# 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

# ENGLISH ABSTRACT

WORDS! SOOOOO MANY WORDS!

# Dansk sammenfatning

ORD! SAAAAAAAAAA MANGE ORD

Part I

SOCIAL FABRIC PROJECT

# PHONE ACTIVITY AND QUANTITATIVE DATA EXTRACTION

THE main objective of part I 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 1.1.

Next, I set out to replicate some recent research results claiming that ⟶ kilder!!! 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 1.2 is the extraction of the many-dimensional data points or *feature vectors* $X$ from the phone logs of the participants, while section 1.3 gives a brief description of the psychological traits $Y$. Finally, section 1.4 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 2.

## 1.1   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.1.1. 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.1.1. 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.

### 1.1.1   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 1.1: 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 1.1 and 1.2, 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 1.2: 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.

### 1.1.2  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 1.1.1 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 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
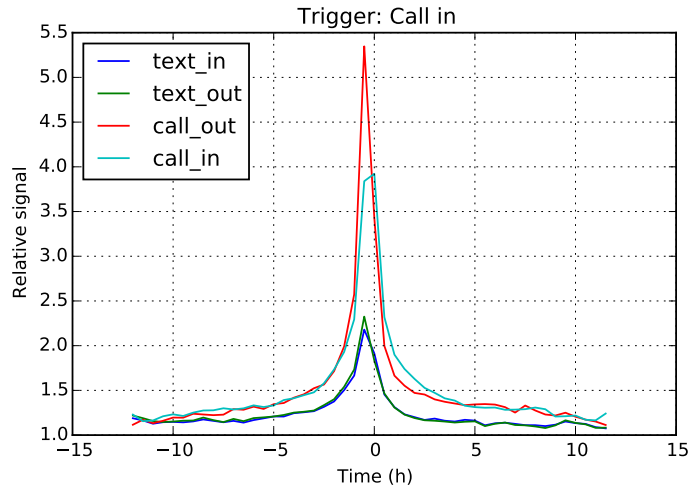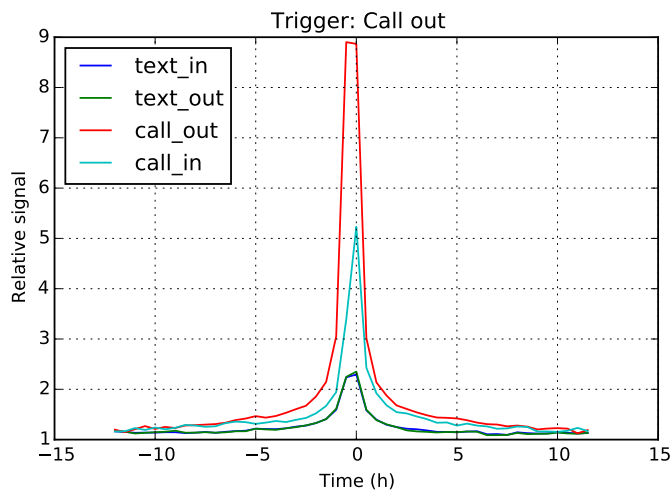yield false measurements of high average speeds when the measurements
are not temporally separated. A lot of the measurements turned out to
be grouped somewhat tightly - for instance, approximately 80% of the
time intervals were below 100 s. This occurs because the Social Fabric
data not only actively 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 1.3. A typical uncertainty on civilian GPS devices is at
most 100 m[21], so because I choose to consider a user active if they travel
at a mean speed of 0.5 m/s, and based on the time spacings shown in figure
1.3, I chose to discard measurements separated by less than 500 s.

An analysis like that of section 1.1.1 reveals that a user's phone activity
is significantly increased around times when they are on the move, as
shown in figure 1.4. Note the asymmetry of the signal, especially visible
in figure 1.4(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 1.4(b) appears to be increasing at around
±24 h, which would seem reasonable assuming people have slightly

Figure 1.3: 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 1.5 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 1.4: Relative increase of activities triggered by GPS activity.

Figure 1.5: GPS-triggered activity increase over an interval of 36 hours using a bin size of 10 minutes.

### 1.1.3  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 1.6, 1.7 and 1.8 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 1.6: 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 1.7: 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 1.8: 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 1.6 through 1.8. This code is included in section A.1.2.

## 1.2 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[6]. 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.1.3.

### 1.2.1 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

$$S_u = -\sum_c \frac{n_c}{n_t} \ln_2 \frac{n_c}{n_t}, \tag{1.1}$$

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.

### 1.2.2   *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[6] 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

$$S = \sum_K \sum_{x_n \in C_k} |x_n - c_k|^2, \tag{1.2}$$

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 (1.2), 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.1.5.

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 1.9. 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 (1.2). A few such iterations are shown in figure 1.10. The difference in end results for the two methods is exemplified in figure 1.11. 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 1.9: A few iterations of Lloyd's algorithm using 'sample' initial-ization. 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 1.10: 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 (1.2).

(a) Sample initialization.



(b) Scattered initialization.

Figure 1.11: 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.

### 1.2.3 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 1.12.

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

### 1.2.4 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.

### 1.2.5 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 1.1.3. 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 1.12: 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 1.2.3. I also chose to compute a 'social entropy' much like (1.1), but weighted by the time the user spends with each acquaintance:

$$E = -\sum_i f_i \ln_2 (f_i), \qquad (1.3)$$

$$f_i = \frac{\text{time spent with } i}{\sum_j \text{time spent with } j}. \qquad (1.4)$$

Note that the denominator of (1.4) 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.

## 1.3 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[7]. **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.

## 1.4 Linear Discriminant Analysis & Gender Prediction

Linear discriminant analysis is basically a dimensionality reduction technique developed by Fisher in 1936 [8] 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

$$y = \omega^T x \tag{1.5}$$

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.

### 1.4.1 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, $\mathcal{N}_i(x) = N(x; \mu_i, \sigma_i^2)$, the joint probability distribution for the distance between the projections will be the convolution

$$P(x) = \int_{-\infty}^{\infty} \mathcal{N}_a(y) \cdot \mathcal{N}_b(x - y) \, dy. \tag{1.6}$$

Computing this for a Gaussian distribution,

$$\mathcal{N}(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \tag{1.7}$$

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

Denoting convolution by $*$ and Fourier transforms by

$$\mathcal{F}(f) = \frac{1}{(2\pi)^{n/2}} \int_{\mathbb{R}^n} f(x) \cdot e^{-i\omega x} \, dx, \tag{1.8}$$

the convolution theorem is derived as follows:

$$\mathcal{F}(f * g) = \frac{1}{(2\pi)^{n/2}} \int_{\mathbb{R}^n} \int_{\mathbb{R}^n} f(y) \cdot g(x - y) \, dy \, e^{-i\omega x} \, d\omega, \tag{1.9}$$

$$= \frac{1}{(2\pi)^{n/2}} \int_{\mathbb{R}^n} f(y) \int_{\mathbb{R}^n} g(x - y) e^{-i\omega x} \, dy \, d\omega, \tag{1.10}$$

$$= \frac{1}{(2\pi)^{n/2}} \int_{\mathbb{R}^n} f(y) \int_{\mathbb{R}^n} g(z) e^{-i\omega(z+y)} \, dz \, d\omega, \tag{1.11}$$

$$= \frac{1}{(2\pi)^{n/2}} \int_{\mathbb{R}^n} f(y) e^{-i\omega y} \int_{\mathbb{R}^n} g(z) e^{-i\omega z} \, dz \, d\omega, \tag{1.12}$$

$$\boxed{\mathcal{F}(f * g) = (2\pi)^{n/2} \mathcal{F}(f) \cdot \mathcal{F}(g),} \tag{1.13}$$

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 (1.8), rather than

$$\mathcal{F}(f) = \int_{\mathbb{R}^n} f(x) \cdot e^{-2\pi i v x} \, dx. \tag{1.14}$$

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

$$\mathcal{N}_a * \mathcal{N}_b = (2\pi)^{n/2} \mathcal{F}^{-1}\left(\mathcal{F}(\mathcal{N}_a) \cdot \mathcal{F}(\mathcal{N}_b)\right). \tag{1.15}$$

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

$$
\begin{aligned}
\mathcal{F}(\mathcal{N}(x)) &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \cdot e^{-i\omega x}\, \mathrm{d}x, \\
&= \frac{1}{2\pi\sigma} \int_{-\infty}^{\infty} e^{-\frac{x^2}{2\sigma^2}} e^{-i\omega(x+\mu)}\, \mathrm{d}x, \\
&= \frac{1}{2\pi\sigma} e^{-i\omega\mu} \int_{-\infty}^{\infty} e^{-\frac{x^2}{2\sigma^2}} (\cos(\omega x) + i\sin(\omega x))\, \mathrm{d}x, \\
&= \underbrace{\frac{1}{2\pi\sigma} e^{-i\omega\mu}}_{a} \underbrace{\int_{-\infty}^{\infty} e^{-\frac{x^2}{2\sigma^2}} \cos(\omega x)\, \mathrm{d}x}_{I(\omega)} .
\end{aligned}
\tag{1.16}
$$

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 (1.16) gives the result

$$
\mathcal{F}(\mathcal{N}) = \frac{1}{\sqrt{2\pi}} e^{-i\omega\mu} e^{-\sigma^2\omega^2/2}.
\tag{1.17}
$$

This can be used in conjunction with (1.13) to obtain

$$
\begin{aligned}
\mathcal{F}(\mathcal{N}_a * \mathcal{N}_b) &= \sqrt{2\pi} \frac{1}{\sqrt{2\pi}} e^{-i\omega\mu_a} e^{-\sigma_a^2\omega^2/2} \cdot \frac{1}{\sqrt{2\pi}} e^{-i\omega\mu_b} e^{-\sigma_b^2\omega^2/2}, \tag{1.18} \\
&= \frac{1}{\sqrt{2\pi}} e^{-i\omega(\mu_a-\mu_b)} e^{-(\sigma_a^2+\sigma_b^2)\omega^2/2}, \tag{1.19}
\end{aligned}
$$

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

$$\mathcal{N}_a * \mathcal{N}_b = \frac{1}{\sqrt{2\pi}\sigma_s} e^{-\frac{(x-\mu_s)^2}{2\sigma_s^2}}. \tag{1.20}$$

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

$$d = \frac{(\mu_a - \mu_b)^2}{\sigma_a^2 + \sigma_b^2}. \tag{1.21}$$

### 1.4.2 Optimizing separation

To maximize the separation, the numerator and denominator, respectively, of (1.21) 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:

$$(\widetilde{\mu}_a - \widetilde{\mu}_b)^2 = \left(w^T (\mu_a - \mu_b)\right)^2, \tag{1.22}$$

$$= w^T (\mu_a - \mu_b)(\mu_a - \mu_b)^T w, \tag{1.23}$$

$$= w^T S_B w, \tag{1.24}$$

and

$$\widetilde{\sigma}_i^2 = \sum_{y \in i} \frac{1}{N} (y - \widetilde{\mu}_i)^2, \tag{1.25}$$

$$= w^T \sum_{y \in i} (x - \mu_i)(x - \mu_i)^T w, \tag{1.26}$$

$$= w^T S_i w, \tag{1.27}$$

$$\widetilde{\sigma}_a^2 + \widetilde{\sigma}_b^2 = w^T S_W w, \tag{1.28}$$

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

$$S_B = (\mu_a - \mu_b)(\mu_a - \mu_b)^T, \tag{1.29}$$

$$S_i = \sum_{y \in i} (x - \mu_i)(x - \mu_i)^T, \tag{1.30}$$

$$S_W = S_a + S_b. \tag{1.31}$$

Hence, the objective is to solve

$$\frac{\mathrm{d}}{\mathrm{d}w}J(w) = \frac{\mathrm{d}}{\mathrm{d}w}\left(\frac{w^T S_B w}{w^T S_W w}\right) = 0, \tag{1.32}$$

$$\frac{\frac{\mathrm{d}\left[w^T S_B w\right]}{\mathrm{d}w}w^T S_W w - w^T S_B w \frac{\mathrm{d}\left[w^T S_W w\right]}{\mathrm{d}w}}{\left(w^T S_W w\right)^2} = 0, \tag{1.33}$$

$$2S_B w \cdot w^T S_W w - w^T S_B w \cdot 2S_W w = 0, \tag{1.34}$$

$$S_B w - \frac{w^T S_B w \cdot S_W w}{w^T S_W w} = 0, \tag{1.35}$$

$$S_B w - S_W w J(w) = 0, \tag{1.36}$$

$$S_B w = S_W w J(w), \tag{1.37}$$

$$S_W^{-1} S_B w = J(w)w. \tag{1.38}$$

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

$$w^* = S_W^{-1}(\mu_a - \mu_b). \tag{1.39}$$

> Vær lige sikker på at du forstår det her.

Figure 1.13 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

$$P(a|d) = \frac{P(d|a)P(a)}{P(d|a)p(a) + P(d|b)P(b)}, \tag{1.40}$$

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 (1.40) is not a precise estimate of the probability of the point representing a given class, the class with the highest posteriori according to (1.40) 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 1.13: 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 (1.39), 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 1.14 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 1.14. This is described in section 2.1

Figure 1.14: 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.

# 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 1.2 and various machine learning schemes, starting with support vector machines (SVMs).

1000 kilder!!!

Uddyb når der er flere modeller.

## 2.1 Support Vector Machines

The purpose of this section is to introduce SVMs and attempt to apply them to the data obtained in 1.2. The introduction is mainly based on introductory texts by Marti Hearst [11] and Christopher Burges [5]. 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 1.4. 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 2.1. 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 $f : \mathbb{R}^N \rightarrow \{\pm1\}$ based on a set of training data in $\mathbb{R}^N \otimes \{\pm1\}$. I'll describe separately the properties of

Figure 2.1: The same points as those shown in figure 1.13, 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 1.14.

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:

$$\mathbf{w} \cdot (\mathbf{x} - \mathbf{p}) = 0. \tag{2.1}$$

The left hand side of (2.1) 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 (2.1) 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

$$\mathbf{w} \cdot (\mathbf{x} - \mathbf{p}) = \mathbf{w} \cdot \mathbf{x} + b = \pm 1, \tag{2.2}$$

for the points that are closest to the separating plane. Those points located on the margin are encircled in figure 2.1. In general then, the meaning of

the sign and magnitude of

$$\mathbf{w} \cdot \mathbf{x} + b \tag{2.3}$$

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 2.1 as

$$\mathbf{w} = \sum_i v_i \mathbf{x}_i, \tag{2.4}$$

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.

### 2.1.1   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 (2.2), the following inequalities hold:

$$\begin{aligned} \mathbf{x}_i \cdot \mathbf{w} + b \geq 1, \, y_i = +1, \\ \mathbf{x}_i \cdot \mathbf{w} + b \leq -1, \, y_i = -1. \end{aligned} \tag{2.5}$$

Multiplying by $y_i$, both simply become

$$y_i \left( \mathbf{x}_i \cdot \mathbf{w} + b \right) - 1 \geq 0. \tag{2.6}$$

The distance between the separating plane and each of the margins in figure 2.1 is $1/|\mathbf{w}|$, so in order to maximize the separation, $|\mathbf{w}|$ must be minimized. For mathematical convenience, $\frac{1}{2}|\mathbf{w}|^2$, rather than $|\mathbf{w}|$ is included in the Lagrangian, which then becomes

$$L = \frac{1}{2}|\mathbf{w}|^2 - \sum_i \alpha_i y_i \left( \mathbf{x}_i \cdot \mathbf{w} + b - 1 \right), \tag{2.7}$$

must be minimized with the constraints

$$\alpha_i \geq 0, \tag{2.8}$$

$$\frac{\partial L}{\partial \alpha_i} = 0. \tag{2.9}$$

A result from convex optimization theory known as Wolfe Duality[20] states that one may instead maximize the above Lagrangian subject to

$$\nabla_w L = \frac{\partial L}{\partial b} = 0, \tag{2.10}$$

which gives conditions

$$\mathbf{w} = \sum_j \alpha_j y_j \mathbf{x}_j, \tag{2.11}$$

$$\sum_j \alpha_j y_j = 0. \tag{2.12}$$

These can be plugged back into (2.7) to obtain

$$L_D = \frac{1}{2} \sum_i \sum_j \alpha_i y_i \alpha_j y_j \mathbf{x}_i \cdot \mathbf{x}_j - \sum_i \alpha_i y_i \Big( \mathbf{x}_i \cdot \sum_j \alpha_j y_j \mathbf{x}_j + b \Big) + \sum_i \alpha_i, \tag{2.13}$$

$$L_D = -\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \mathbf{x}_j + \sum_i \alpha_i. \tag{2.14}$$

A problem with this is that eqs 2.5 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 2.1 were permuted, the classifier shown in the plot would still do a very good job, but eqs. 2.5 would not be satisfiable, causing the method to fail. This is remedied by introducing slack variables[3]

$$\begin{aligned} \mathbf{x}_i \cdot \mathbf{w} + b &\geq 1 - \xi_i, \quad y_i = +1, \\ \mathbf{x}_i \cdot \mathbf{w} + b &\leq -(1 - \xi_i), \quad y_i = -1, \\ \xi_i &\geq 0, \end{aligned} \tag{2.15}$$

which allows the algorithm to misclassify. This should come without a cost to the overall Lagrangian, or one would just end up classifying randomly, so a 'cost term', $C \cdot \sum_i \xi_i$ is added as well. The value of $C$ (as well as a few similar parameters which have yet to be introduced) is usually determined experimentally by simply varying it across some space of possible values and choosing the value resulting in the best performance - see for instance figure 2.3. In the case of the misclassification cost parameter $C$, low values will result in low performance, whereas too large values will result in overfitting. The above can be rewritten exactly as previously, except another set of non-negative Lagrange multipliers $\mu_i$ are added to (2.7) to ensure positivity of the $\xi_i$, resulting in

$$L = \frac{1}{2} |\mathbf{w}|^2 + C \cdot \sum_i \xi_i - \sum_i \alpha_i y_i (\mathbf{x}_i \cdot \mathbf{w} + b - 1 + \xi_i) - \sum_i \mu_i \xi_i. \tag{2.16}$$

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

$$0 \le \alpha_i \le C. \tag{2.17}$$

The methods outlined above can also be used to solve regression, rather than classification, problems[19]. The training data will then be in $\mathbb{R}^{N+1}$ rather than in $\mathbb{R}^N \otimes \{\pm 1\}$, and the value of the decision function in (2.3) is predicted instead of using only its sign. The criterion of correct classification from (2.6) is replaced by the demand that predictions be within some tolerated margin $\epsilon$ of the true value of the training point $y_i$, so (2.5) becomes

$$-\epsilon \le \mathbf{x}_i \cdot \mathbf{w} + b - y_i \le \epsilon \tag{2.18}$$

so when slack variables $\xi_i$ and $\xi_i^*$ (for the lower and upper bound, respectively) like in (2.15) are introduced, the Lagrangian from (2.16) becomes

$$L = \frac{1}{2}|\mathbf{w}|^2 + C\sum_i \left(\xi_i + \xi_i^*\right), \tag{2.19}$$

with constraints

$$\mathbf{x}_i \cdot \mathbf{w} + b - y_i \ge -\left(\epsilon + \xi_i\right), \tag{2.20}$$

$$\mathbf{x}_i \cdot \mathbf{w} + b - y_i \le \epsilon + \xi_i^*, \tag{2.21}$$

$$\xi_i, \xi^* \ge 0. \tag{2.22}$$

The main point to be emphasized here is that the training data $\mathbf{x}_i$ only enter into the dual Lagrangian of (2.14) as inner products. This is essential when extending the SVM model to nonlinear cases, which is the subject of the following section.

### 2.1.2 Generalizing to the non-linear case

The fact that the data $\mathbf{x}_i$ only occur as inner products in (2.14) makes one way of generalizing to non-linearly separable datasets straightforward: Referring back to figure 1.14, 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

$$\Phi : \mathbb{R}^l \to \mathbb{R}^h, \quad h > l, \tag{2.23}$$

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

$$\Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j) = \left(\mathbf{x}_i \cdot \mathbf{x}_j\right)^2. \tag{2.24}$$

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

$$(\mathbf{x} \cdot \mathbf{y})^2 = \left( \sum_{i=1}^{2} x_i y_i \right)^2 = x_1^2 y_1^2 + 2 x_1 y_1 x_2 y_2 + x_2^2 y_2^2, \qquad (2.25)$$

meaning that one possibility for $\Phi$ is

$$\Phi : \mathbf{x} \mapsto \begin{pmatrix} x_1^2 \\ \sqrt{2} x_1 x_2 \\ x_2^2 \end{pmatrix} \qquad (2.26)$$

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

$$\left( \sum_{i=1}^{d} x_i \right)^n = \sum_{\sum_{i=1}^{d} k_i = n} \frac{n!}{\prod_{l=1}^{d} k_l!} \prod_{j=1}^{d} x_j^{k_j}, \qquad (2.27)$$

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 (2.27) runs over all combinations of $d$ integers which sum to $n$. This can be simplified by introducing a function $K$, which simply maps

$$K : n, d \mapsto \left\{ \{k\} \in \mathbb{N}^d \,\middle|\, \sum_{i=1}^{d} k_i = n \right\}, \qquad (2.28)$$

and denoting each of those collections $\{k\}_i$ so each of the coefficients in (2.27) can be written

$$\frac{n!}{\prod_{i=1}^{d} k_i!} = C_{\{k\}}. \qquad (2.29)$$

Then, (2.27) becomes

$$\left( \sum_{i=1}^{d} x_i \right)^n = \sum_{K(n,d)} C_{\{k\}} \prod_{j=1}^{d} x_j^{k_j} \qquad (2.30)$$

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.

$$\dim(\mathbb{R}^h) = |K(n,d)| = \left\| \left( \{k\} \in \mathbb{N}^d \middle| \sum_{i=1}^{d} k_i = n \right) \right\|, \qquad (2.31)$$

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

$$\binom{n + d - 1}{n} = \frac{(n + d - 1)!}{n!(d - 1)!}. \qquad (2.32)$$

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.

$$K(\mathbf{x}, \mathbf{y}) = \Phi(\mathbf{x}) \cdot \Phi(\mathbf{y}). \qquad (2.33)$$

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

$$K(\mathbf{x}, \mathbf{y}) = (\mathbf{x} \cdot \mathbf{y})^n, \qquad (2.34)$$

$$\Phi(\mathbf{x}) = \sum_{K(n,d)} \sqrt{C_{\{k\}}} \prod_{j=1}^{d} x_j^{k_j}, \qquad (2.35)$$

so that (2.33) 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:

$$K(\mathbf{x}, \mathbf{y}) = e^{|\mathbf{x}-\mathbf{y}|^2/2\sigma}. \tag{2.36}$$

Burges [5] shows that the polynomial kernel is valid, so I'll show how the argument extends to the RBF kernel in the following.

Mercer's condition[18] 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

$$\int K(\mathbf{x}, \mathbf{y}) g(\mathbf{x}) g(\mathbf{y}) \, d\mathbf{x} \, d\mathbf{y} \geq 0. \tag{2.37}$$

This can be shown be rewriting (2.36) as

$$K(\mathbf{x}, \mathbf{y}) = e^{(\mathbf{x}-\mathbf{y})\cdot(\mathbf{x}-\mathbf{y})/2\sigma} = e^{|\mathbf{x}|^2/2\sigma} e^{|\mathbf{y}|^2/2\sigma} e^{-\mathbf{x}\cdot\mathbf{y}/\sigma}, \tag{2.38}$$

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

$$e^{-\mathbf{x}\cdot\mathbf{y}/\sigma} = \sum_{i=0}^{\infty} \frac{(-1)^i}{i!\sigma^i} (\mathbf{x} \cdot \mathbf{y})^i, \tag{2.39}$$

but using (2.30) on the dot product gives

$$(\mathbf{x} \cdot \mathbf{y})^i = \left( \sum_{j=1}^{d} x_j y_j \right)^i = \sum_{K(i,d)} C_{\{k\}} \prod_{j=1}^{d} x_j^{k_j} y_j^{k_j} \tag{2.40}$$

so the Taylor expansion becomes

$$e^{-\mathbf{x}\cdot\mathbf{y}/\sigma} = \sum_{i=0}^{\infty} \sum_{K(i,d)} \frac{(-1)^i}{i!\sigma^i} C_{\{k\}} \prod_{j=1}^{d} x_j^{k_j} y_j^{k_j}, \tag{2.41}$$

which can be plugged back into (2.38) to yield

$$K(\mathbf{x}, \mathbf{y}) = \sum_{i=0}^{\infty} \sum_{K(i,d)} \frac{(-1)^i}{i!\sigma^i} C_{\{k\}} e^{|\mathbf{x}|^2/2\sigma} e^{|\mathbf{y}|^2/2\sigma} \prod_{j=1}^{d} x_j^{k_j} y_j^{k_j}. \tag{2.42}$$

The underlying reason for these algebraic shenanigans is that (2.42) is clearly separable so that the integral in (2.37) from Mercer's condition becomes

$$\int K(\mathbf{x}, \mathbf{y}) g(\mathbf{x}) g(\mathbf{y}) \, d\mathbf{x} \, d\mathbf{y} \tag{2.43}$$

$$= \sum_{i=0}^{\infty} \sum_{K(i,d)} \frac{(-1)^i}{i! \sigma^i} C_{\{k\}} \int_{\mathbb{R}^{2d}} e^{|\mathbf{x}|^2/2\sigma} e^{|\mathbf{y}|^2/2\sigma} \prod_{j=1}^{d} x_j^{k_j} y_j^{k_j} g(\mathbf{x}) g(\mathbf{y}) \, d\mathbf{x} \, d\mathbf{y} \tag{2.44}$$

$$= \sum_{i=0}^{\infty} \sum_{K(i,d)} \frac{(-1)^i}{i! \sigma^i} C_{\{k\}} \left( \int_{\mathbb{R}^d} e^{|\mathbf{x}|^2/2\sigma} \prod_{j=1}^{d} x_j^{k_j} g(\mathbf{x}) \, d\mathbf{x} \right) \cdot \left( \int_{\mathbb{R}^d} e^{|\mathbf{y}|^2/2\sigma} \prod_{j=1}^{d} y_j^{k_j} g(\mathbf{y}) \, d\mathbf{y} \right) \tag{2.45}$$

$$= \sum_{i=0}^{\infty} \sum_{K(i,d)} \frac{(-1)^i}{i! \sigma^i} C_{\{k\}} \left( \int_{\mathbb{R}^d} e^{|\mathbf{x}|^2/2\sigma} \prod_{j=1}^{d} x_j^{k_j} g(\mathbf{x}) \, d\mathbf{x} \right)^2 \tag{2.46}$$

$$\geq 0. \tag{2.47}$$

Hence, radial basis functions satisfy Mercer's condition and the kernel described above can be plugged into the dual Lagrangian from (2.14) to obtain

$$L_D = -\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j e^{|\mathbf{x}_i - \mathbf{x}_j|^2/2\sigma} + \sum_i \alpha_i, \tag{2.48}$$

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 2.2 shows the points again, along with the *decision frontier* i.e. the curve which separates regions in which points are classified into separate classes. The danger of overfitting should be clear from figure 2.2. 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.

Figure 2.2: The 'island' scenario of figure 1.14 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.

The $\gamma$ parameter is often fixed by performing a grid search similar to that discussed earlier.  Figure 2.3 shows the resulting heat map from a grid search.

Figure 2.3: 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 giving optimal values of $C = 0.52$ and $\gamma = 0.12$.

### 2.1.3 Implementing a custom weighted radial basis function kernel

In general, the features extracted as described in section 1.2 are too numerous for efficient implementations of most machine learning schemes. [6] work around this by determining the linear correlation between each feature and the target values, and then discarding features whose correlation is below some threshold and letting the rest contribute equally in the SVM. They never explicate this threshold but state that they include features 'significantly related' to the target values, which convention would suggest means a threshold of 0.05. [15] have demonstrated significant improvements in SVM performance by assigning variable importances to each feature, so I implemented a modified kernel function where I don't just discard features with linear correlations below some threshold but also assign to each remaining feature a normalized weight given by its correlation with the output variable.

I do this by defining a diagonal matrix $\hat{M}$ where the $i, i$th element is the linear correlation coefficient between feature $i$ and the target variable, and using the matrix as a metric in the exponent of the usual radial basis

functions so the kernel $K(\mathbf{x}, \mathbf{y})$ becomes

$$e^{\gamma |\mathbf{x}-\mathbf{y}|^2} \rightarrow e^{\gamma (\mathbf{x}-\mathbf{y})^T \hat{M} (\mathbf{x}-\mathbf{y})} \tag{2.49}$$

Note that this does not change the validity of the proof I gave in section 2.1.2, so this weighted RBF kernel also satisfies Mercer's condition and may hence be used as a kernel function. It turned out to be nontrivial to write an implementation of this which had a syntax consistent with that of the default methods available in the library I used, and which provided a similarly simple way to grid search over its parameters ($C$ and $\gamma$ as in the usual RBF kernel SVM along with a threshold parameter). I ended up solving this by writing a metamethod to provide a kernel matrix based on an input list of variable importances as well as a default sharpness parameter denoted $\gamma$.

```python
def make_kernel(importances, gamma = 1.0):
    '''Returns a weighted radial basis function (WRBF) kernel which can be
    passed to an SVM or SVR from the sklearn module.

    Parameters:
    -----------------------
    importances : list
      The importance of each input feature. The value of element i can mean
      e.g. the linear correlation between feature i and target variable y.
      None means feature will be weighted equally.

    gamma : float
      The usual gamma parameter denoting inverse width of the gaussian used.
    '''
    def kernel(x,y, *args, **kwargs):
        d = len(importances) #number of features
        impsum = sum([imp**2 for imp in importances])
        if not impsum == 0:
            normfactor = 1.0/np.sqrt(impsum)
        else:
            normfactor = 0.0
        #Metric to compute distance between points
        metric = dok_matrix((d,d), dtype = np.float64)
        for i in xrange(d):
            metric[i,i] = importances[i]*normfactor
        #
        result = np.zeros(shape = (len(x), len(y)))
        for i in xrange(len(x)):
            for j in xrange(len(y)):
                diff = x[i] - y[j]
                dist = diff.T.dot(metric*diff)
                result[i,j] = np.exp(-gamma*dist)
        return result
    return kernel
```

With that in place, it was a simple matter to implement classifiers and regressors by inheriting from the default classes and overriding the constructor methods like this:

Figure 2.4: Heat map resulting from a grid search over the parameter space of the linear correlation threshold and the default sharpness $\gamma$ of the radial basis functions showing optimal values of 0.25 and 0.15 for $\gamma$ and the threshold, respectively.

```
class WRBFR(svm.SVR):
    '''Weighted radial basis function support vector regressor.'''
    def __init__(self, importances, C = 1.0, epsilon = 0.1,
                 gamma = 0.0):
    kernel = make_kernel(importances = importances, gamma = gamma)
    super(WRBFR, self).__init__(C = C, epsilon = epsilon, kernel = kernel)
```

Figures 2.4 and 2.5 show heat maps resulting from grid searches over the threshold parameter versus $\gamma$ and $C$, respectively. Interestingly, the ideal correlation threshold seems to be well above the value used in [6].

Figure 2.5: Heat map resulting from a grid search over the parameter space of the linear correlation threshold and the cost parameter $C$ of the radial basis functions showing optimal values of 3.6 and 0.13 for $\gamma$ and the threshold, respectively.

## 2.1.4   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

$$\sigma^2 = \frac{1}{N} \sum_i^N (x_i - \mu)^2 , \qquad (2.50)$$

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

$$\hat{\sigma}^2 = \frac{1}{N-1} \sum_i^N (x_i - \mu)^2 . \qquad (2.51)$$

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.

## 2.2 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:

$$I_G = \sum_i f_i (1 - f_i) = 1 - \sum_i f_i^2, \tag{2.52}$$

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

$$I_E = -\sum_i f_i \log_2 f_i + \frac{n_a}{N} \sum_i a_i \log_2 a_i + \frac{n_b}{N} \sum_i b_i \log_2 b_i. \qquad (2.53)$$

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.[4]

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 2.6 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 2.6: 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 2.7. 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 2.1 can be sensitive to irrelevant data[15]. There exist off-the-shelf methods, such as recursive feature elimination (RFE)[10], 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 2.7: 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:

$$w_j = \frac{\langle s - s_p \rangle}{\text{std}(s - s_p)} \tag{2.54}$$

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.

## 2.3   Nearest Neighbour-classifiers

Do iiit!

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

# 2.4 Results

## 2.4.1 Big Five Personality Traits

HARJ

Figure 2.8: Comparison of performance of models using random forest, support vector regression and weighted radial basis function regressor 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.

### 2.4.2   Miscellaneous Traits

Her står der ting om figur 2.9.

Opdater med ML

Figure 2.9: Triangle of Pearson correlation coefficients for various miscellaneous traits.

Part II

WIKIPEDIA-BASED EXPLICIT SEMANTIC ANALYSIS

# 3

# 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

## 3.1 Methods

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

### 3.1.1 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.

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

$$P(B_j|T) = \frac{P(T|B_j)P(B_j)}{P(T)}, \tag{3.1}$$

$$= \frac{\prod_i P(w_i|B_j)P(B_j)}{\sum_j \prod_i P(w_i|B_j)P(B_j)}, \tag{3.2}$$

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.

### 3.1.2   Semantic Analysis

Salton et al proposed in their 1975 paper *A Vector Space Model for Automatic Indexing*[16] an approach where words and text fragments are mapped

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

$$T \rightarrow |V\rangle = \sum_i v_i |i\rangle, v_i \in \mathbb{R}, \tag{3.3}$$

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

$$S(V, W) = \langle \hat{V} | \hat{W} \rangle = \frac{\sum_i v_i w_i}{\left( \sum_i v_i^2 \right) \left( \sum_i w_i^2 \right)}, \tag{3.4}$$

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

$$D(V, W) = 1 - S(V, W). \tag{3.5}$$

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 (3.3) 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 [12]). 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[14]. 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[12].

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 ([9], 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.

## 3.2   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 3.3.

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 3.2.3)

These steps are elaborated upon in the following.

### 3.2.1   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 3.1  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 a 'title' tag is encountered, a callback method is triggered which assigns

---

[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 3.1: A simplified snippet, of a Wikipedia XML dump.

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.2.5. The remainder of the parser is my work, and is included in section A.2.1. 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.

### 3.2.2   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 $O\left(n^2\right)$ 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.2.2.

### 3.2.3 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

$$T_{ij} = \left(1 + \ln\left(\text{count}(w_i, c_j)\right)\right) \ln\left(\frac{n_c}{df_i}\right), \qquad (3.6)$$

where $n_c$ is the total number of concepts and

$$df_i = |\{c_k, w_i \in c_k\}| \qquad (3.7)$$

is the number of concepts whose corresponding article contains the $i$th word. Thus, the first part of (3.6), $1 + \ln\left(\text{count}(w_i, c_j)\right)$ 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 (3.6) 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 of documents containing word $i$ goes to 1, and takes its highest values if

Giver det mening?

word $i$ occurs with high frequency in only a few documents[13]. While (3.6) is not the only expression to have those properties, empirically it tends to achieve superior results in information retrieval[17].

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

$$T_{ij} \to \frac{T_{ij}}{\sqrt{\sum_i T_{ij}^2}}. \tag{3.8}$$

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 [9] 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 3.1.2. The code is included in section A.2.3.

> 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'

## 3.3 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.2.4. 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 3.2. 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. Using this, I wrote a script that harvested tweets mentioning some selected brand names for about 5 months which resulted in a massive dataset of 370 million tweets. This section demonstrates some applications of the methods outlined earlier to this dataset.

```
{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 3.2: Example of a downloaded tweet.

### 3.3.1 Trend Discovery and Monitoring

A simple application of the methods outlined above is a purely exploratory analysis. As an example, using code included in section A.2.6, I extracted the 10 concepts most closely related to every tweet mentioning Carlsberg for each week over a period of a few months. This gave a bar chart like figure 3.1 for each of those weeks. Figure 3.1 implies some significant relation between Carlsberg and Raidió Teilifís Éireann blablabla!!!

Erstat med det med den irske radio

Another application is to manually select a few concepts of interest and then monitor how the number of tweets strongly related to them develops over time. For example I did an exploratory analysis as described above and occasionally saw the concept 'Kim Little' surface. This turned out to be a young football player signed to a team sponsored by Carlsberg, and when I produced a series of bar charts like that shown in figure 3.2 and combined them into an animation, the bar representing Little tended to peak around dates on which news stories featuring her or her team could be found. Of course these method are only effective insofar the concepts of interest actually have an associated Wikipedia article. However, these concepts are merely the basis vectors of the semantic space in which an arbitrary text can be represented, so this procedure can be extended fairly elegantly to estimate the impact of e.g. a press release on social media.



Figure 3.1: Sæt det rigtige ind når du kommer hjem!!!

Figure 3.2: Bar chart showing the concepts most related to tweets about Carlsberg for January 15th 2014.

This is the subject of section 3.3.2.

### 3.3.2 Measuring Social Media Impact

One possible application of the software I wrote to perform ESA is to provide a quantitative measure of the impact of some event on social media. For instance, some corporation or organization might be interested in learning precisely effectively a campaign or press release has reached the general public as it is expressed by social media. While my tweet harvesting script was running, Google posted a blog entry[2] on their experiments with producing psychedelic images with deep neural network called *Deep Dreams* which received widespread attention on social media.

Using the text from the Deep Dreams blog post as a reference text, I converted each English tweet about Google from a period around the blog post to a semantic vector using (3.3) and computed their semantic cosine similarity with the reference text as described by (3.4). Figure 3.3 shows this behaviour in the time around the release of the blog post. A measure of impact should of course take into account the rate at which new tweets occur. Just as one would expect, not only semantic relatedness, but also tweet frequency increase drastically around an interesting event. To make the signal from figure 3.3 independent of the normal tweet frequency, yet sensitive to changes in it, I first find the difference in semantic relatedness between the signal and a reference signal obtained long before the event to be investigated, then modulate the signal by multiplying each bin with the activity $A$ given by the ratio between the current and baseline tweet frequency. As normal Twitter activity tends to vary greatly over a 24 hour period, I obtained a baseline tweet frequency by computing the average tweet rate, weighted by the activity, for each discreet time interval of the day, and then repeating that signal into the period after the event of interest. The baseline along with the observed tweet rate is shown in figure 3.4(a). Modulating in this fashion and computing the difference between the observed signal and the baseline gives an expression of the 'impact rate' of a given time bin, which can be integrated to obtain a measure of the total impact which is independent of the average tweet rate before publication, which is practical if one wishes to compare e.g. the success of media campaigns for companies or organizations of varying sizes. This is shown in figure 3.5. If one wishes to include the number of tweets posted in the impact measure, the result can be computed and visualized nicely by again obtaining a baseline cosine similarity from a period prior to publication and then considering e.g. the cumulative deviation from that following the release, as shown in figure 3.6.

---

[2]The original blog post is available at `http://googleresearch.blogspot.dk/2015/06/inceptionism-going-deeper-into-neural.html`

Figure 3.3: Graph of the mean semantic cosine similarity of tweets around the Deep Dreams press release. There is a clear peak around $t = 0$ and the similarity appears increased in the period following the release. The error bars represent the true standard deviation of the sample mean $\sigma / \sqrt{N}$ for each time bin, each representing a 10-minute interval.



Figure 3.4: Tweet activity around the time the Deep Dreams blog entry was posted. The signal tweet rate increases with a factor of about 5 relative to the baseline rate around the post. Figure 3.4(b) shows the resulting modulation factor to the signal from in figure 3.3.

Figure 3.5: A measure of the total impact on Twitter by a press release can be obtained by integrating over the difference between a signal corresponding of the average tweet cosine similarity modulated by the relative activity, and a corresponding baseline. This measure does not depend on the typical Twitter activity and so can be used to compare the effectiveness of campaigns of varying size.

Figure 3.6: The cumulative deviation from the mean cosine similarity between the Deep Dreams post and tweets in the period around its release.

# APPENDIX

APPENDIX

# A.1   Social Fabric-related Code

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

## A.1.1   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:
```

```
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):
```

```
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.1.2 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
```

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

import glob
import numpy as np
import matplotlib.pyplot as plt

from social_fabric.phonetools import (load, make_filename)

#==============================================================================
display = True
#These are just  default − they're  updated when data is read
interval = 12*60**2
bin_size = 10*60
#==============================================================================


#Numpy−compliant method to convert activity and background to relative signal
get_signal = np.vectorize(lambda act, back : act * 1.0/back if back else 0)

def read_data(filename):
    '''Reads in data to plot. Updates the interval and bin sizes parameters.
     Returns data to be plotted as a hashmap with the following structure:
    {trigger : {event : signal}}.'''
    data = {}
    with open(filename, 'r') as f:
        (a,b) = load(f)
    first = True
    for key in a.keys():
        #Update interval and bin info
        if first:
            global interval
            global bin_size
            interval = a[key].interval
            bin_size = a[key].bin_size
            first = False

        (trigger, event) = key
        act = a[key].normalized()
        back = b[key].normalized()
        signal = get_signal(list(act), list(back))
        if not trigger in data:
            data[trigger] = {}
        data[trigger][event] = signal
    return data


def make_plot(trigger, signals):
    '''Generate a plot of relative user activity signal.'''
    legendstuff = []
    for event, vals in signals.iteritems():
        t = np.arange(−interval, interval, bin_size)
        t = map(lambda x: x*1.0/3600,  t)
        legendstuff.append(event)
        s = list(vals)
        plt.plot(t, s)
    plt.legend(legendstuff, loc='upper left')
    plt.xlabel('Time (h)')
    plt.ylabel('Relative signal')
    plt.title("Trigger: " + trigger)
    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.1.3  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


#===========================================================================
# 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_single_thread_numedgps.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 = []

#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))

if max_users:
```

```python
58        users = users[:max_users]
59
60   # Number of wallclock hours pr bin when fitting autoregressive series
61   hours_pr_ar_bin = 6
62   assert 24%hours_pr_ar_bin == 0
63
64   # Number of hours pr bin when compute daily rythm entropy
65   hours_pr_daily_rythm_bin = 1
66   assert 24%hours_pr_daily_rythm_bin == 0
67
68   # Time zone information
69   cph_tz = pytz.timezone('Europe/Copenhagen')
70
71   #Which data kinds to include
72   include_calls = True
73   include_ar = True
74   include_gps = True
75   include_network = False  #Allow geolocation from network data
76   include_bluetooth = True
77   include_facebook = True
78
79   #Whether to include not−a−number values in final output
80   allow_nan = True
81
82   #Threshold values to discard users with insufficient  data
83   minimum_number_of_texts = 10
84   minimum_number_of_calls = 5
85   minimum_number_of_gps_points = 100
86   minimum_number_of_facebook_friends = 1
87
88   #N_pings required to be considered social
89   bluetooth_social_threshold = 2
90
91   #Whether to output plots of cluster  analysis
92   plot_clusters = False
93
94   #=========================================================================
95   # Define helper methods
96   #=========================================================================
97
98   def get_distance(p, q):
99       return math.sqrt(sum([(p[i]−q[i])**2 for i in xrange(len(p))]))
100
101  def is_sorted(l):
102      return all([l[i+1] >= l[i] for i in xrange(len(l)−1)])
103
104  def get_entropy(event_list):
105      '''Takes a list of contacts from in/ or outgoing call/text events and
106       computes its entropy. event_list must be simply a list of user codes
107       corresponding to events.'''
108      n = len(event_list)
109      counts = Counter(event_list)
110      ent = sum([−v/n * math.log(v/n, 2) for v in counts.values()])
111      return ent
112
113  def next_time(dt, deltahours):
114      '''Accepts a datetime object and returns the next datetime object at which
115       the 'hour' count modulo deltahours is zero.
116       For example, deltahours = 6 gives the next time clock time is
117       0, 6, 12 or 18.
118       Sounds simple but is pretty annoying due to daylight saving time and so on,
119       so take care not to mess with this.'''
```

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

```python
182        return time_series
183
184
185 def timestamps2daily_entropy(timestamps, hours_pr_bin):
186        '''Constructs a histogram of hour-values of the imput timestamps
187         and computes its entropy.'''
188        if not 24%hours_pr_bin == 0:
189            raise ValueError("24 must be divisible by hours_pr_bin.")
190        bins = {}
191        for timestamp in timestamps:
192            hour = epoch2dt(timestamp).hour
193            _bin = int(hour/hours_pr_bin)
194            try:
195                bins[_bin] += 1
196            except KeyError:
197                bins[_bin] = 1
198            #
199        total = sum(bins.values())
200        entropy = sum([-v/total*math.log(v/total, 2) for v in bins.values()])
201        return entropy
202
203 #Make sure we have a blank file to write to.
204 open(output_filename, 'w').close()
205 assert os.stat(output_filename).st_size == 0  #Check that it worked.
206
207 user_counter = 0
208
209 for user in users:
210     user_counter += 1
211     msg = "Processing user %s of %s: %s" % (user_counter, len(users), user)
212     print msg
213     if user_counter % 100 == 0:
214         pushme(msg)
215
216     #Dict to hold all data extracted on current user
217     data = {}
218
219     #Try to load user psych profile- Discard user if they're not in file
220     try:
221         profile = user2profile[user]
222     except KeyError:
223         print "No psychological data on user."
224         continue  # No questionaire data. Shouldn't happen.
225
226     #Read in calls/texts
227     calls = readncheck(userfile_path+user+"/call_log.txt")
228
229     if len(calls) < minimum_number_of_calls:
230         print "too few calls"
231         continue
232     texts = readncheck(userfile_path+user+"/sms_log.txt")
233     if len(texts) < minimum_number_of_texts:
234         print "too few texts"
235         continue
236
237     #Extract list of times for calls, texts and combination
238     call_times = sorted([d['timestamp'] for d in calls])
239     text_times = sorted([d['timestamp'] for d in texts])
240     call_text_times = sorted(call_times + text_times)
241
242     #Get calls from the first three months
243     tmin = call_times[0]
```

```python
244     tmax = call_times[0] + 60*60*24*30*3
245     early_calls = [c['number'] for c in calls if tmin<=c['timestamp']<=tmax]
246
247     #Repeat for texts
248     tmin = text_times[0]
249     tmax = text_times[0] + 60*60*24*30*3
250     early_texts = [t['body'] for t in texts if tmin<=t['timestamp']<=tmax]
251
252     #Get number of unique contacts for first 3 months and append to data.
253     uniques = len(set(early_calls + early_texts))
254     data['n_contacts_first_three_months'] = uniques
255
256     #Compute daily entropy for calls and texts
257     data['call_daily_entropy'] = timestamps2daily_entropy(call_times,
258                                         hours_pr_daily_rythm_bin)
259     data['text_daily_entropy'] = timestamps2daily_entropy(text_times,
260                                         hours_pr_daily_rythm_bin)
261
262     #Compute median and std for call durations
263     call_durations = [c['duration'] for c in calls if not c['duration'] == 0]
264     data['call_duration_med'] = np.median(call_durations)
265     data['call_duration_std'] = np.std(call_durations)
266
267 #==============================================================================
268 # Crunch time-series info
269 #==============================================================================
270     if include_ar:
271         #Grap a sorted list of only times of events caused by user
272         outgoing_stuff = sorted([d['timestamp'] for d in calls
273                                 if d['type'] == 2] + [d['timestamp']
274                                 for d in texts if d['type'] == 2])
275
276         n_params = int(24*7/hours_pr_ar_bin + 1) #1 week plus 1 extra bin
277
278         #-Fit time series and extract parameters
279         try:
280             #Convert into timezone aware datetime objects
281             dts = [epoch2dt(timestamp) for timestamp in outgoing_stuff]
282
283             time_series = make_time_series(dts, hours_pr_ar_bin)
284
285             model = ar_model.AR(time_series)
286             result = model.fit(n_params)
287             #Grab parameters from fitted model
288             params = result.params[1:]
289             while len(params) < n_params:
290                 params.append(float('nan'))
291         except:
292             if not allow_nan:
293                 continue
294             else:
295                 params = [float('nan') for _ in xrange(n_params)]
296
297         #Append AR-coefficients to user data
298         count = 0
299         for par in params:
300             count += 1
301             name = "outgoing_activity_AR_coeff_"+str(count)
302             data[name] = par
303
304         #Get autocorrelation coefficients as well
305         try:
```

```
306                 accs = get_autocorrelation_coefficients(time_series, n_params)[:1]
307             except:
308                 if not allow_nan:
309                     continue
310                 accs = [float('nan') for _ in xrange(n_params + 1)]
311
312             #Append autocorrelation coefficients to user data
313             for i in xrange(len(accs)):
314                 name = "outgoing_activity_acc_"+str(i)
315                 data[name] = accs[i]
316
317             #Repeat with incoming signals. Might be interesting.
318             incoming_stuff = sorted([d['timestamp'] for d in calls
319                             if d['type'] == 1] + [d['timestamp']
320                             for d in texts if d['type'] == 1])
321
322             try:
323                 # Convert into timezone aware datetime objects
324                 dts = [epoch2dt(timestamp) for timestamp in incoming_stuff]
325
326                 time_series = make_time_series(dts, hours_pr_ar_bin)
327                 model = ar_model.AR(time_series)
328                 result = model.fit(n_params)
329                 params = result.params[1:]
330                 while len(params) < n_params:
331                     params.append(float('nan'))
332             except:
333                 if not allow_nan:
334                     continue
335                 else:
336                     params = [float('nan') for _ in xrange(n_params)]
337
338             # Name each of them and append to user data
339             count = 0
340             for par in params:
341                 count += 1
342                 name = "incoming_activity_AR_coeff_"+str(count)
343                 data[name] = par
344
345             #Get autocorrelation coefficients as well
346             try:
347                 accs = get_autocorrelation_coefficients(time_series, n_params)[:1]
348             except:
349                 if not allow_nan:
350                     continue
351                 accs = [float('nan') for _ in xrange(n_params + 1)]
352
353             for i in xrange(len(accs)):
354                 name = "incoming_activity_acc_"+str(i)
355                 data[name] = accs[i]
356
357
358     #==============================================================================
359     #     Crunch call/text info
360     #==============================================================================
361
362     if include_calls:
363         #Add values to temporary data map.
364         d = {'call': call_times, 'text': text_times, 'ct': call_text_times}
365         for label, times in d.iteritems():
366             timegaps = [times[i+1] - times[i] for i in xrange(len(times)-1)]
367             timegaps = filter(lambda x: x < 259200, timegaps) #3 days, tops
```

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

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

```
492            data['text_outgoing'] = sum([t['type'] == 2 for t in texts])/len(texts)
493
494    #========================================================================
495    # Crunch location data
496    #========================================================================
497
498        if include_gps:
499            with open(userfile_path+user+"/gps_log.txt", 'r') as f:
500                raw = [LE(line) for line in f.readlines()]
501
502            if not is_sorted([l['timestamp'] for l in raw]):
503                raw = sort_dicts_by_key(raw, 'timestamp')
504
505            #We only want measurements taken at least 500s apart.
506            prev = 0
507            gps_data = []
508            allowed_providers = ['gps', 'network'] if include_network else ['gps']
509            for line in raw:
510                now = line['timestamp']
511                if line['provider'] in allowed_providers and now - prev >= 500:
512                    #Convert coordinates to km and note it down
513                    x = (longitude2meters*line['lon'] - x_offset)*0.001
514                    y = (latitude2meters*line['lat'] - y_offset)*0.001
515                    gps_data.append({'point' : (x,y), 'timestamp' : now,
516                                     'smarttime' : epoch2dt(now)})
517                    prev = now
518                #
519            #ignore user if there aren't enough data
520            if not len(gps_data) >= minimum_number_of_gps_points:
521                continue
522
523            # We want to investigate each day saparately so start at midnight.
524            first_midnight = next_midnight(gps_data[0]['smarttime'])
525            for i in xrange(len(gps_data)):
526                if gps_data[i]['smarttime'] >= first_midnight:
527                    gps_data = gps_data[i:]
528                    break
529
530
531            #Generate list of radii of smallest enclosing circle, SEC, for each day
532            current_points = []
533            prev = gps_data[0]['timestamp']
534            radii = []
535
536            distances = []
537            early_day = gps_data[0]['smarttime']
538            late_day = next_midnight(early_day)
539            for datum in gps_data:
540                now = datum['smarttime']
541                while not early_day <= now < late_day:
542                    if len(current_points) > 2:
543                        crds = [p['point'] for p in current_points]
544                        circle = make_circle(crds)
545                        r = circle[2] if circle else 0
546                        if circle and r > 0:
547                            if r <= 500:
548                                radii.append(r)
549                            distances.append(sum([get_distance(crds[i], crds[i+1])
550                                                  for i in xrange(len(crds)-1)]))
551                    # Reset counters and update bins
552                    current_points = []
553                    early_day = late_day
```

```
554                        late_day = next_midnight(late_day)
555                    current_points.append(datum)
556
557
558         data['radius_of_gyration_med'] = np.median(radii)
559         data['radius_of_gyration_std'] = np.std(radii)
560         data['travel_med'] = np.median(distances)
561         data['travel_std'] = np.std(distances)
562
563         #Run Lloyd's algorithm to identify clusters
564         #Determine which points are stationary - less movement than 100m
565         stationary_data = []
566         try:
567             for i in xrange(1,len(gps_data)-1):
568                 a,b,c = tuple([gps_data[ind]['point'] for ind in
569                                 [i-1, i, i+1]])
570                 if (get_distance(a, b) < 0.1 and get_distance(b, c) < 0.1):
571                     stationary_data.append(gps_data[i])
572                 #
573             initial_clusters = 50
574             threshold_percent = 0.05
575             points = [elem['point'] for elem in stationary_data]
576
577             minimum_points = int(threshold_percent*len(points))
578
579
580             clusters_scatter = lloyds.lloyds(points, initial_clusters, runs=3,
581                                     init='scatter')
582
583             clusters_sample = lloyds.lloyds(points, initial_clusters, runs=3,
584                                     init='sample')
585
586             #Determine most succesful method
587             locs = lambda c: [p for p in c.values() if len(p)>=minimum_points]
588             if len(locs(clusters_scatter)) >= len(locs(clusters_sample)):
589                 method = 'scatter'
590                 best = clusters_scatter
591             else:
592                 method = 'sample'
593                 best = clusters_sample
594
595             # Number of places which survive cutoff (true = 1, so just sum)
596             n_places = sum([len(pl) >= minimum_points for pl in best.values()])
597
598             #Output plots of clustering
599             if plot_clusters:
600                 if not os.path.isdir('pics'):
601                     os.mkdir('pics')
602                 for ext in ['.pdf', '.png']:
603                     filename = 'pics/'+user+"_"+method+ext
604                     lloyds.draw_clusters(clusters = best,
605                                     threshold = minimum_points,
606                                     show = False,
607                                     filename = filename)
608         except:
609             if not allow_nan:
610                 continue
611             n_places = float('nan')
612         data['n_places'] = n_places
613
614         #Compute location entropy and add to data
615         try:
```

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

```python
678 #===============================================================================
679 # Facebook data
680 #===============================================================================
681
682     if include_facebook:
683         with open(userfile_path+user+'/facebook_log.txt', 'r') as f:
684             n = len(f.readlines())
685             if n < minimum_number_of_facebook_friends:
686                 continue
687             data['number_of_facebook_friends'] = n
688
689 #===============================================================================
690 # Bluetooth data
691 #===============================================================================
692
693     if include_bluetooth:
694         with open(userfile_path+user+"/bluetooth_log.txt", 'r') as f:
695             raw = [LE(line) for line in f.readlines()]
696         #Make sure data is sorted chronologically
697         if not is_sorted([entry['timestamp'] for entry in raw]):
698             raw = sort_dicts_by_key(raw, 'timestamp')
699         #List of times when user was social
700         social_times = []
701         total_social_time = 0
702         total_time = 0
703         #maps from each other user encountered to time spend with said user
704         friend2time_spent = {}
705         #Temporary variables
706         current_time = 0
707         previous_time = 0
708         current_users = []
709         previous_users = []
710         for signal in raw:
711             new_time = signal['timestamp']
712             #Check if line represents a new signal and if so, update values
713             if new_time != current_time:
714                 dt = new_time - current_time
715                 #Determine number of pings
716                 overlap = set(previous_users).intersection(set(current_users))
717                 if len(overlap) >= bluetooth_social_threshold:
718                     social_times.append(previous_time)
719                     if dt <= 7200:
720                         total_social_time += dt
721                         total_time += dt
722                         for friend in overlap:
723                             try:
724                                 friend2time_spent[friend] += dt
725                             except KeyError:
726                                 friend2time_spent[friend] = dt
727                         #
728                 elif dt <= 7200:
729                     total_time += dt
730                 #Update variables
731                 previous_time = current_time
732                 previous_users = current_users
733                 current_users = []
734                 current_time = new_time
735             new_user = signal['name']
736             if new_user=='-1' or not new_user:
737                 continue
738             else:
739                 current_users.append(new_user)
```

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

```
802          f.write("\n")
803
804  #Done.
805  pushme("Data extraction done.")
```

## A.1.4    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 os
18  import random
19  from scipy.sparse import dok_matrix
20  from sklearn import svm
21  from sklearn.ensemble import RandomForestRegressor, RandomForestClassifier
22  from sklearn.cross_validation import (cross_val_score, LeaveOneOut, KFold,
23                                         StratifiedKFold)
24  import sys
25  import traceback
26
27  BASE = os.path.dirname(__file__)
28
29  oldhat = (35/256,39/256,135/256)
30  nude = (203/256,150/256,93/256)
31  wine = (110/256,14/256,14/256)
32  moerkeroed = (156/256,30/256,36/256)
33
34
35  #Empirically determined optimal hyperparameters for SVRs and RFs
36  with open(BASE+'/svr_parameters.json','r') as f:
37      svr_parameters = json.load(f)
38
39  wrbf_parameters = svr_parameters #updater!!!
40
41  with open(BASE+'/rf_parameters.json','r') as f:
42      rf_parameters = json.load(f)
43
44  def _make_colormap(seq):
45      """Return a LinearSegmentedColormap
46       seq: a sequence of floats and RGB-tuples. The floats should be increasing
47       and in the interval (0,1).
48       """
49      seq = [(None,) ∗ 3, 0.0] + list(seq) + [1.0, (None,) ∗ 3]
50      cdict = {'red': [], 'green': [], 'blue': []}
```

```
51      for i, item in enumerate(seq):
52          if isinstance(item, float):
53              r1, g1, b1 = seq[i − 1]
54              r2, g2, b2 = seq[i + 1]
55              cdict['red'].append([item, r1, r2])
56              cdict['green'].append([item, g1, g2])
57              cdict['blue'].append([item, b1, b2])
58      return mcolors.LinearSegmentedColormap('CustomMap', cdict)

color_map = _make_colormap([oldhat, moerkeroed, 0.33, moerkeroed, nude, 0.67, nude])

big_five = ['openness', 'conscientiousness', 'extraversion', 'agreeableness',
            'neuroticism']

#_default_features = ["n_texts",
#                        "ct_iet_std ",
#                        " call_cir ",
#                        "call_entropy",
#                        " text_cir ",
#                        "n_calls ",
#                        "text_latency ",
#                        "call_outgoing",
#                        "fraction_social_time ",
#                        "text_outgoing",
#                        " call_iet_std ",
#                        "n_text_contacts",
#                        " call_night_activity ",
#                        "call_iet_med",
#                        "outgoing_activity_AR_coeff_2",
#                        "text_entropy",
#                        " ct_cir ",
#                        "text_response_rate",
#                        "n_ct_contacts",
#                        "social_entropy",
#                        "n_call_contacts ",
#                        "n_ct",
#                        " text_iet_std ",
#                        "ct_iet_med",
#                        "ct_entropy",
#                        "text_iet_med",
#                        "call_response_rate",
#                        "number_of_facebook_friends"]

_default_features = ['call_iet_med', 'text_iet_med', 'social_entropy',
'call_entropy', 'travel_med', 'n_places', 'text_latency',
'call_night_activity']

def split_ntiles(values, n):
    '''Determines the values that separate the imput list into n equal parts.
     this is a generalization of the notion of median (in the case n = 2) or
     quartiles (n=4).
     Usage: ntiles([5,6,7], 2) gives [6] for instance.'''
    result = []
    for i in xrange(1,n):
        percentile = 100/n * i
        result.append(np.percentile(values, percentile,
                                    interpolation='linear'))
    return result

def determine_ntile(value, ntiles):
    '''Determines which n-tile the input value belongs to.
     Usage: determine_ntile([7,9,13], 10) gives 2 (third quartile).
```

```python
113          This uses zero indexing so data split into e.g. quartiles will give results
114          like 0,1,2,3 - NOT 1,2,3,4.'''
115       #Check if value is outside either extreme, meaning n-tile 1 or n.
116       if value >= ntiles[-1]:
117           return len(ntiles) #Remember the length is n-1
118       elif value < ntiles[0]:
119           return 0 #Values was in the first n-tile
120       # Define possible region and search for where value is between two elements
121       left = 0
122       right = len(ntiles)-2
123       #Keep checking th middle of the region and updating region
124       ind = (right + left)//2
125       while not ntiles[ind] <= value < ntiles[ind + 1]:
126           #Check if lower bound tile is on the left
127           if value < ntiles[ind]:
128               right = ind - 1
129           else:
130               left = ind + 1
131           ind = (right + left)//2
132       # Being between ntiles 0 and 1, means second n-tile and so on.
133       return ind + 1


135   def assign_labels(Y, n):
136       '''Accepts a list and an int n and returns a list of discrete labels
137        corresponding to the ntile each original y-value was in.'''
138       ntiles = split_ntiles(Y, n)
139       labels = [determine_ntile(y, ntiles) for y in Y]
140       return labels


142   def normalize_data(list_):
143       '''Normalizes input data to the range [-1, 1]'''
144       lo, hi = min(list_), max(list_)
145       if lo == hi:
146           z = len(list_)*[0]
147           return z
148       else:
149           return [2*(val-lo)/(hi - lo) - 1 for val in list_]



152   def read_data(filename, trait, n_classes = None, normalize = True,
153                   features='default', interpolate = True):
154       '''This reads in a preprocessed datafile, splits psych profile data into
155       n classes if specified, filters desired psychological traits and
156       features and returns as a tuple (X,Y, indexdict), which can be fed to a'
157       number of off-the-shelf ML schemes.
158       If trait=='Sex', female and male are converted to 0 and 1, respectively.
159       indexdict maps each element of the feature vectors to their label, as in
160       {42 : 'distance_travelled_pr_day'} etc.

162       Args:
163         filename : str
164           Name of the file containing the data.
165         trait : str
166           The psychological trait to extract data on.
167         n_classes : int
168           Number of classes to split data into. Default is None,
169           i.e. just keep the decimal values. Ignored if trait == 'Sex', as data
170           only has two discreet values.
171         normalize : bool
172           Whether to hard normalize data to [-1, 1].
173         features : str/list
174           Which features to read in. Can also be 'all'
```

```
175            or 'default', meaning the ones I've pragmatically found to be
176            reasonable.
177         interpolate : bool:
178            Whether to replace NaN's with the median value of
179            the feature in question.'''
180     if trait == 'sex':
181         n_classes = None
182     #Read in the raw data
183     with open(filename, 'r') as f:
184         raw = [json.loads(line) for line in f.readlines()]
185     #Get list of features to be included – everything if nothing's specified
186     included_features = []
187     if features == 'default':
188         included_features = _default_features
189     elif features == 'all':
190         included_features = raw[0]['data'].keys() if features=='all' else features
191     else:
192         included_features = features
193
194     #Remove any features that only have NaN values.
195     for i in xrange(len(included_features)−1,−1,−1):
196         feat = included_features[i]
197         if all(math.isnan(line['data'][feat]) for line in raw):
198             del included_features[i]
199
200
201     # −−−−− Handle feature vectors −−−−−
202     # Dict mapping indices to features
203     indexdict = {ind : feat for ind, feat in enumerate(included_features)}
204     # N_users x N_features array to hold data
205     rows = len(raw)
206     cols = len(included_features)
207     X = np.ndarray(shape = (rows, cols))
208     for i in xrange(rows):
209         line = raw[i]
210         #Construct data matrix
211         for j, feat in indexdict.iteritems():
212             val = line['data'][feat]
213             X[i, j] = val
214         #
215     #Replace NaNs with median values
216     if interpolate:
217         for j in xrange(cols):
218             #Get median of feature j
219             med = np.median([v for v in X[:, j] if not math.isnan(v)])
220             if math.isnan(med):
221                 raise ValueError('''Feature %s contains only NaN's and should
222                                 have been removed.''' % indexdict[j])
223             for i in xrange(rows):
224                 if math.isnan(X[i,j]):
225                     X[i,j] = med
226
227     if normalize:
228         for j in xrange(cols):
229             col = X[:, j]
230             X[:, j] = normalize_data(col)
231
232     # −−−−− Handle class info −−−−−
233     trait_values = []
234     for line in raw:
235         #Add value of psychological trait
236         psych_trait = line['profile'][trait]
```

```
237            if trait == 'Sex':
238                if psych_trait == 'Female':
239                    psych_trait = 0
240                elif psych_trait == 'Male':
241                    psych_trait = 1
242                else:
243                    raise ValueError('My code is binary gender normative, sorry.')
244            trait_values.append(psych_trait)
245        Y = []
246        if n_classes == None:
247            Y = trait_values
248        else:
249            ntiles = split_ntiles(trait_values, n_classes)
250            Y = [determine_ntile(tr, ntiles) for tr in trait_values]
251
252        return (X, Y, indexdict)
253
254    def get_strong_indices(X, Y, threshold):
255        '''Returns the indices of the features vectors X whose linear correlation
256         with Y is above the imput threshold.'''
257        correlations = get_correlations(X, Y)
258        indices = [i for i in xrange(len(correlations))
259                    if correlations[i] >= threshold]
260        return indices
261
262    def reduce_data(X, indices):
263        '''Return feature vectors including only the feature numbers in the input
264         list of indices.'''
265        reduced_X = np.array([[x[i] for i in indices] for x in X])
266        return reduced_X
267
268    def plot_stuff(input_filename, output_filename=None, color=moerkeroed):
269        with open(input_filename, 'r') as f:
270            d = json.load(f)
271        x = d['x']
272        y = d['y']
273        yerr = d['mean_stds']
274        plt.plot(x,y, color=color, linestyle='dashed',
275                marker='o')
276        plt.errorbar(x, y, yerr=yerr, linestyle="None", marker="None",
277                    color=color)
278        if output_filename:
279            plt.savefig(output_filename)
280
281    def get_TPRs_and_FPRs(X, Y, forest = None, verbose = False):
282        '''Accepts a list of feature vectors and a list of labels and returns a
283         tuple of true positive and false positive rates (TPRs and FPRs,
284         respectively) for various confidence thresholds.'''
285        kf = LeaveOneOut(n=len(Y))
286
287        results = []
288        thresholds = []
289
290        counter = 0
291        for train, test in kf:
292            counter += 1
293            if counter % 10 == 0 and verbose:
294                print "Testing on user %s of %s..." % (counter, len(Y))
295
296            result = {}
297            train_data = [X[i] for i in train]
298            train_labels = [Y[i] for i in train]
```

```python
299        test_data = [X[i] for i in test]
300        test_labels = [Y[i] for i in test]
301
302        if not forest:
303            forest = RandomForestClassifier(
304                                n_estimators = 1000,
305                                n_jobs=-1,
306                                criterion='entropy')
307
308        forest.fit(train_data, train_labels)
309        result['prediction'] = forest.predict(test_data)[0]
310        result['true'] = test_labels[0]
311        confidences = forest.predict_proba(test_data)[0]
312        result['confidences'] = confidences
313        thresholds.append(max(confidences))
314
315        results.append(result)
316
317    #ROC curve stuff - false and true positive rates
318    TPRs = []
319    FPRs = []
320
321    unique_thresholds = sorted(list(set(thresholds)), reverse=True)
322
323    for threshold in unique_thresholds:
324        tn = 0
325        fn = 0
326        tp = 0
327        fp = 0
328        for result in results:
329            temp = result['prediction']
330            if temp == 1 and result['confidences'][1] >= threshold:
331                pred = 1
332            else:
333                pred = 0
334            if pred == 1:
335                if result['true'] == 1:
336                    tp += 1
337                else:
338                    fp += 1
339                #
340            elif pred == 0:
341                if result['true'] == 0:
342                    tn += 1
343                else:
344                    fn += 1
345                #
346            #
347        TPRs.append(tp/(tp + fn))
348        FPRs.append(fp/(fp + tn))
349    return (TPRs, FPRs)
350
351 def make_roc_curve(TPRs, FPRs, output_filename = None):
352    '''Accepts a list of true and false positive rates (TPRs and FPRs,
353     respectively) and generates a ROC-curve.'''
354    predcol = moerkeroed
355    basecol = oldhat
356    fillcol = nude
357
358    fig = plt.figure()
359    ax = fig.add_subplot(1,1,1)
360
```

```
361        TPRs = [0] + TPRs + [1]
362        FPRs = [0] + FPRs + [1]
363
364        area = 0.0
365        for i in xrange(len(TPRs)−1):
366            dx = FPRs[i+1] − FPRs[i]
367            y = 0.5*(TPRs[i] + TPRs[i+1])
368            under_curve = dx*y
369            baseline = dx*0.5*(FPRs[i] + FPRs[i+1])
370            area += under_curve − baseline
371
372        baseline = FPRs
373        ax.fill_between(x = FPRs, y1 = TPRs, y2 = baseline, color = fillcol,
374                        interpolate = True, alpha=0.8)
375        ax.plot(baseline, baseline, color = basecol, linestyle = 'dashed',
376                linewidth = 1.0, label = 'Baseline')
377        ax.plot(FPRs, TPRs, color=predcol, linewidth = 1,
378                        label = 'Prediction')
379
380        plt.xlabel('False positive rate.')
381        plt.ylabel('True positive rate.')
382
383        handles, labels = ax.get_legend_handles_labels()
384        hest = mpatches.Patch(color=fillcol)
385
386        labels += ['Area = %.3f' % area]
387        handles += [hest]
388        ax.legend(handles, labels, loc = 'lower right')
389 #      plt.legend(handles = [tp_line, base])
390        if output_filename:
391            plt.savefig(output_filename)
392        plt.show()
393
394 def rank_features(X, Y, forest, indexdict, limit = None):
395     '''Ranks the features of a given dataset and classifier.
396      indexdict should be a map from indices to feature names like
397      {0 : 'average_weigth'} etc.
398      if limit is specified, this method returns only the top n ranking features.
399      Returns a dict like {'feature name' : (mean importances, std)}.'''
400     importances = forest.feature_importances_
401     stds=np.std([tree.feature_importances_ for tree in forest.estimators_],
402                axis=0)
403     indices = np.argsort(importances)[::−1]
404     if limit:
405         indices = indices[:limit]
406     d = {indexdict[i] : (importances[i], stds[i]) for i in indices}
407     return d
408
409 def check_performance(X, Y, clf, strata = None):
410     '''Checks forest performance compared to baseline.'''
411     N_samples = len(Y)
412     #Set up validation indices
413     if not strata: #Do leave−one−out validation
414         skf = KFold(N_samples, n_folds=N_samples, shuffle = False)
415     else: #Do stratified K−fold
416         skf = StratifiedKFold(Y, n_folds=strata)
417
418     #Evaluate classifier performance
419     scores = []
420     for train, test in skf:
421         train_data = [X[ind] for ind in train]
422         train_labels = [Y[ind] for ind in train]
```

```
423            test_data = [X[ind] for ind in test]
424            test_labels = [Y[ind] for ind in test]
425
426            #Check performance of input forest
427            clf.fit(train_data, train_labels)
428            score = clf.score(test_data, test_labels)
429            scores.append(score)
430
431        #Compute baseline
432        most_common_label = max(Counter(Y).values())
433        baseline = float(most_common_label)/N_samples
434
435        #Compare results with prediction baseline
436        score_mean = np.mean(scores)/baseline
437        score_std = np.std(scores)/baseline
438
439        return (score_mean, score_std)
440
441    def check_regressor(X, Y, reg, strata = None):
442        '''Checks the performance of a regressor against mean value baseline.'''
443        N_samples = len(Y)
444        #Set up validation indices
445        if not strata:  #Do leave−one−out validation
446            skf = KFold(N_samples, n_folds=N_samples,
447                                        shuffle = False)
448        else:  #Do stratified K−fold
449            skf = StratifiedKFold(Y, n_folds=strata)
450
451        #Evaluate performance
452        model_abs_errors = []
453        baseline_abs_errors = []
454        for train, test in skf:
455            train_data = [X[ind] for ind in train]
456            train_labels = [Y[ind] for ind in train]
457            test_data = [X[ind] for ind in test]
458            test_labels = [Y[ind] for ind in test]
459
460            #Check performance of input forest
461            reg.fit(train_data, train_labels)
462            base = np.mean(train_labels)
463            for i in xrange(len(test_data)):
464                pred = reg.predict(test_data[i])
465                true = test_labels[i]
466                model_abs_errors.append(np.abs(pred − true))
467                baseline_abs_errors.append(np.abs(base − true))
468            #
469        return (np.mean(model_abs_errors), np.mean(baseline_abs_errors))
470
471    class _RFNN(object):
472        __metaclass__ = abc.ABCMeta
473        '''Abstract class for random forest nearest neightbor predictors.
474         This should never be instantiated.'''
475
476        def __init__(self, forest, n_neighbors):
477            self.forest = forest
478            self.n_neighbors = n_neighbors
479            self.X = None
480            self.Y = None
481
482        def fit(self, X, Y):
483            '''Fits model to training data.
484
```

```python
        Args
        -----------
        X : List
           List of training feature vectors.

        Y : List
           List of training labels or values to be predicted.'''
        if not len(X) == len(Y):
            raise ValueError("Training input and output lists must have "
                             "same length.")
        if not self.n_neighbors <= len(X):
            raise ValueError("Fewer data points than neighbors.")

        self.forest.fit(X, Y)
        self.X = X
        self.Y = Y

    def _rf_similarity(self, a, b):
        '''Computes a similarity measure for two points using a trained random
        forest classifier.'''
        if self.X == None or self.Y == None:
            raise NotImplementedError("Model has not been fittet to data yet.")

        #Feature vectors must be single precision.
        a = np.array([a], dtype = np.float32)
        b = np.array([b], dtype = np.float32)
        hits = 0
        tries = 0
        for estimator in self.forest.estimators_:
            tries += 1
            tree = estimator.tree_
            # Check whether the points end up on the same leaf for this  tree
            if tree.apply(a) == tree.apply(b):
                hits += 1
            #
        return hits/tries

    def find_neighbors(self, point):
        '''Determine the n nearest nieghbors for the given point.
         Returns a list of n tuples like (yval, similarity).
         The tuples are sorted descending by similarity.'''
        if self.X == None or self.Y == None:
            raise NotImplementedError("Model has not been fittet to data yet.")

        #Get list  of tuples like  (y, similarity) for the n 'nearest' points
        nearest = [(None,float('-infinity')) for _ in xrange(self.n_neighbors)]
        for i in xrange(len(self.X)):
            similarity = self._rf_similarity(self.X[i], point)
            # update top n list  if  more similar than the  furthest  neighbor
            if similarity > nearest[-1][1]:
                nearest.append((self.Y[i], similarity))
                nearest.sort(key = lambda x: x[1], reverse = True)
                del nearest[-1]
            #
        return nearest

    #Mandatory methods – must be overridden
    @abc.abstractmethod
    def predict(self, point):
        pass

    @abc.abstractmethod
```

```python
547      def score(self, X, Y):
548          pass
549
550  def _reservoir_sampler(start = 1):
551      '''Generator of the probabilities need to do reservoir sampling. The point
552       it that this can be used to iterate through a list, discarding each element
553       for the following element with probability P_n and ending up with a random
554       element from the list.'''
555      n = start
556      while True:
557          p = 1/n
558          r = random.uniform(0,1)
559          if r < p:
560              yield True
561          else:
562              yield False
563          n += 1
564
565  class RFNNClassifier(_RFNN):
566      '''Random Forest Nearest Neighbor Classifier.
567
568      Parameters
569      ----
570      n_neighbors : int
571        Number of neighbors to consider.
572
573      forest : RandomForestClassifier
574        The forest which will provide a
575        distance measure on which determine nearest neighbors.
576
577      weighting : str
578        How to weigh the votes of different neighbors.
579        'equal' means each neighbor has an equivalent vote.
580        'linear' mean votes are weighed by their similarity to the input point.
581      '''
582      def predict(self, point):
583          '''Predicts the label of a given point.'''
584          neighbortuples = self.find_neighbors(point)
585          if self.weighting == 'equal':
586              #Simple majority vote. Select randomly if it's a tie.
587              predictions = [t[0] for t in neighbortuples]
588              best = 0
589              winner = None
590              switch = _reservoir_sampler(start = 2)
591              for label, votes in Counter(predictions).iteritems():
592                  if votes > best:
593                      best = votes
594                      winner = label
595                      switch = _reservoir_sampler(start = 2)
596                  elif votes == best:
597                      if switch.next():
598                          winner = label
599                      else:
600                          pass
601                  else:
602                      pass
603                  #
604              return winner
605
606          #Weigh votes by their similarity to the input point
607          elif self.weighting == 'linear':
608              #The votes are weighted by their similarity
```

```python
609                d = {}
610                for yval, similarity in neighbortuples:
611                    try:
612                        d[yval] += similarity
613                    except KeyError:
614                        d[yval] = similarity
615                best = float('-infinity')
616                winner = None
617                for k, v in d.iteritems():
618                    if v > best:
619                        best = v
620                        winner = k
621                    else:
622                        pass
623                return winner
624
625        def score(self, X, Y):
626            if not len(X) == len(Y):
627                raise ValueError("Training data and labels must have same length.")
628
629            hits = 0
630            n = len(X)
631            for i in xrange(n):
632                pred = self.predict(X[i])
633                if pred == Y[i]:
634                    hits += 1
635                #
636            return hits/n
637
638
639        def __init__(self, forest = None, n_neighbors = 3, weighting = 'equal',
640                     n_jobs = 1):
641            #Make sure we have a forest * classifier *
642            if forest == None:
643                forest = RandomForestClassifier(n_estimators = 1000,
644                                                criterion = 'entropy',
645                                                n_jobs = n_jobs)
646            if not isinstance(forest, RandomForestClassifier):
647                raise TypeError("Forest must be a classifier")
648
649            self.weighting = weighting
650
651            #Call parent constructor
652            super(RFNNClassifier, self).__init__(forest, n_neighbors)
653
654 class RFNNRegressor(_RFNN):
655        '''Random Forest Nearest Neighbor Regressor.
656
657        Parameters
658        ----
659        n_neighbors : int
660          Number of neighbors to consider.
661
662        forest : RandomForestRegressor
663          The forest which will provide a
664          distance measure on which determine nearest neighbors.
665
666        weighting : str
667          How to weigh the votes of different neighbors.
668          'equal' means each neighbor has an equivalent weight.
669          'linear' mean votes are weighed by their similarity to the input point.
670        '''
```

```python
671     def predict(self, point):
672         # lists of the y vaues and similarities of nearest neighbors
673         neighbortuples = self.find_neighbors(point)
674         yvals, similarities = zip(*neighbortuples)
675
676         # Weigh each neighbor y value equally is that's how we roll
677         if self.weighting == 'equal':
678             weight = 1.0/len(yvals)
679             result = 0.0
680             for y in yvals:
681                 result += y*weight
682             return result
683             #
684         # Otherwise, weigh neighbors by similarity
685         elif self.weighting == 'linear':
686             weight = 1.0/(sum(similarities))
687             result = 0.0
688             for i in xrange(len(yvals)):
689                 y = yvals[i]
690                 similarity = similarities[i]
691                 result += y*similarity*weight
692             return result
693
694
695     def score(self, X, Y):
696         if not len(X) == len(Y):
697             raise ValueError("X and Y must be same length.")
698         errors = [Y[i] - self.predict(X[i]) for i in xrange(len(X))]
699         return np.std(errors)
700
701
702     def __init__(self, forest = None, n_neighbors = 3, weighting = 'equal'):
703         #Check forest type.
704         if forest == None:
705             forest = RandomForestRegressor(n_estimators = 1000, n_jobs = -1)
706         if not isinstance(forest, RandomForestRegressor):
707             raise TypeError("Must use Random Forest Regressor to initialize.")
708         # Set params
709         self.weighting = weighting
710         # Done. Call parent constructor
711         super(RFNNRegressor, self).__init__(forest, n_neighbors)
712
713
714 class _BaselineRegressor(object):
715     '''Always predicts the mean of the training set.'''
716     def __init__(self, guess=None):
717         self.guess = guess
718     def fit(self, xtrain, ytrain):
719         '''Find the average of input lidt of target values and guess on that
720          from now on.'''
721         self.guess = np.mean(ytrain)
722     def predict(self, x):
723         return self.guess
724
725
726 class _BaselineClassifier(object):
727     '''Always predicts the most common label in the training set'''
728     def __init__(self, guess=None):
729         self.guess = guess
730     def fit(self, xtrain, ytrain):
731         '''Find the most common label and guess on that from now on.'''
732         countmap = Counter(ytrain)
```

```python
733              best = 0
734              for label, count in countmap.iteritems():
735                  if count > best:
736                      best = count
737                      self.guess = int(label)
738                  #
739          #
740      def predict(self, x):
741          return self.guess


744  def _worker(X, Y, score_type, train_percentage, classifier, clf_args, n_groups,
745              replace, threshold):
746      '''Worker method for parallelizing bootstrap evaluations.'''
747      #Create bootstrap sample
748      try:
749          rand = np.random.RandomState()  #Ensures PRNG works in children
750          indices = rand.choice(xrange(len(X)), size = len(X), replace = replace)
751          xsample = [X[i] for i in indices]
752          ysample = [Y[i] for i in indices]
753
754          #Generate training and testing set
755          cut = int(train_percentage*len(X))
756          xtrain = xsample[:cut]
757          xtest = xsample[cut:]
758          ytrain = ysample[:cut]
759          ytest = ysample[cut:]
760
761          #Discard features with too low correlation with output vectors
762          inds = get_strong_indices(X=xtrain, Y=ytrain, threshold = threshold)
763          xtrain = reduce_data(xtrain, inds)
764          xtest = reduce_data(xtest, inds)
765
766          #Create regressor if we're doing regression
767          if classifier == 'RandomForestRegressor':
768              clf = RandomForestRegressor(**clf_args)
769          elif classifier == 'SVR':
770              clf = svm.SVR(**clf_args)
771          elif classifier == 'baseline_mean':
772              clf = _BaselineRegressor()
773          elif classifier == 'WRBFR':
774              importances = get_correlations(xtrain, ytrain)
775              clf_args['importances'] = importances
776              clf = WRBFR(**clf_args)
777
778          #Create classifier and split dataset into labels
779          elif classifier == 'RandomForestClassifier':
780              clf = RandomForestClassifier(**clf_args)
781              ysample = assign_labels(ysample, n_groups)
782              ytrain = ysample[:cut]
783              ytest = ysample[cut:]
784          elif classifier == 'SVC':
785              clf = svm.SVC(**clf_args)
786              ysample = assign_labels(ysample, n_groups)
787              ytrain = ysample[:cut]
788              ytest = ysample[cut:]
789          elif classifier == 'baseline_most_common_label':
790              clf = _BaselineClassifier()
791              ysample = assign_labels(ysample, n_groups)
792              ytrain = ysample[:cut]
793              ytest = ysample[cut:]
794          elif classifier == 'WRBFC':
```

```
795              importances = get_correlations(xtrain, ytrain)
796              clf_args['importances'] = importances
797              clf = WRBFC(**clf_args)
798              ysample = assign_labels(ysample, n_groups)
799              ytrain = ysample[:cut]
800              ytest = ysample[cut:]
801
802          #Fail if none of the above classifiers were specified
803          else:
804              raise ValueError('Regressor or classifier not defined.')
805
806          #Fit the classifier or regressor
807          clf.fit(xtrain, ytrain)
808
809          #Compute score and append to output list
810          if score_type == 'mse':
811              scores = [(ytest[i] - clf.predict(xtest[i]))**2
812                      for i in xrange(len(xtest))]
813              return np.mean(scores)
814
815          elif score_type == 'fraction_correct':
816              n_correct = sum([ytest[i] == int(clf.predict(xtest[i]))
817                              for i in xrange(len(ytest))])
818              score = n_correct/len(ytest)
819              return score
820          elif score_type == 'over_baseline':
821              #Get score
822              score = sum([ytest[i] == int(clf.predict(xtest[i]))
823                              for i in xrange(len(ytest))])
824              #Get baseline
825              baselineclf = _BaselineClassifier()
826              baselineclf.fit(xtrain, ytrain)
827              baseline = sum([ytest[i] == baselineclf.predict(xtest[i])
828                              for i in xrange(len(ytest))])
829              return score/baseline
830
831          #Fail if none of the above performance metrics were specified
832          else:
833              raise ValueError('Score type not defined.')
834
835          #Job's done!
836          return None
837      except:
838          raise Exception("".join(traceback.format_exception(*sys.exc_info())))
839
840
841  def bootstrap(X, Y, classifier, score_type = 'mse', train_percentage = 0.8,
842              clf_args = {}, iterations = 1000, n_groups = 3, n_jobs = 1,
843              replace = False, threshold = 0.0):
844      '''Performs bootstrap resampling to evaluate the performance of some
845      classifier or regressor. Note that this takes the *complete dataset* as
846      arguments as well as arguments specifying which predictor to use and which
847      function to estimate the distribution of.
848      This seems to be the most straightforward generalizable implementation
849      which can be parallelized, as passing e.g. the scoring function directly
850      clashed with the mechanisms implemented to work around the GIL for
851      multiprocessing for obscure reasons.
852
853      Parameters:
854      ----------------
855      X : list
856        All feature vectors in the complete dataset.
```

```
857
858      Y : list
859        All 'true' labels or output values in the complete dataset.
860
861      classifier : str
862        Which classifier to use to predict the test set. Allowed values:
863        'RandomForestRegressor', 'baseline_mean', 'SVR',
864        'RandomForestClassifier', 'SVC', 'baseline_most_common_label, WRBFR,
865        WRBFC'
866
867      score_type : str
868        String signifying which function to estimate the distribution of.
869        Allowed values: 'mse', 'fraction_correct', 'over_baseline'
870
871      train_percentage : float
872        The percentage [0:1] of each bootstrap sample to be used for training.
873
874      clf_args : dict
875        optional arguments to the constructor method of the regressor/classifier.
876
877      iterations : int
878        Number of bootstrap samples to run.
879
880      n_jobs : int
881        How many cores (maximum) to use.
882
883      replace : bool
884        Whether to sample with replacement when obtaining the bootstrap samples.
885
886      threshold : float
887        features whose linear correlations with vectors in the output space
888        are below this value are discarded.
889      '''
890
891      if not len(X) == len(Y):
892          raise ValueError("X and Y must have equal length.")
893
894      #Arguments to pass to worker processes
895      d = {'X' : X, 'Y' : Y, 'train_percentage' : train_percentage,
896           'classifier' : classifier, 'clf_args' : clf_args,
897           'score_type' : score_type, 'n_groups' : n_groups,
898           'replace' : replace, 'threshold' : threshold}
899
900      #Make job queue
901      pool = multiprocessing.Pool(processes = n_jobs)
902      jobs = [pool.apply_async(_worker, kwds = d) for _ in xrange(iterations)]
903      pool.close()  #run
904      pool.join()  #Wait for  remaining jobs
905
906      #Make sure no children died too early
907      if not all(job.successful() for job in jobs):
908          raise RuntimeError('Some jobs failed.')
909
910      return [j.get() for j in jobs]
911
912  def get_correlations(X, Y, absolute = True):
913      '''Given a list of feature vectors X and labels or values Y, returns a list
914       of correlation coefficients for each dimension of the feature vectors.'''
915      n_feats = len(X[0])
916      correlations = []
917      for i in xrange(n_feats):
918          temp = np.corrcoef([x[i] for x in X], Y)
```

```
919          correlation = temp[0,1]
920          if math.isnan(correlation):
921              correlation = 0
922          correlations.append(correlation)
923      if absolute:
924          correlations = [np.abs(c) for c in correlations]
925      return correlations
926
927
928  def make_kernel(importances, gamma = 1.0):
929      '''Returns a weighted radial basis function (WRBF) kernel which can be
930      passed to an SVM or SVR from the sklearn module.
931
932      Parameters:
933      ----------------------
934      importances : list
935        The importance of each input feature. The value of element i can mean
936        e.g. the linear correlation between feature i and target variable y.
937        None means feature will be weighted equally.
938
939      gamma : float
940        The usual gamma parameter denoting inverse width of the gaussian used.
941      '''
942      def kernel(x,y, *args, **kwargs):
943          d = len(importances) #number of features
944          impsum = sum([imp**2 for imp in importances])
945          if not impsum == 0:
946              normfactor = 1.0/np.sqrt(impsum)
947          else:
948              normfactor = 0.0
949          #Metric to compute distance between points
950          metric = dok_matrix((d,d), dtype = np.float64)
951          for i in xrange(d):
952              metric[i,i] = importances[i]*normfactor
953          #
954          result = np.zeros(shape = (len(x), len(y)))
955          for i in xrange(len(x)):
956              for j in xrange(len(y)):
957                  diff = x[i] - y[j]
958                  dist = diff.T.dot(metric*diff)
959                  result[i,j] = np.exp(-gamma*dist)
960          return result
961      return kernel
962
963  class WRBFR(svm.SVR):
964      '''Weighted radial basis function support vector regressor.'''
965      def __init__(self, importances, C = 1.0, epsilon = 0.1,
966                   gamma = 0.0):
967          kernel = make_kernel(importances = importances, gamma = gamma)
968          super(WRBFR, self).__init__(C = C, epsilon = epsilon, kernel = kernel)
969
970  class WRBFC(svm.SVC):
971      '''Weighted radial basis function support vector classifier.'''
972      def __init__(self, importances, C = 1.0, gamma = 0.0):
973          kernel = make_kernel(importances = importances, gamma = gamma)
974          super(WRBFC, self).__init__(C = C, kernel = kernel)
975
976  if __name__ == '__main__':
977
978      random.seed(42)
979      X, Y, ind_dict = read_data('../data.json', trait = 'extraversion',
```

```
980  #                                 features = ['call_iet_med', 'text_iet_med', 'social_entropy', 'call_entropy', 'travel_med', ↵
          'n_places', 'text_latency', ' call_night_activity ']
981                                   features = 'all'
982                                )
983      corrs = get_correlations(X, Y)
984
985
986  #   for i in xrange(len(corrs)):
987  #        print corrs[i], ind_dict[i]
988
989      cut = int(0.6*len(X))
990      xtrain = X[:cut]
991      ytrain = Y[:cut]
992      xtest = X[cut:]
993      ytest = Y[cut:]
994
995      inds = get_strong_indices(xtrain, ytrain, threshold = 0.08)
996      ind_dict = {i : ind_dict[inds[i]] for i in xrange(len(inds))}
997      xtrain = reduce_data(xtrain, inds)
998      xtest = reduce_data(xtest, inds)
999
1000     importances = get_correlations(xtrain, ytrain)
1001     clf = WRBFR(importances = importances, C = 26.5, epsilon=0.11, gamma = 0.23)
1002     clf_baseline = _BaselineRegressor()
1003 #    clf  = svm.SVR(C = 26.5, epsilon=0.11, gamma = 0.23)
1004     clf.fit(xtrain, ytrain)
1005     clf_baseline.fit(xtrain, ytrain)
1006     scores_baseline = []
1007     scores = []
1008     for i in xrange(len(xtest)):
1009         pred = clf.predict(xtest[i])
1010         pred_baseline = clf_baseline.predict(xtest[i])
1011         print pred
1012         scores.append((pred − ytest[i])**2)
1013         scores_baseline.append((pred_baseline − ytest[i])**2)
1014     print np.mean(scores)**0.5
1015     print np.mean(scores_baseline)**0.5
1016
1017 #   X =      [[1,2,7,0],[3,1,6,0.01],[6,8,1,0],[10,8,2,0.01]]
1018 #   Y =  [1,1,0,0]
1019 #
1020
1021 #   corrs  = get_correlations(X, Y)
1022 #   print corrs
1023 #
1024 #   C = 70
1025 #   gamma = 3.75
1026 #
1027 #   kernel = make_kernel(corrs, 0.0)
1028 #
1029 #   clf  = svm.SVC(kernel = kernel)
1030 #   clf  = svm.SVC()
1031 #   clf . fit (xtrain , ytrain)
1032 #
1033 #   hits  = 0
1034 #
1035 #   for  i  in  xrange(len(xtest)):
1036 #        if  clf .predict( xtest [ i ])  == ytest[i ]:
1037 #             hits  += 1
1038 #        #
1039 #   print 100.0*hits /len( ytest )
1040 #
```

```
1041  #     for i in xrange(len(corrs)):
1042  #         print ind_dict[i], corrs[i]
1043  #
1044  #     test = WRBFC(threshold=0.1, gamma = 0.2, importances=corrs)
1045  #     test.fit(xtrain, ytrain)
1046  #
1047  #     preds = []
1048  #     for i in xrange(len(ytest)):
1049  #         pred = test.predict(xtest[i])[0]
1050  #         print pred, ytest[i]
1051  #         preds.append(pred)
1052  #
1053  #     print len(Counter(preds).keys())
1054
1055  #     print len(X[0])
1056  ##     print i
1057  #     print [el for el in i.values() if 'init' in el]
```

## *A.1.5   Lloyd's Algorithm*

```python
1   # −∗− coding: utf−8 −∗−
2   from __future__ import division
3
4   import numpy as np
5   import matplotlib
6   matplotlib.use('Agg') #ugly hack to allow plotting from terminal
7   import matplotlib.pyplot as plt
8   import random
9   from copy import deepcopy
10
11  def _dist(p,q):
12      return sum([(p[i]−q[i])**2 for i in xrange(len(p))])
13
14  def _lloyds_single_run(X, K, max_iterations, init):
15      # Initialize  with a subset of the data
16      if init == 'sample':
17          initials = random.sample(X, K)
18      # Or initialize  with random points across the same range as data
19      elif init == 'scatter':
20          vals = zip(∗X)
21          xmin = min(vals[0])
22          xmax = max(vals[0])
23          ymin = min(vals[1])
24          ymax = max(vals[1])
25          initials = [(random.uniform(xmin, xmax),
26                       random.uniform(ymin, ymax)) for _ in xrange(K)]
27      # Or yell  RTFM at user
28      else:
29          raise ValueError('Invalid initialization mode!')
30
31      #Contruct hashmap mapping integers up to K to centroids
32      centroids = dict(enumerate(initials))
33      converged = False
34      iterations = 0
35
36      while not converged and iterations < max_iterations:
37          clusters = {i : [] for i in xrange(K)}
```

```python
38            #Make sure clusters and centroids have identical keys, or we're doomed.
39            assert set(clusters.keys()) == set(centroids.keys())
40            prev_centroids = deepcopy(centroids)
41
42            ### STEP ONE −update clusters
43            for x in X:
44                #Check distances to all centroids
45                bestind = −1
46                bestdist = float('inf')
47                for ind, centroid in centroids.iteritems():
48                    dist = _dist(x, centroid)
49                    if dist < bestdist:
50                        bestdist = dist
51                        bestind = ind
52                    #
53                clusters[bestind].append(x)
54
55            ### STEP TWO −update centroids
56            for ind, points in clusters.iteritems():
57                if not points:
58                    pass #Cluster's empty − nothing to update
59                else:
60                    centroids[ind] = np.mean(points, axis = 0)
61
62            ### We're converged when all old centroids = new centroids.
63            converged = all([_dist(prev_centroids[k],centroids[k]) == 0
64                    for k in xrange(K)])
65            iterations += 1
66            #
67    return {tuple(centroids[i]) : clusters[i] for i in xrange(K)}
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]
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.1.6   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:
```

```python
54                     c = _make_diameter(p, q)
55                 else:
56                     c = _make_circle_two_points(points[0 : i + 1], p, q)
57     return c


# Two boundary points known
def _make_circle_two_points(points, p, q):
    diameter = _make_diameter(p, q)
    if all(_is_in_circle(diameter, r) for r in points):
        return diameter

    left = None
    right = None
    for r in points:
        cross = _cross_product(p[0], p[1], q[0], q[1], r[0], r[1])
        c = _make_circumcircle(p, q, r)
        if c is None:
            continue
        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])):
            left = c
        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])):
            right = c
    return left if (right is None or (left is not None and left[2] <= right[2])) else right


def _make_circumcircle(p0, p1, p2):
    # Mathematical algorithm from Wikipedia: Circumscribed circle
    ax = p0[0];  ay = p0[1]
    bx = p1[0];  by = p1[1]
    cx = p2[0];  cy = p2[1]
    d = (ax * (by - cy) + bx * (cy - ay) + cx * (ay - by)) * 2.0
    if d == 0.0:
        return None
    x = ((ax * ax + ay * ay) * (by - cy) + (bx * bx + by * by) * (cy - ay) + (cx * cx + cy * cy) * (ay - by)) / d
    y = ((ax * ax + ay * ay) * (cx - bx) + (bx * bx + by * by) * (ax - cx) + (cx * cx + cy * cy) * (bx - ax)) / d
    return (x, y, math.hypot(x - ax, y - ay))


def _make_diameter(p0, p1):
    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)


_EPSILON = 1e-12


def _is_in_circle(c, p):
    return c is not None and math.hypot(p[0] - c[0], p[1] - c[1]) < c[2] + _EPSILON


# Returns twice the signed area of the triangle defined by (x0, y0), (x1, y1), (x2, y2)
def _cross_product(x0, y0, x1, y1, x2, y2):
    return (x1 - x0) * (y2 - y0) - (y1 - y0) * (x2 - x0)

#pts = [(0.0, 0.0), (6.0, 8.0)]
#
#test = make_circle(pts)
#
#print test
```

# A.2   Source Code for Explicit Semantic Analysis

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

## *A.2.1   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
52              count += 1
53              yield filename
54
55  make_filename = filename_generator(temp_dir)
56
57  #Format {right title : redirected title }, e.g. {because : ([cuz, cus])}
58  redirects = {}
59
60  #Minimum number of links/words required to keep an article.
61  from shared import min_links_out, min_words
62
63  #Open log file for writing and import logging function
64  logfile = open(os.path.basename(__file__)+'.log', 'w')
65  log = shared.logmaker(logfile)
66
67  class WikiHandler(SAX.ContentHandler):
68      '''ContentHandler class to process XML and deal with the WikiText.
69      It works basically like this:
70      It traverses the XML file, keeping track of the type of data being read and
71      adding any text to its input buffer. When event handlers register a page
72      end, the page content is processed, the processed content is placed in the
73      output buffer, and the input buffer is flushed.
74      Whenever a set number of articles have been processed, the output buffer is
75      written to a file. The point of this approach is to
76      limit memory consumption.'''
77
78      def __init__(self):
79          SAX.ContentHandler.__init__(self)
80          self.current_data = None
81          self.title = ''
82          self.input_buffer = []
83          self.output_buffer = {}
84          self.article_counter = 0
85          self.links = []
86          self.categories = []
87          self.redirect = None
88          self.verbose = False
89          #Harvest unique words here
90          self.words = set([])
91          #keeps track of ingoing article links. format {to : set([from])}
92          self.linkhash = {}
93
94      def flush_input_buffer(self):
95          '''Deletes info on the currently processed article.
96           This is called when a page end event is registered.'''
97          self.input_buffer = []
98          self.current_data = None
99          self.title = ''
100         self.links = []
101         self.categories = []
102         self.redirect = None
103
104     def flush_output_buffer(self):
105         '''Flushes data gathered so far to a file and resets.'''
106         self.output_buffer = {}
107         self.words = set([])
108         self.linkhash = {}
109
110     def startElement(self, tag, attrs):
111         '''Eventhandler for element start - keeps track of current datatype.'''
112         self.current_data = tag
113         #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)
```

```
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.2.2   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
```

```python
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.2.3   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
def matrix_chopper(matrix, dim):
    '''Generator to split a huge matrix into small submatrices, which can
     then be stored in individual files.
     This is handy both when constructing the matrix (building the whole
     matrix without saving to files in the process takes about 50 gigs RAM),
     and when applying it, as this allows one to load only the submatrix
     relevant to a given word.'''
    ind = 0
    counter = 0
    rows = matrix.get_shape()[0]
    while ind < rows:
        end = min(ind+dim, rows)
        #Return pair of submatrix number and the submatrix itself
        yield counter, sps.vstack([matrix.getrow(i)\
                                for i in xrange(ind, end)], format = 'csr')
        counter += 1
        ind += dim

def writeout():
    '''Saves the matrix as small submatrrices in separate files.'''
    for n, submatrix in matrix_chopper(mtx, row_chunk_size):
        filename = matrix_dir+str(n)+extensions['matrix']
        #Update submatrix if it's already partially calculated
        log("Writing out chunk %s" % n)
        try:
            with open(filename, 'r') as f:
                submatrix = submatrix + shared.mload(f)
            #
        except IOError:
            pass #File doesn't exist yet, so no need to change mtx

            #Dump the submatrix to file
        with open(filename, 'w') as f:
            shared.mdump(submatrix, f)
    return None

log("Constructing matrix.")
filelist = glob.glob(temp_dir + '*'+extensions['content'])
files_read = 0
for filename in filelist:
    with open(filename, 'r') as f:
        content = shared.load(f)

    #Loop over concepts (columns) as so we don't waste time with rare words
    for concept, entry, in content.iteritems():
        #This is the column index (concept w. index j)
        j = concept2index[concept]

        #Convert concept 'countmap' like so: {word : n}
        wordmap = Counter(entry['text'].split()).iteritems()

        #Add them all to the matrix
        for word, count in wordmap:
            #Find row index of the current word
            i = word2index[word]

            #Add the number of times word i occurs in concept j to the matrix
            mtx[i,j] = count
        #
    #Update file count
    files_read += 1
    log("Processed content file no. %s of %s - %s"
```

```python
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.2.4   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
import os
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
     def __init__(self, max_tweets=-1, verbose = False, tweets_pr_file = 10**5):
         #Set parameters
         self.max_tweets = max_tweets  #-1 for unlimited stream
         self.verbose = verbose
         self.tweets_pr_file = tweets_pr_file

         #Internal parameters to keep track of harvest status
         self.files_saved = 0
         self.harvested_tweets = []
         self.current_filenumber = 0
         self.current_date = date.today()

     def filename_maker(self):
         #Update counter and date if neccessary
         if not self.current_date == date.today():
             self.current_date = date.today()
             self.current_filenumber = 0
         else:
             pass  #Date hasn't changed. Proceed.
         filename = str(self.current_date) + "-data%s.json" % self.current_filenumber
         self.current_filenumber += 1
         return filename

 #Simple logging function
     def log(self, text):
         string = text+" at "+time.asctime()+"\n"
         if self.verbose:
             print string
         with open('tweetlog.log', 'a') as logfile:
             logfile.write(string)
         #Must return true so I can log errors without breaking the stream.
         return True

     def listen(self, search_term):
         #Make a listener
         listener = tweepy.StreamListener()
         #Override relevant methods.
         listener.on_status = self.process
         listener.on_error = lambda status_code: self.log("Error: "+status_code)
         listener.on_timeout = lambda: self.log("Timeout.")

         twitterStream = tweepy.Stream(auth, listener)
         twitterStream.filter(track=search_term)

     def mine(self, search_term, n = None):
         '''Mine a predefined number of tweets using input search word'''
         if n == None:
             n = self.max_tweets

         api = tweepy.API(auth)
         tweets = tweepy.Cursor(api.search, q=search_term).items(n)
         for tweet in tweets:
             self.process(tweet)

     def process(self, tweet):
         self.harvested_tweets.append(tweet)
         if self.verbose:
             print "Holding %s tweets." % len(self.harvested_tweets)

         #Write to file if buffer is full
         if len(self.harvested_tweets) == self.tweets_pr_file:
```

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

```python
182         for word, count in countmap:
183             try:
184                 ind = self.word2index[word]
185                 #Which file to look in
186                 file_number = int(ind/row_chunk_size)
187                 filename = matrix_dir+str(file_number)+extensions['matrix']
188
189                 #And which row to extract
190                 row_number = ind % row_chunk_size
191
192                 #Do it! Do it naw!
193                 with open(filename, 'r') as f:
194                     temp = shared.mload(f)
195                     result = result + count*temp.getrow(row_number)
196             except KeyError:
197                 pass    #No data on this word -> discard
198
199         #Done. Return row vector as a 1x#concepts CSR matrix
200         return result
201
202     def interpret_text(self, text, display_concepts = 10):
203         '''Attempts to guess the core concepts of the given text fragment'''
204         #Compute the interpretation vector for the text fragment
205         vec = self.interpretation_vector(text)
206
207         #Magic, don't touch
208         top_n = vec.data.argsort()[:len(vec.data)-1-display_concepts:-1]
209
210         #List top scoring concepts and their TD-IDF
211         concepts = [self.index2concept[vec.indices[i]] for i in top_n]
212         return concepts
213 #         scores = [vec.data[i] for i in top_n]
214 #         #Return as dict {concept : score}
215 #         return dict(zip(concepts, scores))
216
217     def interpret_file(self, filename):
218         with open(filename, 'r') as f:
219             data = self.clean(f.read())
220         return self.interpret_text(data)
221
222     def interpret_input(self):
223         text = raw_input("Enter text fragment: ")
224         topics = self.interpret_text(text)
225         print "Based on your input, the most probable topics of your text are:"
226         print topics[:self.display_concepts]
227
228     def scalar(self, v1, v2):
229         #Compute their inner product and make sure it's a scalar
230         dot = v1.dot(v2.transpose())
231         assert dot.shape == (1,1)
232
233         if dot.data:
234             scal = dot.data[0]
235         else:
236             scal = 0    #Empty sparse matrix means zero
237
238         #Normalize and return
239         sim = scal/(norm(v1.data)*norm(v2.data))
240         return sim
241
242     def cosine_similarity(self, text1, text2):
243         '''Determines cosine similarity between input texts.
```

```
244              Returns float in [0,1]'''
245
246          #Determine intepretation vectors
247          v1 = self.interpretation_vector(text1)
248          v2 = self.interpretation_vector(text2)
249
250          #Compute the normalized dot product and return
251          return self.scalar(v1, v2)
252
253
254      def cosine_distance(self, text1, text2):
255          return 1−self.cosine_similarity(text1, text2)
256
257  if __name__ == '__main__':
258      th = TweetHarvester(verbose=True, max_tweets=10)
259      th.mine('carlsberg', n=10)
260      temp = [t._json for t in th.harvested_tweets if t._json['lang'] == 'en']
261      js = temp[4]
262      with open('tweet_example.json', 'w') as f:
263          pprint(js, stream=f)
264
265  #     if len(sys.argv) > 1:
266  #         fn  = sys.argv[1]
267  #     else :
268  #         fn = 'interpret_me.txt'
269  #         with open(fn, 'r') as f:
270  #             data = f.read()
271  #         #
272  #     data = sa.clean(data)
273  #     guesses = sa. interpret_text (data)
274  #
275  #     if len(sys.argv) > 2:
276  #         output_filename = sys.argv[2]
277  #     else :
278  #         output_filename = 'guesses.txt'
279  #     with open(output_filename, 'w') as f:
280  #         for line in guesses:
281  #             f .write(line .encode('utf8'))
282  #             f .write ('\n')
```

### A.2.5   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)
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
81                  if text[2] == "x":
82                      return unichr(int(code[1:], 16))
83                  else:
84                      return unichr(int(code))
85              else:                    # named entity
86                  return unichr(name2codepoint[code])
87          except UnicodeDecodeError:
88              return text # leave as is
89
90      return re.sub("&#?(\w+);", fixup, text)
91
92  def drop_spans(matches, text):
93      """Drop from text the blocks identified in matches"""
94      matches.sort()
95      res = ''
96      start = 0
97      for s, e in  matches:
98          res += text[start:s]
99          start = e
100     res += text[start:]
101     return res
102
103 ###Compile regexps for text cleanup:
104 #Construct patterns for  elements to be discarded:
105 discard_elements = set([
106         'gallery', 'timeline', 'noinclude', 'pre',
107         'table', 'tr', 'td', 'th', 'caption',
108         'form', 'input', 'select', 'option', 'textarea',
109         'ul', 'li', 'ol', 'dl', 'dt', 'dd', 'menu', 'dir',
110         'ref', 'references', 'img', 'imagemap', 'source'
111         ])
112 discard_element_patterns = []
113 for tag in discard_elements:
114     pattern = re.compile(r'<\s*%s\b[^>]*>.*?<\s*/\s*%s>' % (tag, tag), re.DOTALL | re.IGNORECASE)
115     discard_element_patterns.append(pattern)
116
117 #Construct patterns to  recognize HTML tags
118 selfclosing_tags = set([ 'br', 'hr', 'nobr', 'ref', 'references' ])
119 selfclosing_tag_patterns = []
120 for tag in selfclosing_tags:
121     pattern = re.compile(r'<\s*%s\b[^/]*/\s*>' % tag, re.DOTALL | re.IGNORECASE)
122     selfclosing_tag_patterns.append(pattern)
123
124 #Construct patterns for  tags  to  be ignored
125 ignored_tags = set([
126         'a', 'b', 'big', 'blockquote', 'center', 'cite', 'div', 'em',
127         'font', 'h1', 'h2', 'h3', 'h4', 'hiero', 'i', 'kbd', 'nowiki',
128         'p', 'plaintext', 's', 'small', 'span', 'strike', 'strong',
129         'sub', 'sup', 'tt', 'u', 'var',
130 ])
131 ignored_tag_patterns = []
132 for tag in ignored_tags:
133     left = re.compile(r'<\s*%s\b[^>]*>' % tag, re.IGNORECASE)
134     right = re.compile(r'<\s*/\s*%s>' % tag, re.IGNORECASE)
135     ignored_tag_patterns.append((left, right))
136
137 #Construct patterns to  recognize math and code
138 placeholder_tags = {'math':'formula', 'code':'codice'}
139 placeholder_tag_patterns = []
```

```python
140  for tag, repl in placeholder_tags.items():
141      pattern = re.compile(r'<\s*%s(\s*| [^>]+?)>.*?<\s*/\s*%s\s*>' % (tag, tag), re.DOTALL | re.IGNORECASE)
142      placeholder_tag_patterns.append((pattern, repl))
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
```

```python
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.2.6   Application Examples

This section contains code used for the examples of applications of explicit semantic analysis described in section 3.3.

### Trend Discovery and Monitoring

This is the script used to generate the results seen in section 3.3.1.

```
1  # −∗− coding: utf−8 −∗−
2  from __future__ import division
3
4  import matplotlib as mpl
5  mpl.use('Agg')
6  import matplotlib.pyplot as plt
7  import sys
8  sys.path.insert(0, '../../')
9  import codecs
10 sys.stdout = codecs.getwriter('utf8')(sys.stdout)
11 sys.stderr = codecs.getwriter('utf8')(sys.stderr)
12 from collections import Counter
13 from cunning_linguistics import SemanticAnalyser as SA
14 import json
15 import time
16 import datetime
17 from shared import moerkeroed#, wine, oldhat, nude
18
19 INPUT_FILENAME = 'carlsberg_short.txt'
20 #INPUT_FILENAME = 'carlsberg_filtered_tweets.txt'
21 OUTPUT_DIR = 'animation/'
22 EXT = '.pdf'
23 SHOW = True
24 USE_ONE_IN = 1 # e.g.  if  5 only every 5th tweet will  be used. 1 for  all
25 PROCESS = False  #Whether to process data from scratch or read in old  data
26
27
28 #INPUT_FILENAME = 'carlsberg_filtered_tweets.txt'
29 MAX_TWEETS = float('inf')
30
31 TOP_N = 10  #Number of best matched concepts to include in analysis
32 IGNORE = ['Carlsberg Group']  #Concepts to ignore
33 #If  not empty, track only concepts in this  list
34 TRACK = ['Carlsberg Foundation',
35         'Carlsberg Polska',
36         'Carlsberg Srbija',
37         'Old Lions',
38         'Copenhagen',
39         'Carlsberg Laboratory',
40 #        'Raidió  Teilifís  Éireann',
```

```
41            'Kim Little']
42
43  _display = TOP_N + len(IGNORE)
44
45  _now = time.time()
46
47  def tweet2epoch(t):
48      '''Extracts the time of input tweets creation and returns as datetime.'''
49      epoch = time.mktime(time.strptime(t['created_at'],'%a %b %d %H:%M:%S +0000 %Y'))
50      return epoch
51
52  def epoch2dtindex(epoch, span = 'day'):
53      dt = datetime.datetime.fromtimestamp(epoch)
54      year, week, weekday = datetime.date.isocalendar(dt)
55      day, month = dt.day, dt.month
56      if span == 'week':
57          return "%s-%02d" % (year, week)
58      elif span == 'day':
59          return "%s-%02d-%02d" % (year, month, day)
60
61  def pad_labels(X, limit = 15):
62      '''Limits label length to <limit> characters.'''
63      for i in xrange(len(X)):
64          if len(X[i]) <= limit:
65              continue
66          else:
67              X[i] = X[i][:limit]+"..."
68
69  def process():
70      tweets = []
71      n_read = 0
72      with open(INPUT_FILENAME, 'r') as f:
73          for line in f.readlines():
74              tweet = json.loads(line)
75              n_read += 1
76              if not n_read % USE_ONE_IN == 0:
77                  continue
78              if tweet['retweeted'] or not tweet['lang'] == 'en':
79                  continue
80              else:
81                  tweets.append({'created_at' : tweet2epoch(tweet),
82                                 'text' : tweet['text']})
83                  if len(tweets) >= MAX_TWEETS:
84                      break
85                  #
86              #
87          #
88
89      # Create a semantic analyser
90      sa = SA()
91
92      data = {}   #This will map index yyyy_ww to cobined list of top concepts
93      counter = 0
94      for tweet in tweets:
95          concepts = sa.interpret_text(tweet['text'], display_concepts = _display)
96          filtered = [c for c in concepts if not c in IGNORE][:TOP_N]
97          index = epoch2dtindex(tweet['created_at'])
98          # If tracking, include only tracked concepts
99          if TRACK:
100             filtered = [c for c in filtered if c in TRACK]
101         try:
102             data[index] += filtered
```

```python
103            except KeyError:
104                data[index] = filtered
105            counter += 1
106            print "Processed %d of %d tweets..." % (counter, len(tweets))
107
108        #Change data into {index : (sorted topics, their counts)}
109        for index, _list in data.iteritems():
110            # Most often used concepts in entire period coresponding to index
111            d = dict(Counter(_list))
112            top_concepts = sorted(d, key = d.get)[::-1][:TOP_N]
113            #If not tracking specific concepts, just save top n concepts and counts
114            if not TRACK:
115                X = sorted(top_concepts)
116                Y = [d[x] for x in X]
117            #If tracking, use NaN for concepts that didn't occur in top n.
118            else:
119                X = TRACK
120                Y = [d[x] if x in top_concepts else float('nan') for x in TRACK]
121            data[index] = (X, Y)
122
123
124        return data
125
126    if __name__ == '__main__':
127        if PROCESS:
128            data = process()
129
130            with open('processed_data.json', 'w') as f:
131                json.dump(data, f, indent = 4)
132            #
133        else:
134            with open('processed_data.json', 'r') as f:
135                data = json.load(f)
136
137    #=============================================================================
138    # Plotting ...
139    #=============================================================================
140
141        #Set pylab params to make plots look similar
142        ymin = float('inf')
143        ymax = float('-inf')
144        for _,Y in data.values():
145            if not Y:
146                continue
147            thismin = min(Y)
148            thismax = max(Y)
149            if thismin < ymin:
150                ymin = thismin
151            if thismax > ymax:
152                ymax = thismax
153
154        ymin = 0  # Don't truncate
155        yair = (ymax - ymin)*0.05
156
157        xair = 0.5
158        xmin = 0
159        xmax = len(TRACK) if TRACK else TOP_N
160
161
162        for index, datum in data.iteritems():
163            #Enforce axes 'n stuff
164            fig = plt.figure(figsize = (3.25,3))
```

```
165        ax = fig.add_subplot(111)
166        plt.axis((xmin − xair, xmax + xair, ymin − yair, ymax + yair))
167        ax.set_autoscale_on(False)
168
169        X, Y = datum
170        pad_labels(X, limit = 20)
171        plt.xticks([n+0.4 for n in xrange(len(X))], X, rotation = 50,
172                   ha='right', fontsize = 8)
173        plt.bar(range(len(X)), Y, color = moerkeroed)
174        plt.gcf().subplots_adjust(bottom = 0.42, left = 0.3)
175        plt.title(index)
176
177        plt.savefig(OUTPUT_DIR+index+EXT, dpi = 600)
178
179        if SHOW:
180            plt.show()
181
182        plt.close()
183
184
185    print "Runtime (m): ", (time.time()−_now)/60
```

## Social Media Impact

This is the code used to extract and analyse the data described in section 3.3.2.

```
1   # −∗− coding: utf−8 −∗−
2   '''This script computes the semantic vector corresponding to an input reference
3   text, loads in tweets contained in files saved in the specified period
4   (change filelist to whatever files contain your tweets), and then computes
5   the cosine similarity of each tweet with the reference text and saves the
6   result.'''
7   from __future__ import division
8
9   import sys
10  sys.path.insert(0, '../../')
11  from cunning_linguistics import SemanticAnalyser
12  import datetime
13  from dateutil.parser import parse
14  import json
15  from matplotlib import pylab as plt
16  import numpy as np
17  import pytz
18  import glob
19  import time
20  import multiprocessing
21
22
23  REFERENCE = 'reference_google.txt'
24  PUBLISHED = datetime.datetime(2015, 06, 17)  #Date ref was published
25  EARLY = 16  #Number of days to include around pub date
26  LATE = 42
27  OUTPUT_FILENAME = 'deep_dreams.txt'
28  KEYWORD = 'google'
29  N_JOBS = 8
30
```

```python
31
32  #15−06−17
33  with open(REFERENCE, 'r') as f:
34      reference_text = f.read()
35
36  now = time.time()
37
38  epoch = datetime.datetime(1970,1,1, tzinfo = pytz.utc)
39  def dt2epoch(dt):
40      utc_date = dt.astimezone(pytz.utc)
41      delta = utc_date − epoch
42      return delta.total_seconds()
43
44  def percentof(small, large):
45      return str(100*small/large) + "%"
46
47  #Get canonical date string
48  def timeparse(string):
49      dt = parse(string)
50      (y,m,d) = (str(dt.year), str(dt.month).zfill(2), str(dt.day).zfill(2))
51      return "%s-%s-%s" % (y,m,d)
52
53  with open(REFERENCE, 'r') as f:
54      reference_text = f.read()
55
56  #required entries  in tweets
57  ineedthese = ['lang', 'text', 'created_at']
58
59  def worker(filename):
60      Y = []
61      X = []
62      #Make an analyser
63      sa = SemanticAnalyser()
64      reference_vector = sa.interpretation_vector(reference_text)
65      with open(filename, 'r') as f:
66          tweets = [json.loads(line) for line in f.readlines()]
67      for tweet in tweets:
68          if not all(stuff in tweet.keys() for stuff in ineedthese):
69              continue
70          if not tweet['lang'] == 'en':
71              continue
72          text = tweet['text'].lower()
73          if not KEYWORD in text:
74              continue
75          t = dt2epoch(parse(tweet['created_at']))
76          this_vector = sa.interpretation_vector(text)
77          similarity = sa.scalar(this_vector, reference_vector)
78          if np.isnan(similarity):
79              continue
80          X.append(t)
81          Y.append(similarity)
82          #
83      print "Processed file: ", filename
84      d = {'X' : X,  'Y': Y}
85      return d
86
87  if __name__ == '__main__':
88      filelist = set([])
89
90      for w in np.arange(−EARLY, LATE+1, 1):
91          delta = datetime.timedelta(days = w)
92          dt = PUBLISHED + delta
```

```
 93          prefix = dt.strftime("%Y-%m-%d")
 94          pattern = 'tweets/' + prefix + '-data*'
 95          filelist.update((glob.glob(pattern)))
 96
 97      pool = multiprocessing.Pool(processes = N_JOBS)
 98      jobs = [pool.apply_async(worker, kwds = {'filename' : fn})
 99              for fn in filelist]
100      pool.close() #run
101      pool.join() #Wait for remaining jobs
102
103      #Make sure no children died too early
104      if not all(job.successful() for job in jobs):
105          raise RuntimeError('Some jobs failed.')
106
107      X = []
108      Y = []
109      for d in [j.get() for j in jobs]:
110          X += d['X']
111          Y += d['Y']
112      inds = np.argsort(X)
113      X = [X[i] for i in inds]
114      Y = [Y[i] for i in inds]
115
116      with open(OUTPUT_FILENAME, 'w') as f:
117          json.dump(X, f)
118          f.write('\n')
119          json.dump(Y, f)
120
121  #    plt.plot(X, Y)
122  #    plt.show()
```

# 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] Stephen Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.

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

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

[6] 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.

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

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

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

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

[11] 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.

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

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

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

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

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

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

[18] AJ Smola and B Schölkopf. *Learning with kernels*. 1998.

[19] AJ Smola and B Schölkopf. A tutorial on support vector regression. *Statistics and computing*, 2004.

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

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