

Appendix A

Estimating probability distribution of good and bad events.

For much research on decision making there is not a well established context and at first sight, the problem of estimating the distribution of the probability of good or bad things happening, based on peoples experience, across a wide range of contexts and concepts, seems impossible. Ideally, but impractically, we would observe someone throughout their lifetime and for each of a large range of contexts and concepts, record the number of times good and bad things happened (dividing by the number of times they could have happened). However, fortunately, a more practical solution is provided by the recent phenomenon of the internet weblog or blog. Blogs are short descriptions of people's life experiences (good, bad and indifferent). They are also searchable. Obviously, the probability of good or bad things occurring is dependent on the context or concept in question, for example, birthdays, weddings and Christmas tend to be associated with high probabilities of good things happening, while fights, house fires, and earthquakes tend to be associated with a high probability of bad things happening.

In order to establish our prior of inference we used marginalisation. Marginal probability can be obtained by summing (or integrating, more generally) the conditional probabilities over all outcomes and in this way we can find the probability distribution of these probabilities. The contexts and concepts (henceforth concept) that we searched for were 1500 nouns that had been pre-rated using Osgood's semantic differential framework (Osgood, Suci, & Tannenbaum, 1957) and grouped under the broad headings of Behaviour (actions that one person can perform on another person), Identities (different kinds of individual), Settings (places or times where social interactions might take place) and Modifiers (emotions, traits, and statuses that might characterise people) (Francis & Heise, 2006), offering a broad selection for analysis. For each of the concepts we calculated the probability of good and bad events happening. This was carried out using automated Matlab scripts querying a blog search engine, technorati.com (together with an alternative blog search engine, blogscope.com and a different sort of data source; Reuters-21578 data set; a collection of documents that appeared on Reuters newswire in 1987 (Lewis, 1997)).

Using this method it is straightforward to find how many blogs, out of the millions indexed, contain a given concept (say, knife: C_i). In order to calculate the probability of a good or bad event occurring with a particular concept, we used a set of words that are relatively unambiguous in their goodness and badness (e.g. happy, evil); the words used are shown in Table 1 and were chosen because they are the best and

worst rated concepts in the modifier group for the evaluation dimension in the pre rated data set. The distribution of the frequencies of good and bad modifiers (taken from the British National corpus using their simple search (www.natcorp.ox.ac.uk)) was not significantly different ($t(10)=0.55$, $p=.60$).

Table 1: Good and bad modifier words used for the internet blog search.

Good modifiers	Bad modifiers
good	bad
amused	suicidal
polite	evil
relaxed	abusive
pleased	cruel
helpful	depressed
delighted	miserable
friendly	rude
generous	hurt
honest	mean
happy	unhappy

For every concept we calculated the number of blogs that contained the concept on its own and then the concept together with one or more of the good/bad modifiers: *Gi* and *Bi*. This resulted in a probability that a good or bad event will occur for each concept (*Gi/Ci* or *Bi/Ci*). In order to carry out the search, three statements were constructed for each concept as follows:

- The concept on its own e.g.

bully

- The concept in conjunction with a disjunctive list of good modifiers e.g.

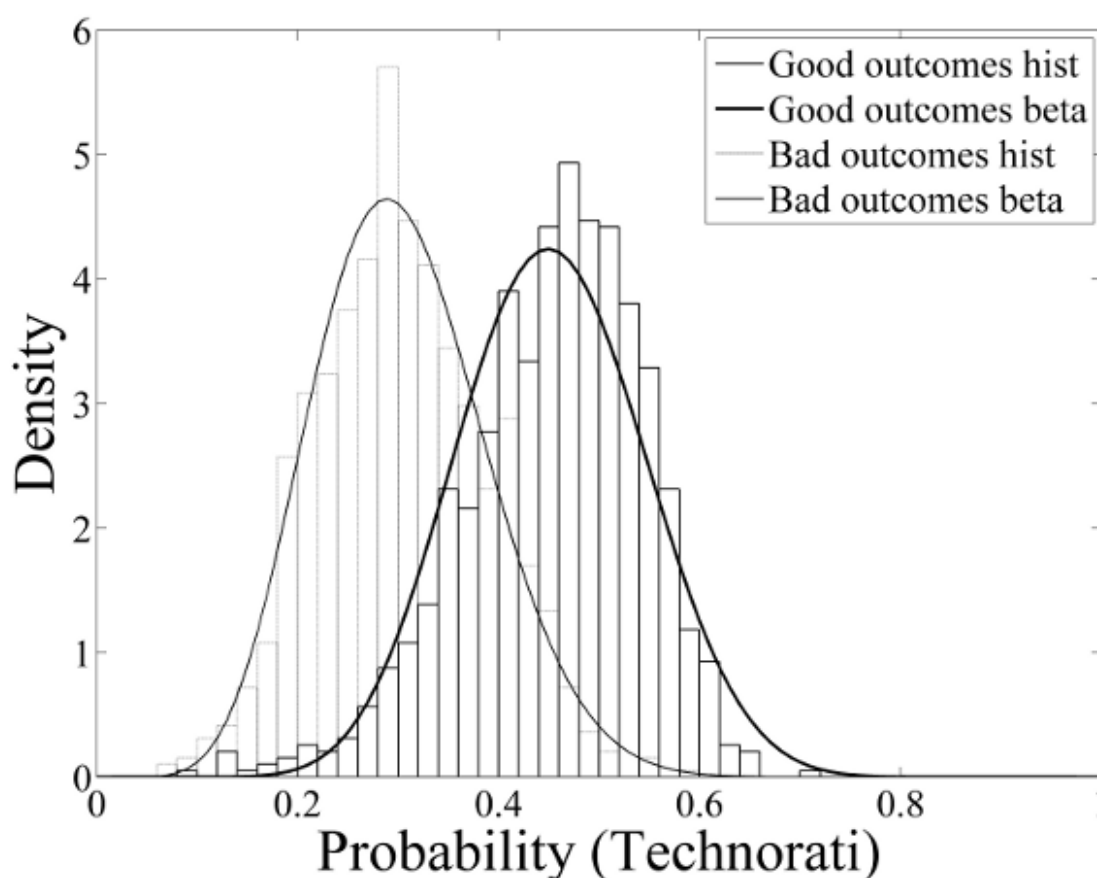
bully and (good or amused or polite or relaxed or pleased or helpful or delighted or friendly or generous or honest or happy)

- The concept in conjunction with a disjunctive list of bad modifiers e.g.

bully and (bad or suicidal or evil or abusive or cruel or depressed or

miserable or rude or hurt or mean or unhappy)

in this way, a blog containing a concept with a modifier appearing more than once was only counted once. Across all concepts, we get a distribution of probabilities and to characterise these two distributions (one for good events and one for bad) we fitted the best beta distribution using maximum likelihood. Figure A1 shows the distribution of good and bad events together with the best fitting beta distribution; the left panel shows the distributions from the technorati search engine where significant uncertainty associated with the distribution of probabilities can be seen, the average probability is less than half, good things are more common than bad, and the data is well summarised by a beta distribution. The middle panel shows distributions using the alternative blog search engine, blogscope, and the right panel shows the equivalent distribution obtained when using searches based the Reuters-21578 data set, where the same characteristics are maintained.



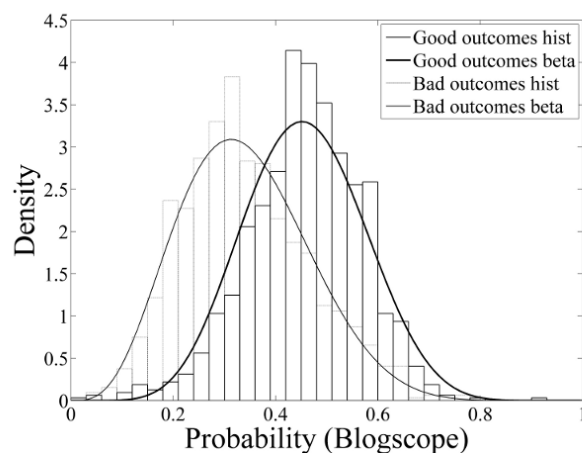


Figure 1: The probability distribution of the probability of good (thick line) and bad (thin line) outcomes across 1500 different "word" contexts. The left panel shows the data from analysing 133 million blogs using the technorati search engine, the middle panel the results from the blogscope search engine and the right hand panel shows the results from analysing 19,043 articles from Reuters-21578. The lines superimposed show the best fitting (maximum likelihood) beta distributions. The exact details across the data sets are different but there are three robust properties: the mean is less than 50%, there is a fair amount of spread, and they show "good" events are more probable than bad, even in the slightly pessimistic world described in blogs and newspapers.

Given a particular $beta(\alpha, \beta)$ prior and a Bernoulli probability statement (s, N) the posterior probability distribution is $beta(\alpha + s.N, \beta + (1-s).N)$ (MacKay, 2003). This gives two posteriors, one associated with ignorance and one with (good and bad) inference. To combine them, we used an evidence based Bayesian model averaging framework (based on MacKay, 2003). This last step is important as it provides the means to identify the extent to which a situation is like something we have encountered before and therefore the weight to apply to it. Essentially the two priors (ignorance and inference) are simply models of the world and their relative probability is determined by the compatibility of the probability statement and the probability distribution associated with each prior. In statements compatible with previous experience, the effective prior is dominated by the prior of inference and for statements incompatible with previous experience, the effective prior is dominated by the prior of ignorance. This is achieved automatically by the application of the rules of probability theory.

Please refer to the manuscript for a breakdown of the calculations used.

Implementation

Below is an implementation, in MATLAB code, of the main calculations for the probability weighting function contained in the manuscript. The full set of routines

together with a test data file is contained in the following zip file [probabilityWeightingFunction.zip](#). In order to run the probability weighting function, download the zip file and unpack it to a convenient folder and then run `weightingFunction.m`. Two plots will be displayed, one showing the good and bad curves and the other showing the mean of the good and bad curves.

```
% getPosterior - calculate a posterior distribution from the parameters
%                  below
%
% inputs
%
% A0              - alpha value for ignorance model
% B0              - beta value for ignorance model
% A1              - prior alpha value for inference model
% B1              - prior beta value for inference model
% modelPrior0     - prior for ignorance model
% modelPrior1     - prior for inference model
% nInc            - number of iterations for the data (e.g. coin tosses)
%
% outputs
%
% posteriorPdf    - matrix containing nInc distributions

function [posteriorPdf]=getPosterior(A0,B0,A1,B1,modelPrior0,modelPrior1,nInc)

    % set up a vector of probabilities between 0 and 1
    probs=0:.001:1;
    % note the number
    nProbs=length(probs);
    % initialise an array for the resulting distributions
    posteriorPdf=zeros(nInc,nProbs);

    for nToss = 0:nInc
        % for the number of trials requested
        % calculate the beta integral using gammas (could probably use
        % matlabs beta function) and the data and prior values for the
        % uniform and empirical priors
        % ...grab the probability densities for each too

        evidence0=gamma(A0+B0)/(gamma(A0)*gamma(B0)) * ...
            gamma(nToss+A0)*gamma(nInc-nToss+B0)/gamma(nInc+A0+B0);

        pdf0=betapdf(probs,nToss+A0,nInc-nToss+B0);

        evidence1=gamma(A1+B1)/(gamma(A1)*gamma(B1)) * ...
            gamma(nToss+A1)*gamma(nInc-nToss+B1)/gamma(nInc+A1+B1);

        pdf1=betapdf(probs,nToss+A1,nInc-nToss+B1);

        hypothesis1=evidence1*modelPrior1 / ...
            (evidence1*modelPrior1+evidence0*modelPrior0);
```

```
hypothesis0=evidence0*modelPrior0 / ...  
    (evidence1*modelPrior1+evidence0*modelPrior0);  
  
posteriorPdf(nToss+1,:)=(pdf0.*hypothesis0)+(pdf1.*hypothesis1);  
  
end
```

References

- Francis, C., & Heise, D. (2006). Mean Affective Ratings of 1,500 Concepts by Indiana University Undergraduates in 2002-3 [Computer file] Distributed at Affect Control Theory Website, Program Interact
<<http://www.indiana.edu/socpsy/ACT/interact/JavaInteract.html>>
- Lewis, D. D. (1997). Reuters-21578 text categorization test collection. Distribution 1.0. 1997. URL www.daviddlewis.com/resources/testcollections.
- MacKay, D. J. C. (2003). *Information theory, inference and learning algorithms*. Cambridge University Press.
- Osgood, C. E., Suci, G., & Tannenbaum, P. (1957). *The measurement of meaning*. Urbana: University of Illinois Press.
-