir('pics'):
604                      os.mkdir('pics')
605                  for ext in ['.pdf', '.png']:
606                      filename = 'pics/'+user+"_"+method+ext
607                      lloyds.draw_clusters(clusters = best,
608                                  threshold = minimum_points,
609                                  show = False,
610                                  filename = filename)
611          except:
612              if not allow_nan:
613                  return None
614              n_places = float('nan')
615          data['n_places'] = n_places
```

```python
616
617        #Compute location entropy and add to data
618        try:
619            n_points = sum(len(location) for location in best.values())
620            data['location_entropy'] = sum([-len(p)/n_points*math.log(len(p)/n_points)
621                                    for p in best.values()])
622        except ValueError:
623            data['location_entropy'] = float('nan')
624
625        '''Guess where people live. Probably where they spend weeknights...
626         It's important to avoid selection bias here (people probably turn off
627         their phone when sleeping at home but not while partying at DTU, which
628         means fewer data points at their actual home).
629         This is rectified by excluding points that aren't logged monday to
630         thursday and only recording one 'late' or 'early' data point pr date.
631         These points are labelled 'weird' and are used to determine the user's
632         home.'''
633        try:
634            weirdpoints = []
635            latedays = []
636            earlydays = []
637            for datum in stationary_data:
638                now = datum['smarttime']
639                #ignore weekends
640                if now.weekday() > 3:
641                    continue
642                thisdate = (now.year, now.month, now.day)
643                if now.hour >= 20 and not thisdate in latedays:
644                    weirdpoints.append(datum['point'])
645                    latedays.append(thisdate)
646                elif now.hour <= 7 and not thisdate in earlydays:
647                    weirdpoints.append(datum['point'])
648                    earlydays.append(thisdate)
649
650            best_score = 0
651            home = None
652            for key, val in best.iteritems():
653                score = len(set(val).intersection(weirdpoints))
654                if score > best_score:
655                    home = key
656                    best_score = score
657                #
658            # Estimate how much user spends at home
659            ordered_gps = sort_dicts_by_key(gps_data, 'timestamp')
660            is_home = lambda p: get_distance(home, p) <= 0.200
661            time_home = 0
662            time_away = 0
663            for i in xrange(len(ordered_gps)-1):
664                a = ordered_gps[i]
665                b = ordered_gps[i+1]
666                dt = b['timestamp'] - a['timestamp']
667                if dt > 7200:
668                    continue
669                elif is_home(a['point']) and is_home(b['point']):
670                    time_home += dt
671                elif (not is_home(a['point'])) and (not is_home(b['point'])):
672                    time_away += dt
673                #
674            data['home_away_time_ratio'] = time_home/time_away
675        except:
676            if not allow_nan:
677                return None
```

```
678              data['home_away_time_ratio'] = float('nan')
679
680
681    #===============================================================================
682    # Facebook data
683    #===============================================================================
684
685        if include_facebook:
686            with open(userfile_path+user+'/facebook_log.txt', 'r') as f:
687                n = len(f.readlines())
688                if n < minimum_number_of_facebook_friends:
689                    return None
690                data['number_of_facebook_friends'] = n
691
692    #===============================================================================
693    # Bluetooth data
694    #===============================================================================
695
696        if include_bluetooth:
697            with open(userfile_path+user+"/bluetooth_log.txt", 'r') as f:
698                raw = [LE(line) for line in f.readlines()]
699            #Make sure data is sorted chronologically
700            if not is_sorted([entry['timestamp'] for entry in raw]):
701                raw = sort_dicts_by_key(raw, 'timestamp')
702            #List of times when user was social
703            social_times = []
704            total_social_time = 0
705            total_time = 0
706            #maps from each other user encountered to time spend with said user
707            friend2time_spent = {}
708            #Temporary variables
709            current_time = 0
710            previous_time = 0
711            current_users = []
712            previous_users = []
713            for signal in raw:
714                new_time = signal['timestamp']
715                #Check if line represents a new signal and if so, update values
716                if new_time != current_time:
717                    dt = new_time - current_time
718                    #Determine number of pings
719                    overlap = set(previous_users).intersection(set(current_users))
720                    if len(overlap) >= bluetooth_social_threshold:
721                        social_times.append(previous_time)
722                        if dt <= 7200:
723                            total_social_time += dt
724                            total_time += dt
725                            for friend in overlap:
726                                try:
727                                    friend2time_spent[friend] += dt
728                                except KeyError:
729                                    friend2time_spent[friend] = dt
730                        #
731                    elif dt <= 7200:
732                        total_time += dt
733                    #Update variables
734                    previous_time = current_time
735                    previous_users = current_users
736                    current_users = []
737                    current_time = new_time
738                new_user = signal['name']
739                if new_user=='-1' or not new_user:
```

```
740                continue
741            else:
742                current_users.append(new_user)
743            #
744        # Add fraction of time spent social to output
745        data['fraction_social_time'] = total_social_time/total_time
746        # Compute social entropy
747        normfac = 1.0/sum(friend2time_spent.values())
748        ent = sum(-t*normfac*math.log(t*normfac)
749                    for t in friend2time_spent.values())
750        data['social_entropy'] = ent
751
752        data['bluetooth_daily_entropy']=timestamps2daily_entropy(social_times,
753                                           hours_pr_daily_rythm_bin)
754
755        #Ensure time span is suficcient to make a time series
756        if not (social_times[-1] - social_times[0] > 24*3600*7
757                +1+3600*hours_pr_ar_bin):
758            return None
759
760        #Fit AR-series and append parameters to output
761        try:
762            dts = [epoch2dt(timestamp) for timestamp in social_times]
763            time_series = make_time_series(dts, hours_pr_ar_bin)
764            model = ar_model.AR(time_series)
765            n_params = int(24*7/hours_pr_ar_bin + 1) #1 week plus 1 extra bin
766            result = model.fit(maxlag = None, ic = None)
767            params = result.params#[1:]
768            while len(params) < n_params:
769                params.append(float('nan'))
770        except:
771            if not allow_nan:
772                return None
773            params = [float('nan') for _ in xrange(n_params)]
774
775        count = 0
776        for par in params:
777            count += 1
778            name = "bluetooth_activity_AR_coeff_"+str(count)
779            data[name] = par
780
781        # Compute autocorrelation coeffs and append to output
782        try:
783            accs = get_autocorrelation_coefficients(time_series, n_params)[:1]
784        except:
785            if not allow_nan:
786                return None
787            accs = [float('nan') for _ in xrange(n_params + 1)]
788        for i in xrange(len(accs)):
789            name = "bluetooth_activity_acc_"+str(i)
790            data[name] = accs[i]
791
792 #==============================================================================
793 # Wrap up user
794 #==============================================================================
795
796    # Double check thata doesn't containing nan values
797    if any(np.isnan(value) for value in data.values()) and not allow_nan:
798        return None #Discard user due to insufficient data
799
800    #Collect results
801    final = {'user' : user, 'data' : data, 'profile' : profile}
```

```
802        return final
803
804
805    if __name__ == '__main__':
806        user_counter = 0
807        #Make job queue
808        pool = multiprocessing.Pool(processes = n_jobs)
809        jobs = []
810        for user in users:
811            user_counter += 1
812            args = {'user' : user, 'user_counter' : user_counter}
813            jobs.append(pool.apply_async(process_user, kwds = args))
814
815        pool.close()   #run
816        pool.join()   #Wait for  remaining jobs
817
818        #Write results
819        with open(output_filename, 'a') as f:
820            for job in jobs:
821                result = job.get()
822                if not result:
823                    continue
824                json.dump(result, f)
825                f.write("\n")
826
827
828        #Done.
829        pushme("Data extraction done.")
```

## A.2.4   Social Fabric Code

```
1    # −∗− coding: utf−8 −∗−
2    """This module aims to allow sharing of some common methods and settings
3    when testing and tweaking various machine learning schemes.
4    Always import settings and the like from here!"""
5
6    from __future__ import division
7    import abc
8    from collections import Counter
9    import itertools
10   import json
11   import math
12   import matplotlib.colors as mcolors
13   import matplotlib.pyplot as plt
14   import matplotlib.patches as mpatches
15   import multiprocessing
16   import numpy as np
17   import random
18   from scipy.sparse import dok_matrix
19   from sklearn import svm
20   from sklearn.ensemble import RandomForestRegressor, RandomForestClassifier
21   from sklearn.cross_validation import (cross_val_score, LeaveOneOut, KFold,
22                                         StratifiedKFold)
23   import sys
24   from time import time
25   import traceback
26
```

```python
27   oldhat = (35/256,39/256,135/256)
28   nude = (203/256,150/256,93/256)
29   wine = (110/256,14/256,14/256)
30   moerkeroed = (156/256,30/256,36/256)
31
32   def _make_colormap(seq):
33       """Return a LinearSegmentedColormap
34        seq: a sequence of floats and RGB-tuples. The floats should be increasing
35        and in the interval (0,1).
36        """
37       seq = [(None,) * 3, 0.0] + list(seq) + [1.0, (None,) * 3]
38       cdict = {'red': [], 'green': [], 'blue': []}
39       for i, item in enumerate(seq):
40           if isinstance(item, float):
41               r1, g1, b1 = seq[i - 1]
42               r2, g2, b2 = seq[i + 1]
43               cdict['red'].append([item, r1, r2])
44               cdict['green'].append([item, g1, g2])
45               cdict['blue'].append([item, b1, b2])
46       return mcolors.LinearSegmentedColormap('CustomMap', cdict)
47
48   color_map = _make_colormap([oldhat, moerkeroed, 0.33, moerkeroed, nude, 0.67, nude])
49
50   big_five = ['openness', 'conscientiousness', 'extraversion', 'agreeableness',
51               'neuroticism']
52
53   #_default_features = ["n_texts",
54   #                       "ct_iet_std ",
55   #                       " call_cir ",
56   #                       "call_entropy",
57   #                       " text_cir ",
58   #                       "n_calls ",
59   #                       "text_latency ",
60   #                       "call_outgoing",
61   #                       "fraction_social_time ",
62   #                       "text_outgoing",
63   #                       " call_iet_std ",
64   #                       "n_text_contacts",
65   #                       " call_night_activity ",
66   #                       "call_iet_med",
67   #                       "outgoing_activity_AR_coeff_2",
68   #                       "text_entropy",
69   #                       " ct_cir ",
70   #                       "text_response_rate",
71   #                       "n_ct_contacts",
72   #                       "social_entropy",
73   #                       "n_call_contacts ",
74   #                       "n_ct",
75   #                       " text_iet_std ",
76   #                       "ct_iet_med",
77   #                       "ct_entropy",
78   #                       "text_iet_med",
79   #                       "call_response_rate",
80   #                       "number_of_facebook_friends"]
81
82   _default_features = ['call_iet_med', 'text_iet_med', 'social_entropy',
83   'call_entropy', 'travel_med', 'n_places', 'text_latency',
84   'call_night_activity']
85
86   def split_ntiles(values, n):
87       '''Determines the values that separate the imput list into n equal parts.
88        this is a generalization of the notion of median (in the case n = 2) or
```

```python
89          quartiles (n=4).
90          Usage: ntiles([5,6,7], 2) gives [6] for instance.'''
91      result = []
92      for i in xrange(1,n):
93          percentile = 100/n * i
94          result.append(np.percentile(values, percentile,
95                                  interpolation='linear'))
96      return result
97
98  def determine_ntile(value, ntiles):
99      '''Determines which n-tile the input value belongs to.
100         Usage: determine_ntile([7,9,13], 10) gives 2 (third quartile).
101         This uses zero indexing so data split into e.g. quartiles will give results
102         like 0,1,2,3 - NOT 1,2,3,4.'''
103     #Check if value is outside either extreme, meaning n-tile 1 or n.
104     if value >= ntiles[-1]:
105         return len(ntiles) #Remember the length is n-1
106     elif value < ntiles[0]:
107         return 0 #Values was in the first n-tile
108     # Define possible region and search for where value is between two elements
109     left = 0
110     right = len(ntiles)-2
111     #Keep checking th middle of the region and updating region
112     ind = (right + left)//2
113     while not ntiles[ind] <= value < ntiles[ind + 1]:
114         #Check if lower bound tile is on the left
115         if value < ntiles[ind]:
116             right = ind - 1
117         else:
118             left = ind + 1
119         ind = (right + left)//2
120     # Being between ntiles 0 and 1, means second n-tile and so on.
121     return ind + 1
122
123 def assign_labels(Y, n):
124     '''Accepts a list and an int n and returns a list of discrete labels
125      corresponding to the ntile each original y-value was in.'''
126     ntiles = split_ntiles(Y, n)
127     labels = [determine_ntile(y, ntiles) for y in Y]
128     return labels
129
130 def normalize_data(list_):
131     '''Normalizes input data to the range [-1, 1]'''
132     lo, hi = min(list_), max(list_)
133     if lo == hi:
134         z = len(list_)*[0]
135         return z
136     else:
137         return [2*(val-lo)/(hi - lo) - 1 for val in list_]
138
139
140 def read_data(filename, trait, n_classes = None, normalize = True,
141             features='default', interpolate = True):
142     '''This reads in a preprocessed datafile, splits psych profile data into
143     n classes if specified, filters desired psychological traits and
144     features and returns as a tuple (X,Y, indexdict), which can be fed to a'
145     number of off-the-shelf ML schemes.
146     If trait=='Sex', female and male are converted to 0 and 1, respectively.
147     indexdict maps each element of the feature vectors to their label, as in
148     {42 : 'distance_travelled_pr_day'} etc.
149
150     Args:
```

```python
151         filename : str
152           Name of the file containing the data.
153         trait : str
154           The psychological trait to extract data on.
155         n_classes : int
156           Number of classes to split data into. Default is None,
157           i.e. just keep the decimal values. Ignored if trait == 'Sex', as data
158           only has two discreet values.
159         normalize : bool
160           Whether to hard normalize data to [-1, 1].
161         features : str/list
162           Which features to read in. Can also be 'all'
163           or 'default', meaning the ones I've pragmatically found to be
164           reasonable.
165
166         interpolate : bool:
167           Whether to replace NaN's with the median value of
168           the feature in question.'''
169     if trait == 'sex':
170         n_classes = None
171     #Read in the raw data
172     with open(filename, 'r') as f:
173         raw = [json.loads(line) for line in f.readlines()]
174     #Get list of features to be included - everything if nothing's specified
175     included_features = []
176     if features == 'default':
177         included_features = _default_features
178     elif features == 'all':
179         included_features = raw[0]['data'].keys() if features=='all' else features
180     else:
181         included_features = features
182
183     #Remove any features that only have NaN values.
184     for i in xrange(len(included_features)-1,-1,-1):
185         feat = included_features[i]
186         if all(math.isnan(line['data'][feat]) for line in raw):
187             del included_features[i]
188
189
190     # ----- Handle feature vectors -----
191     # Dict mapping indices to features
192     indexdict = {ind : feat for ind, feat in enumerate(included_features)}
193     # N_users x N_features array to hold data
194     rows = len(raw)
195     cols = len(included_features)
196     X = np.ndarray(shape = (rows, cols))
197     for i in xrange(rows):
198         line = raw[i]
199         #Construct data matrix
200         for j, feat in indexdict.iteritems():
201             val = line['data'][feat]
202             X[i, j] = val
203         #
204     #Replace NaNs with median values
205     if interpolate:
206         for j in xrange(cols):
207             #Get median of feature j
208             med = np.median([v for v in X[:, j] if not math.isnan(v)])
209             if math.isnan(med):
210                 raise ValueError('''Feature %s contains only NaN's and should
211                                  have been removed.''' % indexdict[j])
212             for i in xrange(rows):
```

```python
213                    if math.isnan(X[i,j]):
214                        X[i,j] = med
215
216        if normalize:
217            for j in xrange(cols):
218                col = X[:, j]
219                X[:, j] = normalize_data(col)
220
221        # ----- Handle class info -----
222        trait_values = []
223        for line in raw:
224            #Add value of psychological trait
225            psych_trait = line['profile'][trait]
226            if trait == 'Sex':
227                if psych_trait == 'Female':
228                    psych_trait = 0
229                elif psych_trait == 'Male':
230                    psych_trait = 1
231                else:
232                    raise ValueError('My code is binary gender normative, sorry.')
233            trait_values.append(psych_trait)
234        Y = []
235        if n_classes == None:
236            Y = trait_values
237        else:
238            ntiles = split_ntiles(trait_values, n_classes)
239            Y = [determine_ntile(tr, ntiles) for tr in trait_values]
240
241        return (X, Y, indexdict)
242
243
244    def plot_stuff(input_filename, output_filename=None, color=moerkeroed):
245        with open(input_filename, 'r') as f:
246            d = json.load(f)
247        x = d['x']
248        y = d['y']
249        yerr = d['mean_stds']
250        plt.plot(x,y, color=color, linestyle='dashed',
251                marker='o')
252        plt.errorbar(x, y, yerr=yerr, linestyle="None", marker="None",
253                    color=color)
254        if output_filename:
255            plt.savefig(output_filename)
256
257    def get_TPRs_and_FPRs(X, Y, forest = None, verbose = False):
258        '''Accepts a list of feature vectors and a list of labels and returns a
259         tuple of true positive and false positive rates (TPRs and FPRs,
260         respectively) for various confidence thresholds.'''
261        kf = LeaveOneOut(n=len(Y))
262
263        results = []
264        thresholds = []
265
266        counter = 0
267        for train, test in kf:
268            counter += 1
269            if counter % 10 == 0 and verbose:
270                print "Testing on user %s of %s..." % (counter, len(Y))
271
272            result = {}
273            train_data = [X[i] for i in train]
274            train_labels = [Y[i] for i in train]
```

```python
275            test_data = [X[i] for i in test]
276            test_labels = [Y[i] for i in test]
277
278            if not forest:
279                forest = RandomForestClassifier(
280                                        n_estimators = 1000,
281                                        n_jobs=-1,
282                                        criterion='entropy')
283
284            forest.fit(train_data, train_labels)
285            result['prediction'] = forest.predict(test_data)[0]
286            result['true'] = test_labels[0]
287            confidences = forest.predict_proba(test_data)[0]
288            result['confidences'] = confidences
289            thresholds.append(max(confidences))
290
291            results.append(result)
292
293        #ROC curve stuff - false and true positive rates
294        TPRs = []
295        FPRs = []
296
297        unique_thresholds = sorted(list(set(thresholds)), reverse=True)
298
299        for threshold in unique_thresholds:
300            tn = 0
301            fn = 0
302            tp = 0
303            fp = 0
304            for result in results:
305                temp = result['prediction']
306                if temp == 1 and result['confidences'][1] >= threshold:
307                    pred = 1
308                else:
309                    pred = 0
310                if pred == 1:
311                    if result['true'] == 1:
312                        tp += 1
313                    else:
314                        fp += 1
315                    #
316                elif pred == 0:
317                    if result['true'] == 0:
318                        tn += 1
319                    else:
320                        fn += 1
321                    #
322                #
323            TPRs.append(tp/(tp + fn))
324            FPRs.append(fp/(fp + tn))
325        return (TPRs, FPRs)
326
327 def make_roc_curve(TPRs, FPRs, output_filename = None):
328     '''Accepts a list of true and false positive rates (TPRs and FPRs,
329      respectively) and generates a ROC-curve.'''
330     predcol = moerkeroed
331     basecol = oldhat
332     fillcol = nude
333
334     fig = plt.figure()
335     ax = fig.add_subplot(1,1,1)
336
```

```python
337          TPRs = [0] + TPRs + [1]
338          FPRs = [0] + FPRs + [1]
339
340          area = 0.0
341          for i in xrange(len(TPRs)−1):
342              dx = FPRs[i+1] − FPRs[i]
343              y = 0.5*(TPRs[i] + TPRs[i+1])
344              under_curve = dx*y
345              baseline = dx*0.5*(FPRs[i] + FPRs[i+1])
346              area += under_curve − baseline
347
348          baseline = FPRs
349          ax.fill_between(x = FPRs, y1 = TPRs, y2 = baseline, color = fillcol,
350                          interpolate = True, alpha=0.8)
351          ax.plot(baseline, baseline, color = basecol, linestyle = 'dashed',
352                  linewidth = 1.0, label = 'Baseline')
353          ax.plot(FPRs, TPRs, color=predcol, linewidth = 1,
354                          label = 'Prediction')
355
356          plt.xlabel('False positive rate.')
357          plt.ylabel('True positive rate.')
358
359          handles, labels = ax.get_legend_handles_labels()
360          hest = mpatches.Patch(color=fillcol)
361
362          labels += ['Area = %.3f' % area]
363          handles += [hest]
364          ax.legend(handles, labels, loc = 'lower right')
365 #        plt.legend(handles = [tp_line, base])
366          if output_filename:
367              plt.savefig(output_filename)
368          plt.show()
369
370  def rank_features(X, Y, forest, indexdict, limit = None):
371      '''Ranks the features of a given dataset and classifier.
372       indexdict should be a map from indices to feature names like
373       {0 : 'average_weigth'} etc.
374       if limit is specified, this method returns only the top n ranking features.
375       Returns a dict like {'feature name' : (mean importances, std)}.'''
376      importances = forest.feature_importances_
377      stds=np.std([tree.feature_importances_ for tree in forest.estimators_],
378                  axis=0)
379      indices = np.argsort(importances)[::−1]
380      if limit:
381          indices = indices[:limit]
382      d = {indexdict[i] : (importances[i], stds[i]) for i in indices}
383      return d
384
385  def check_performance(X, Y, clf, strata = None):
386      '''Checks forest performance compared to baseline.'''
387      N_samples = len(Y)
388      #Set up validation indices
389      if not strata: #Do leave−one−out validation
390          skf = KFold(N_samples, n_folds=N_samples, shuffle = False)
391      else: #Do stratified K−fold
392          skf = StratifiedKFold(Y, n_folds=strata)
393
394      #Evaluate classifier performance
395      scores = []
396      for train, test in skf:
397          train_data = [X[ind] for ind in train]
398          train_labels = [Y[ind] for ind in train]
```

```python
399            test_data = [X[ind] for ind in test]
400            test_labels = [Y[ind] for ind in test]
401
402            #Check performance of input forest
403            clf.fit(train_data, train_labels)
404            score = clf.score(test_data, test_labels)
405            scores.append(score)
406
407        #Compute baseline
408        most_common_label = max(Counter(Y).values())
409        baseline = float(most_common_label)/N_samples
410
411        #Compare results with prediction baseline
412        score_mean = np.mean(scores)/baseline
413        score_std = np.std(scores)/baseline
414
415        return (score_mean, score_std)
416
417 def check_regressor(X, Y, reg, strata = None):
418        '''Checks the performance of a regressor against mean value baseline.'''
419        N_samples = len(Y)
420        #Set up validation indices
421        if not strata:  #Do leave−one−out validation
422            skf = KFold(N_samples, n_folds=N_samples,
423                                    shuffle = False)
424        else:  #Do stratified K−fold
425            skf = StratifiedKFold(Y, n_folds=strata)
426
427        #Evaluate performance
428        model_abs_errors = []
429        baseline_abs_errors = []
430        for train, test in skf:
431            train_data = [X[ind] for ind in train]
432            train_labels = [Y[ind] for ind in train]
433            test_data = [X[ind] for ind in test]
434            test_labels = [Y[ind] for ind in test]
435
436            #Check performance of input forest
437            reg.fit(train_data, train_labels)
438            base = np.mean(train_labels)
439            for i in xrange(len(test_data)):
440                pred = reg.predict(test_data[i])
441                true = test_labels[i]
442                model_abs_errors.append(np.abs(pred − true))
443                baseline_abs_errors.append(np.abs(base − true))
444            #
445        return (np.mean(model_abs_errors), np.mean(baseline_abs_errors))
446
447 class _RFNN(object):
448        __metaclass__ = abc.ABCMeta
449        '''Abstract class for random forest nearest neightbor predictors.
450         This should never be instantiated.'''
451
452        def __init__(self, forest, n_neighbors):
453            self.forest = forest
454            self.n_neighbors = n_neighbors
455            self.X = None
456            self.Y = None
457
458        def fit(self, X, Y):
459            '''Fits model to training data.
460
```

```python
461            Args
462            ------------
463            X : List
464              List of training feature vectors.
465
466            Y : List
467              List of training labels or values to be predicted.'''
468          if not len(X) == len(Y):
469              raise ValueError("Training input and output lists must have "
470                              "same length.")
471          if not self.n_neighbors <= len(X):
472              raise ValueError("Fewer data points than neighbors.")
473
474          self.forest.fit(X, Y)
475          self.X = X
476          self.Y = Y
477
478      def _rf_similarity(self, a, b):
479          '''Computes a similarity measure for two points using a trained random
480           forest classifier.'''
481          if self.X == None or self.Y == None:
482              raise NotImplementedError("Model has not been fittet to data yet.")
483
484          #Feature vectors must be single precision.
485          a = np.array([a], dtype = np.float32)
486          b = np.array([b], dtype = np.float32)
487          hits = 0
488          tries = 0
489          for estimator in self.forest.estimators_:
490              tries += 1
491              tree = estimator.tree_
492              # Check whether the points end up on the same leaf for this  tree
493              if tree.apply(a) == tree.apply(b):
494                  hits += 1
495              #
496          return hits/tries
497
498      def find_neighbors(self, point):
499          '''Determine the n nearest nieghbors for the given point.
500           Returns a list of n tuples like (yval, similarity).
501           The tuples are sorted descending by similarity.'''
502          if self.X == None or self.Y == None:
503              raise NotImplementedError("Model has not been fittet to data yet.")
504
505          #Get list  of tuples like  (y, similarity ) for the n 'nearest' points
506          nearest = [(None, float('-infinity')) for _ in xrange(self.n_neighbors)]
507          for i in xrange(len(self.X)):
508              similarity = self._rf_similarity(self.X[i], point)
509              # update top n list  if  more similar than the furthest neighbor
510              if similarity > nearest[-1][1]:
511                  nearest.append((self.Y[i], similarity))
512                  nearest.sort(key = lambda x: x[1], reverse = True)
513                  del nearest[-1]
514              #
515          return nearest
516
517      #Mandatory methods − must be overridden
518      @abc.abstractmethod
519      def predict(self, point):
520          pass
521
522      @abc.abstractmethod
```

```
523    def score(self, X, Y):
524        pass
525
526 def _reservoir_sampler(start = 1):
527    '''Generator of the probabilities need to do reservoir sampling. The point
528     it that this can be used to iterate through a list, discarding each element
529     for the following element with probability P_n and ending up with a random
530     element from the list.'''
531    n = start
532    while True:
533        p = 1/n
534        r = random.uniform(0,1)
535        if r < p:
536            yield True
537        else:
538            yield False
539        n += 1
540
541 class RFNNClassifier(_RFNN):
542    '''Random Forest Nearest Neighbor Classifier.
543
544     Parameters
545     ----
546     n_neighbors : int
547        Number of neighbors to consider.
548
549     forest : RandomForestClassifier
550        The forest which will provide a
551        distance measure on which determine nearest neighbors.
552
553     weighting : str
554        How to weigh the votes of different neighbors.
555        'equal' means each neighbor has an equivalent vote.
556        'linear' mean votes are weighed by their similarity to the input point.
557     '''
558    def predict(self, point):
559        '''Predicts the label of a given point.'''
560        neighbortuples = self.find_neighbors(point)
561        if self.weighting == 'equal':
562            #Simple majority vote. Select randomly if it's a tie.
563            predictions = [t[0] for t in neighbortuples]
564            best = 0
565            winner = None
566            switch = _reservoir_sampler(start = 2)
567            for label, votes in Counter(predictions).iteritems():
568                if votes > best:
569                    best = votes
570                    winner = label
571                    switch = _reservoir_sampler(start = 2)
572                elif votes == best:
573                    if switch.next():
574                        winner = label
575                    else:
576                        pass
577                else:
578                    pass
579                #
580            return winner
581
582        #Weigh votes by their similarity to the input point
583        elif self.weighting == 'linear':
584            #The votes are weighted by their similarity
```

```
585              d = {}
586              for yval, similarity in neighbortuples:
587                  try:
588                      d[yval] += similarity
589                  except KeyError:
590                      d[yval] = similarity
591              best = float('-infinity')
592              winner = None
593              for k, v in d.iteritems():
594                  if v > best:
595                      best = v
596                      winner = k
597                  else:
598                      pass
599              return winner
600
601      def score(self, X, Y):
602          if not len(X) == len(Y):
603              raise ValueError("Training data and labels must have same length.")
604
605          hits = 0
606          n = len(X)
607          for i in xrange(n):
608              pred = self.predict(X[i])
609              if pred == Y[i]:
610                  hits += 1
611              #
612          return hits/n
613
614
615      def __init__(self, forest = None, n_neighbors = 3, weighting = 'equal',
616                  n_jobs = 1):
617          #Make sure we have a forest * classifier *
618          if forest == None:
619              forest = RandomForestClassifier(n_estimators = 1000,
620                                              criterion = 'entropy',
621                                              n_jobs = n_jobs)
622          if not isinstance(forest, RandomForestClassifier):
623              raise TypeError("Forest must be a classifier")
624
625          self.weighting = weighting
626
627          #Call parent constructor
628          super(RFNNClassifier, self).__init__(forest, n_neighbors)
629
630  class RFNNRegressor(_RFNN):
631      '''Random Forest Nearest Neighbor Regressor.
632
633      Parameters
634      ----
635      n_neighbors : int
636        Number of neighbors to consider.
637
638      forest : RandomForestRegressor
639        The forest which will provide a
640        distance measure on which determine nearest neighbors.
641
642      weighting : str
643        How to weigh the votes of different neighbors.
644        'equal' means each neighbor has an equivalent weight.
645        'linear' mean votes are weighed by their similarity to the input point.
646      '''
```

```python
647    def predict(self, point):
648        # lists  of the y vaues and similarities  of nearest neighbors
649        neighbortuples = self.find_neighbors(point)
650        yvals, similarities = zip(*neighbortuples)
651
652        # Weigh each neighbor y value equally is that's how we roll
653        if self.weighting == 'equal':
654            weight = 1.0/len(yvals)
655            result = 0.0
656            for y in yvals:
657                result += y*weight
658            return result
659            #
660        # Otherwise, weigh neighbors by similarity
661        elif self.weighting == 'linear':
662            weight = 1.0/(sum(similarities))
663            result = 0.0
664            for i in xrange(len(yvals)):
665                y = yvals[i]
666                similarity = similarities[i]
667                result += y*similarity*weight
668            return result
669
670
671    def score(self, X, Y):
672        if not len(X) == len(Y):
673            raise ValueError("X and Y must be same length.")
674        errors = [Y[i] - self.predict(X[i]) for i in xrange(len(X))]
675        return np.std(errors)
676
677
678    def __init__(self, forest = None, n_neighbors = 3, weighting = 'equal'):
679        #Check forest type.
680        if forest == None:
681            forest = RandomForestRegressor(n_estimators = 1000, n_jobs = -1)
682        if not isinstance(forest, RandomForestRegressor):
683            raise TypeError("Must use Random Forest Regressor to initialize.")
684        # Set params
685        self.weighting = weighting
686        # Done. Call parent constructor
687        super(RFNNRegressor, self).__init__(forest, n_neighbors)
688
689
690 class _BaselineRegressor(object):
691     '''Always predicts the mean of the training set.'''
692     def __init__(self, guess=None):
693         self.guess = guess
694     def fit(self, xtrain, ytrain):
695         '''Find the average of input lidt of target values and guess on that
696          from now on.'''
697         self.guess = np.mean(ytrain)
698     def predict(self, x):
699         return self.guess
700
701
702 class _BaselineClassifier(object):
703     '''Always predicts the most common label in the training set'''
704     def __init__(self, guess=None):
705         self.guess = guess
706     def fit(self, xtrain, ytrain):
707         '''Find the most common label and guess on that from now on.'''
708         countmap = Counter(ytrain)
```

```python
709                best = 0
710                for label, count in countmap.iteritems():
711                    if count > best:
712                        best = count
713                        self.guess = int(label)
714                    #
715            #
716        def predict(self, x):
717            return self.guess




722    def _worker(X, Y, score_type, train_percentage, classifier, clf_args, n_groups,
723                replace):
724        '''Worker method for parallelizing bootstrap evaluations.'''
725        #Create bootstrap sample
726        try:
727            rand = np.random.RandomState() #Ensures PRNG works in children
728            indices = rand.choice(xrange(len(X)), size = len(X), replace = replace)
729            rand.randint
730            xsample = [X[i] for i in indices]
731            ysample = [Y[i] for i in indices]
732            #Create regressor if we're doing regression
733            if classifier == 'RandomForestRegressor':
734                clf = RandomForestRegressor(**clf_args)
735            elif classifier == 'SVR':
736                clf = svm.SVR(**clf_args)
737            elif classifier == 'baseline_mean':
738                clf = _BaselineRegressor()
739
740            #Create classifier and split dataset into labels
741            elif classifier == 'RandomForestClassifier':
742                clf = RandomForestClassifier(**clf_args)
743                ysample = assign_labels(ysample, n_groups)
744            elif classifier == 'SVC':
745                clf = svm.SVC(**clf_args)
746                ysample = assign_labels(ysample, n_groups)
747            elif classifier == 'baseline_most_common_label':
748                clf = _BaselineClassifier()
749                ysample = assign_labels(ysample, n_groups)
750            #Fail if none of the above classifiers were specified
751            else:
752                raise ValueError('Regressor or classifier not defined.')
753            #Generate training and testing set
754            cut = int(train_percentage*len(X))
755            xtrain = xsample[:cut]
756            ytrain = ysample[:cut]
757            xtest = xsample[cut:]
758            ytest = ysample[cut:]
759
760            #Fit the classifier or regressor
761            clf.fit(xtrain, ytrain)
762
763            #Compute score and append to output list
764            if score_type == 'mse':
765                scores = [(ytest[i] - clf.predict(xtest[i]))**2
766                          for i in xrange(len(xtest))]
767                return np.mean(scores)
768
769            elif score_type == 'fraction_correct':
770                n_correct = sum([ytest[i] == int(clf.predict(xtest[i]))
```

```python
771                              for i in xrange(len(ytest))])
772                 score = n_correct/len(ytest)
773                 return score
774            elif score_type == 'over_baseline':
775                 #Get score
776                 score = sum([ytest[i] == int(clf.predict(xtest[i]))
777                              for i in xrange(len(ytest))])
778                 #Get baseline
779                 baselineclf = _BaselineClassifier()
780                 baselineclf.fit(xtrain, ytrain)
781                 baseline = sum([ytest[i] == baselineclf.predict(xtest[i])
782                              for i in xrange(len(ytest))])
783                 return score/baseline
784
785             #Fail  if  none of the  above performance metrics were specified
786             else:
787                 raise ValueError('Score type not defined.')
788
789             #Job's  done!
790             return None
791        except:
792            raise Exception("".join(traceback.format_exception(*sys.exc_info())))
793
794
795 def bootstrap(X, Y, classifier, score_type = 'mse', train_percentage = 0.8,
796               clf_args = {}, iterations = 1000, n_groups = 3, n_jobs = 1,
797               replace = True):
798     '''Performs bootstrap resampling to evaluate the performance of some
799     classifier or regressor. Note that this takes the *complete dataset* as
800     arguments as well as arguments specifying which predictor to use and which
801     function to estimate the distribution of.
802     This seems to be the most straightforward generalizable implementation
803     which can be parallelized, as passing e.g. the scoring function directly
804     clashed with the mechanisms implemented to work around the GIL for
805     multiprocessing for obscure reasons.
806
807     Parameters:
808     ----------------
809     X : list
810        All feature vectors in the complete dataset.
811
812     Y : list
813        All 'true' labels or output values in the complete dataset.
814
815     classifier : str
816        Which classifier to use to predict the test set. Allowed values:
817        'RandomForestRegressor', 'baseline_mean', 'SVR',
818        'RandomForestClassifier', 'SVC', 'baseline_most_common_label'
819
820     score_type : str
821        String signifying which function to estimate the distribution of.
822        Allowed values: 'mse', 'fraction_correct', 'over_baseline'
823
824     train_percentage : float
825        The percentage [0:1] of each bootstrap sample to be used for training.
826
827     clf_args : dict
828        optional arguments to the constructor method of the regressor/classifier.
829
830     iterations : int
831        Number of bootstrap samples to run.
832
```

```
833    n_jobs : int
834      How many cores (maximum) to use.
835
836    replace : bool
837      Whether to sample with replacement when obtaining the bootstrap samples.
838    '''
839
840    if not len(X) == len(Y):
841        raise ValueError("X and Y must have equal length.")
842
843    #Arguments to pass to worker processes
844    d = {'X' : X, 'Y' : Y, 'train_percentage' : train_percentage,
845        'classifier' : classifier, 'clf_args' : clf_args,
846        'score_type' : score_type, 'n_groups' : n_groups,
847        'replace' : replace}
848
849    #Make job queue
850    pool = multiprocessing.Pool(processes = n_jobs)
851    jobs = [pool.apply_async(_worker, kwds = d) for _ in xrange(iterations)]
852    pool.close()  #run
853    pool.join()  #Wait for remaining jobs
854
855    #Make sure no children died too early
856    if not all(job.successful() for job in jobs):
857        raise RuntimeError('Some jobs failed.')
858
859    return [j.get() for j in jobs]
860
861 def get_correlations(X, Y):
862    '''Given a list of feature vectors X and labels or values Y, returns a list
863     of correlation coefficients for each dimension of the feature vectors.'''
864    n_feats = len(X[0])
865    correlations = []
866    for i in xrange(n_feats):
867        temp = np.corrcoef([x[i] for x in X], Y)
868        correlation = temp[0,1]
869        if math.isnan(correlation):
870            correlation = 0
871        correlations.append(correlation)
872    return correlations
873
874
875 def make_kernel(correlations = None, gamma = 1.0, threshold = 0.0):
876    '''Returns a weighted radial basis function (WRBF) kernel.'''
877    def kernel(x,y, *args, **kwargs):
878        if correlations == None:
879            _corrs = np.ones(shape = (len(x),), dtype = np.float64)
880        else:
881            _corrs = correlations
882        d = len(_corrs)  #number of features
883        #Strong (above threshold) correlations
884        strong = [np.abs(c) if np.abs(c) >= threshold else 0.0
885                for c in _corrs]
886        normfactor = 1.0/np.sqrt(sum([e**2 for e in strong]))
887        #Metric to compute distance between points
888        metric = dok_matrix((d,d), dtype = np.float64)
889        for i in xrange(d):
890            metric[i,i] = strong[i]*normfactor
891        #
892        result = np.zeros(shape = (len(x), len(y)))
893        for i in xrange(len(x)):
894            for j in xrange(len(y)):
```

```
895                   dist = x[i] − y[j]
896                   result[i,j] = np.exp(−gamma*np.dot(dist,dist))
897           return result
898       return kernel
899
900   if __name__ == '__main__':
901       pass
902       X, Y, ind_dict = read_data('../data.json', trait = 'openness',
903                       features = ['call_iet_med', 'text_iet_med', 'social_entropy', 'call_entropy', '↩
              travel_med', 'n_places', 'text_latency', 'call_night_activity'],
904                       n_classes = 3
905                       )
906   #     X =      [[1,2,7,0],[3,1,6,0.01],[6,8,1,0],[10,8,2,0.01]]
907   #     Y =   [1,1,0,0]
908
909       cut = int(0.8*len(X))
910
911       xtrain = X[:cut]
912       ytrain = Y[:cut]
913       xtest = X[cut:]
914       ytest = Y[cut:]
915
916       corrs = get_correlations(X, Y)
917       print corrs
918
919       C = 70
920       gamma = 3.75
921
922       kernel = make_kernel(corrs, 0.05)
923
924       clf = svm.SVC(kernel = kernel)
925       clf.fit(xtrain, ytrain)
926
927       hits = 0
928
929       for i in xrange(len(xtest)):
930           if clf.predict(xtest[i]) == ytest[i]:
931               hits += 1
932           #
933       print 100.0*hits/len(ytest)
934
935       for i in xrange(len(corrs)):
936           print ind_dict[i], corrs[i]
937
938
939
940   #     print len(X[0])
941   ##     print i
942   #     print [el for el in i.values() if 'init' in el]
```

## A.2.5   Lloyd's Algorithm

```
1   # −*− coding: utf−8 −*−
2   from __future__ import division
3
4   import numpy as np
5   import matplotlib
```

```python
matplotlib.use('Agg') #ugly hack to allow plotting from terminal
import matplotlib.pyplot as plt
import random
from copy import deepcopy

def _dist(p,q):
    return sum([(p[i]−q[i])**2 for i in xrange(len(p))])

def _lloyds_single_run(X, K, max_iterations, init):
    # Initialize  with a subset of the data
    if init == 'sample':
        initials = random.sample(X, K)
    # Or initialize  with random points across the same range as data
    elif init == 'scatter':
        vals = zip(*X)
        xmin = min(vals[0])
        xmax = max(vals[0])
        ymin = min(vals[1])
        ymax = max(vals[1])
        initials = [(random.uniform(xmin, xmax),
                     random.uniform(ymin, ymax)) for _ in xrange(K)]
    # Or yell  RTFM at user
    else:
        raise ValueError('Invalid initialization mode!')

    #Contruct hashmap mapping integers up to K to centroids
    centroids = dict(enumerate(initials))
    converged = False
    iterations = 0

    while not converged and iterations < max_iterations:
        clusters = {i : []  for i in xrange(K)}
        #Make sure clusters and centroids have identical  keys, or  we're doomed.
        assert set(clusters.keys()) == set(centroids.keys())
        prev_centroids = deepcopy(centroids)

        ### STEP ONE −update clusters
        for x in X:
            #Check distances to all  centroids
            bestind = −1
            bestdist = float('inf')
            for ind, centroid in centroids.iteritems():
                dist = _dist(x, centroid)
                if dist < bestdist:
                    bestdist = dist
                    bestind = ind
                #
            clusters[bestind].append(x)

        ### STEP TWO −update centroids
        for ind, points in clusters.iteritems():
            if not points:
                pass  #Cluster's  empty − nothing to update
            else:
                centroids[ind] = np.mean(points, axis = 0)

        ### We're converged when all old centroids = new centroids.
        converged = all([_dist(prev_centroids[k],centroids[k]) == 0
                    for k in xrange(K)])
        iterations += 1
        #
    return {tuple(centroids[i]) : clusters[i] for i in xrange(K)}
```

```python
68
69  def lloyds(X, K, runs = 1, max_iterations = float('inf'), init = 'sample'):
70      '''Runs Lloyd's algorithm to identify K clusters in the dataset X.
71      X is a list of points like [[x1,y1],[x2,y2]---].
72      Returns a hash of centroids mapping to points in the corresponding cluster.
73      The objective is to minimize the sum of distances from each centroid to
74      the points in the corresponding cluster. It might only converge on a local
75      minimum, so the configuration with the lowest score (sum of distances) is
76      returned.
77      init denotes initialization mode, which can be 'sample', using a randomly
78      select subset of the input data, or 'scatter', using random points selected
79      from the same range as the data as initial centroids.
80
81      Parameters
82      ---------------
83      X : array_like
84        list of points. 2D example: [[3,4],[3.4, 7.2], ...]
85
86      K : int
87        Number of centroids
88
89      runs : int
90        Number of times to run the entire algorithm. The result with the lowest
91        score will be returned.
92
93      max_iterations : int or float
94        Number of steps to allow each run. Default if infinit, i.e. the algorithm
95        runs until it's fully converged.
96
97      init : str
98        Initialization mode. 'sample' means use a random subset of the data as
99        starting centroids. 'scatter' means place starting centroids randomly in
100       the entire x-y range of the dataset.
101
102     Returns
103     --------------
104     result : dict
105       A dictionary in which each key is a tuple of coordinated corresponding to
106       a centroid, and each value is a list of points belonging to that cluster.
107       '''
108
109     record = float('inf')
110     result = None
111     for _ in xrange(runs):
112         clusters = _lloyds_single_run(X, K, max_iterations = max_iterations,
113                                       init = init)
114         #Determine how good the clusters came out
115         score = 0
116         for centroid, points in clusters.iteritems():
117             score += sum([_dist(centroid, p) for p in points or [] ])
118         if score < record:
119             result = clusters
120             record = score
121         #
122     return result
123
124
125 def _makecolor():
126     i = 0
127     cols = ['b', 'g', 'r', 'c', 'm', 'y']
128     while True:
129         yield cols[i]
```

```python
130            i = (i+1)%len(cols)
131
132
133  def draw_clusters(clusters, threshold = 0, show = True, filename = None):
134      '''Accepts a dict mapping cluster centroids to cluster points and makes
135       a color-coded plot of them. Clusters containing fewer points than the
136       threshold are plottet in black.'''
137      colors = _makecolor()
138      plt.figure()
139      for centroid, points in clusters.iteritems():
140          if not points:
141              continue
142          if len(points) < threshold:
143              style = ['k,']
144          else:
145              color = colors.next()
146              style = [color+'+']
147              #Plot centroids
148              x,y = centroid
149              plt.plot(x,y, color = color, marker = 'd', markersize = 12)
150          #plot points
151          plt.plot(*(zip(*points)+style))
152      if filename:
153          plt.savefig(filename, bbox_inches = 'tight')
154      if show:
155          plt.show()
156
157
158  if __name__ == '__main__':
159      points = [[random.uniform(−10,10), random.uniform(−10,10)] for _ in xrange(10**3)]
160      clusters = lloyds(X = points, K = 6, runs = 1)
161      draw_clusters(clusters = clusters, filename = 'lloyds_example.pdf')
```

## A.2.6   Smallest Enclosing Circle

```python
1   # −*− coding: utf−8 −*−
2   #
3   # Smallest enclosing circle
4   #
5   # Copyright (c) 2014 Project Nayuki
6   # http :// www.nayuki.io/page/smallest−enclosing−circle
7   #
8   # This program is free software: you can redistribute it and/or modify
9   # it under the terms of the GNU General Public License as published by
10  # the Free Software Foundation, either version 3 of the License, or
11  # (at your option) any later version.
12  #
13  # This program is distributed in the hope that it will be useful,
14  # but WITHOUT ANY WARRANTY; without even the implied warranty of
15  # MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the
16  # GNU General Public License for more details.
17  #
18  # You should have received a copy of the GNU General Public License
19  # along with this program (see COPYING.txt).
20  # If not, see <http ://www.gnu.org/licenses/>.
21  #
22
```

```
23  import math, random
24
25
26  # Data conventions: A point is a pair of floats (x, y). A circle is a triple of floats (center x, center y, radius).
27
28  #
29  # Returns the smallest circle that encloses all the given points. Runs in expected O(n) time, randomized.
30  # Input: A sequence of pairs of floats or ints, e.g. [(0,5), (3.1,-2.7)].
31  # Output: A triple of floats representing a circle.
32  # Note: If 0 points are given, None is returned. If 1 point is given, a circle of radius 0 is returned.
33  #
34  def make_circle(points):
35      '''Accepts list of points as tuples and returns (x, y, r).'''
36      # Convert to float and randomize order
37      shuffled = [(float(p[0]), float(p[1])) for p in points]
38      random.shuffle(shuffled)
39
40      # Progressively add points to circle or recompute circle
41      c = None
42      for (i, p) in enumerate(shuffled):
43          if c is None or not _is_in_circle(c, p):
44              c = _make_circle_one_point(shuffled[0 : i + 1], p)
45      return c
46
47
48  # One boundary point known
49  def _make_circle_one_point(points, p):
50      c = (p[0], p[1], 0.0)
51      for (i, q) in enumerate(points):
52          if not _is_in_circle(c, q):
53              if c[2] == 0.0:
54                  c = _make_diameter(p, q)
55              else:
56                  c = _make_circle_two_points(points[0 : i + 1], p, q)
57      return c
58
59
60  # Two boundary points known
61  def _make_circle_two_points(points, p, q):
62      diameter = _make_diameter(p, q)
63      if all(_is_in_circle(diameter, r) for r in points):
64          return diameter
65
66      left = None
67      right = None
68      for r in points:
69          cross = _cross_product(p[0], p[1], q[0], q[1], r[0], r[1])
70          c = _make_circumcircle(p, q, r)
71          if c is None:
72              continue
73          elif cross > 0.0 and (left is None or _cross_product(p[0], p[1], q[0], q[1], c[0], c[1]) > _cross_product(←
          p[0], p[1], q[0], q[1], left[0], left[1])):
74              left = c
75          elif cross < 0.0 and (right is None or _cross_product(p[0], p[1], q[0], q[1], c[0], c[1]) < _cross_product(←
          p[0], p[1], q[0], q[1], right[0], right[1])):
76              right = c
77      return left if (right is None or (left is not None and left[2] <= right[2])) else right
78
79
80  def _make_circumcircle(p0, p1, p2):
81      # Mathematical algorithm from Wikipedia: Circumscribed circle
82      ax = p0[0]; ay = p0[1]
```

```
83        bx = p1[0];  by = p1[1]
84        cx = p2[0];  cy = p2[1]
85        d = (ax * (by - cy) + bx * (cy - ay) + cx * (ay - by)) * 2.0
86        if d == 0.0:
87            return None
88        x = ((ax * ax + ay * ay) * (by - cy) + (bx * bx + by * by) * (cy - ay) + (cx * cx + cy * cy) * (ay - by)) / d
89        y = ((ax * ax + ay * ay) * (cx - bx) + (bx * bx + by * by) * (ax - cx) + (cx * cx + cy * cy) * (bx - ax)) / d
90        return (x, y, math.hypot(x - ax, y - ay))
91
92
93    def _make_diameter(p0, p1):
94        return ((p0[0] + p1[0]) / 2.0, (p0[1] + p1[1]) / 2.0, math.hypot(p0[0] - p1[0], p0[1] - p1[1]) / 2.0)
95
96
97    _EPSILON = 1e-12
98
99    def _is_in_circle(c, p):
100       return c is not None and math.hypot(p[0] - c[0], p[1] - c[1]) < c[2] + _EPSILON
101
102
103   # Returns twice the signed area of the triangle defined by (x0, y0), (x1, y1), (x2, y2)
104   def _cross_product(x0, y0, x1, y1, x2, y2):
105       return (x1 - x0) * (y2 - y0) - (y1 - y0) * (x2 - x0)
106
107   #pts = [(0.0, 0.0), (6.0, 8.0)]
108   #
109   #test = make_circle(pts)
110   #
111   #print test
```

# BIBLIOGRAPHY

[1] Yoshua Bengio and Yves Grandvalet. No Unbiased Estimator of the Variance of K-Fold Cross-Validation. 5:1089–1105, 2004.

[2] Avrim Blum, Adam Kalai, and John Langford. Beating the Hold-out: Bounds for {K}-fold and Progressive Cross-Validation. *Proceedings of the 12th Annual Conference on Computational Learning Theory*, (c):203–208, 1999.

[3] L Breiman. Random forests. *Machine learning*, pages 5–32, 2001.

[4] Cjc Christopher J C Burges. A Tutorial on Support Vector Machines for Pattern Recognition. *Data Mining and Knowledge Discovery*, 2(2):121–167, 1998.

[5] Yves Alexandre De Montjoye, Jordi Quoidbach, Florent Robic, and Alex Pentland. Predicting personality using novel mobile phone-based metrics. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 7812 LNCS:48–55, 2013.

[6] JM John M Digman. Personality structure: Emergence of the five-factor model. *Annual review of psychology*, 41:417–440, 1990.

[7] RA Fisher. The use of multiple measurements in taxonomic problems. *Annals of eugenics*, 1936.

[8] E Gabrilovich and S Markovitch. Computing Semantic Relatedness Using Wikipedia-based Explicit Semantic Analysis. *IJCAI*, 2007.

[9] Isabelle Guyon. Gene Selection for Cancer Classification. pages 389–422, 2002.

[10] Marti a. Hearst, Susan T. Dumais, Edgar Osuna, John Platt, and Bernhard Schölkopf. Support vector machines. *IEEE Intelligent Systems and their Applications*, 13(4):18–28, 1998.

[11] Thomas K Landauer, Peter W. Foltz, and Darrell Laham. An introduction to latent semantic analysis. *Discourse Processes*, 25(2-3):259–284, 1998.

[12] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schutze. An Introduction to Information Retrieval. *Online*, (c):569, 2009.

[13] Preslav Nakov. Latent semantic analysis of textual data. *Proceedings of the conference on Computer systems and technologies - CompSysTech '00*, pages 5031–5035, 2000.

[14] AD Pietersma. Kernel Learning in Support Vector Machines using Dual-Objective Optimization. *Proceedings of the 23rd . . .* , 2011.

[15] G. Salton, a. Wong, and C. S. Yang. A vector space model for automatic indexing. *Communications of the ACM*, 18(11):613–620, 1975.

[16] Gerard Salton and Christopher Buckley. Term-weighting approaches in automatic text retrieval. *Information Processing & Management*, 24(5):513–523, 1988.

[17] Philip Wolfe. A duality theorem for nonlinear programming. *Quarterly of applied mathematics*, 19(3):239–244, 1961.

[18] PA Zandbergen and SJ Barbeau. Positional accuracy of assisted gps data from high-sensitivity gps-enabled mobile phones. *Journal of Navigation*, 2011.