

# Moderating probability distributions for unrepresented uncertainty: Application to sentiment analysis via deep learning

David R. Bickel

# ▶ To cite this version:

David R. Bickel. Moderating probability distributions for unrepresented uncertainty: Application to sentiment analysis via deep learning. 2020. hal-02437780

# HAL Id: hal-02437780 https://hal.science/hal-02437780

Preprint submitted on 13 Jan 2020

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Moderating probability distributions for unrepresented uncertainty: Application to sentiment analysis via deep learning

January 13, 2020

David R. Bickel Ottawa Institute of Systems Biology Department of Biochemistry, Microbiology and Immunology Department of Mathematics and Statistics University of Ottawa 451 Smyth Road Ottawa, Ontario, K1H 8M5 +01 (613) 562-5800, ext. 8670

dbickel@uottawa.ca

#### Abstract

The probability distributions that statistical methods use to represent uncertainty fail to capture all of the uncertainty that may be relevant to decision making. A simple way to adjust probability distributions for the uncertainty not represented in their models is to average the distributions with a uniform distribution or another distribution of maximum uncertainty. A decision theoretic framework leads to averaging the distributions by taking the means of the logit transforms of the probabilities. That method does not prevent convergence to the truth, as does taking the means of the probabilities themselves. The mean-logit approach to moderating distributions is applied to natural language processing performed by a deep neural network.

**Keywords:** big data; data science; deep learning; deep neural network; discounting probability distributions; maximum entropy; unknown loss function

# 1 Introduction

Statistical models do not incorporate all uncertainty into their probability distributions. As Cox (2001) noted, statistical models only provide lower bounds on uncertainty about a quantity of interest. That is clear in frequentist inference, for each hypothesis test or confidence interval relies on assumptions that remain assumptions even when they pass statistical tests, for the absence of evidence against those assumptions is not evidence for their truth. Bayesian models also underrepresent uncertainty since they could only incorporate all the uncertainty about the models if all reasonable models, their priors, and a hyperprior over the models could be specified with certainty. Even machine learning algorithms fail to capture all uncertainty. For example, neural networks that minimize log loss in effect provide nonparametric maximum likelihood estimates of sampling distributions. They fail to fully quantify the uncertainty in classifications or other predictions, for estimated sampling distributions, as opposed to posterior predictive distributions, neglect the error in the estimates.

Many decision makers are aware that the statistical models they use do not fully represent all uncertainty. They may manage the unrepresented uncertainty by either ignoring it, by hoping it does not negatively impact decisions, or by compensating for it. Such compensation can be informal, without guidance from theory, or formal, with guidance from theory (e.g., Augustin et al., 2014, §4.7).

A simple way to implement the last approach is to adjust reported probabilities for unrepresented uncertainty by combining the reported distribution with a uniform distribution or other distribution of maximum uncertainty, as described in Section 2. A method of distribution combination suitable for that purpose is then proposed in Section 3. Using the framework of the previous two sections, Section 4 specifies the proposed method of moderating a probability distribution to the extent of unrepresented uncertainty and then records its convergence to the truth as the sample size increases. That method is illustrated in Section 5 by using it with a deep neural network to classify movie reviews into two categories: reviews expressing a positive sentiment and those expressing a negative sentiment.

# 2 Moderating a distribution by combining it with another distribution

Let  $P_0$  denote the raw distribution, a probability distribution that may require moderation in order to incorporate unrepresented uncertainty. Let  $P_1$  denote the moderating distribution, a probability distribution that would be used for decision making in the extreme case of maximal uncertainty. If uncertainty is measured as entropy and if the domain is finite, then  $P_1$  could be the probability mass function than maximizes the entropy subject to some constraints. In the absence of constraints, that  $P_1$  would be the uniform distribution.

Finally, let  $P_{\mu}$  denote a  $\mu$ -moderated probability distribution, where  $\mu \in [0, 1]$  is the degree of moderation. A method of distribution moderation is a function that transforms P into  $P_{\mu}$  such that  $P_{\mu}$  weakly converges to  $P_0$  as  $\mu \to 0$  and to  $P_1$  as  $\mu \to 1$ .

**Example 1.** In classification problems, machine learning algorithms assign a probability to each of K categories  $y(1), \ldots, y(K)$ . Let  $P_0(y(k) | (x, y), x_t)$  denote the reported probability that the category  $y_t$  of the independent variable  $x_t$  is y(k), where (x, y) is  $(x_1, y_1), \ldots, (x_n, y_n)$ , a training data set of n pairs of independent variables and their categories, and where each  $t = n + 1, n + 2, \ldots$  is the index of a data pair beyond the training data. Let  $P_1(y(k)) = 1/K$  for  $k = 1, \ldots, K$ .

The  $\mu$ -moderated probability distribution is  $P_{\mu}(\bullet | (x, y), x_t)$ , a combination of the probability mass functions  $P_0(\bullet | (x, y), x_t)$  and  $P_1(\bullet)$ . Kittler et al. (1998) reviews many methods of averaging classification probabilities over machine learning algorithms. A simple method of combination uses the weighted arithmetic mean (Cooke, 1991), which in this case is

$$P_{\mu,\text{arithmetic}}(y(k) | (x, y), x_t) = (1 - \mu) P_0(y(k) | (x, y), x_t) + \mu P_1(y(k)).$$
(1)

Cooke (1991) had considered the weighted arithmetic mean as a way to combine probability distributions elicited from experts.  $\blacktriangle$ 

**Example 2.** Let  $\mathcal{P}$  stand for either  $\{P_0, P_1\}$  or the set of mixture distributions of  $P_0$  and  $P_1$ :

$$\mathcal{P} = \{ (1-w) P_0 + w P_1 : w \in [0,1] \}.$$
(2)

With either of those cases of  $\mathcal{P}$  as a set of distributions generated by  $P_0$  and  $P_1$ , methods of transforming a set of distributions to a single distributions lead to ways to combine  $P_0$  with  $P_1$  in order to moderate  $P_0$ .

Several such methods are reviewed by Troffaes (2007) and, in a Bayesian setting, by Bickel (2015). Weighted versions of those methods could be used to transform  $\mathcal{P}$  to  $P_{\mu}$  by making the weights depend on  $\mu$  in such a way that  $P_{\mu}$  satisfies the above definition of a  $\mu$ -moderated probability distribution.

The simplest version would simply equate the weight w in equation (2) with  $\mu$ , yielding  $P_{\mu} = (1 - \mu) P_0 + \mu P_1$ . In the special case of supervised classification, that is equivalent to equation (1). An alternative version is described in Section 4.

## 3 Combining distributions into an adversarial distribution

As noted in Example 2, there are many methods for transforming a set  $\mathcal{P}$  of distributions to a single distribution. The transformation method proposed in this section will be applied to that example in Section 4.

The method is based on the following method of combining the probabilities of that y(k) is in some sense the true hypothesis among K mutually exclusive hypotheses  $y(1), \ldots, y(K)$ . For example, y(k) could represent the hypothesis that a parameter of interest is in an interval of possible parameter values or the hypothesis that the next dependent variable is equal to category y(k). The set of possible probabilities to be combined is denoted by  $\mathcal{P}(y(k))$ .

Consider future decision makers who must decide whether or not to accept y(k) on the basis of its combined probability. A scale-free, reciprocal invariant distribution of the loss functions of the decision makers leads to

$$\operatorname{logit}^{-1}\left((1-c_k)\operatorname{logit}\overline{P}(y(k))+c_k\operatorname{logit}\underline{P}(y(k))\right)$$

as the *adversarial probability*, the minimax optimal value of the probability to report to the decision makers (Bickel, 2019), where  $\underline{P} = \inf \mathcal{P}(y(k)), \ \overline{P}(y(k)) = \sup \mathcal{P}(y(k)),$ 

logit 
$$P(y) = \log \frac{P(y)}{1 - P(y)}$$
,

and  $c_k \in [0, 1]$  is the degree of caution toward accepting y(k).

Consider instead the k-independent degree  $c(P_0) \in [0,1]$  of caution toward any probability specified by  $P_0$  rather than by  $P_1$ , where  $P_0$  and  $P_1$  are any two distributions defined on the same domain. If  $\mathcal{P}(y(k)) = \{P_0(y(k)), P_1(y(k))\}$ , then the corresponding  $c(P_0)$ -adversarial probability is

$$\hat{P}_{c(P_0)}(y(k)) = \text{logit}^{-1}\left((1 - c(P_0)) \text{logit} P_0(y(k)) + c(P_1) \text{logit} P_1(y(k))\right).$$
(3)

That method of combining two probabilities suggests considering  $\tilde{P}_{c(P_0)}$  as the combination of  $P_0$  and  $P_1$ , where  $\tilde{P}_{c(P_0)}$  is the function such that  $\tilde{P}_{c(P_0)}(y(k))$  satisfies equation (3) for all  $k = 1, \ldots, K$ . Since  $\tilde{P}_{c(P_0)}$  need not satisfy  $\sum_{k=1}^{K} \tilde{P}_{c(P_0)}(y(k)) = 1$ ,  $\tilde{P}_{c(P_0)}$  can only be a guess at the combined probability distribution, which is denoted by  $P_{c(P_0)}$ .

If  $P_{c(P_0)}$  is the initial measure to be updated by a genuine probability distribution  $P_{c(P_0)}$  that satisfies intuitively reasonable regularity, locality, transitivity, and weak scaling conditions, then  $P_{c(P_0)}$  maximizes entropy in the sense that it minimizes the entropy relative to  $\tilde{P}_{c(P_0)}$ , that is,

$$P_{c(P_{0})} = \arg \inf_{P':\sum_{k=1}^{K} P'(y(k)) = 1} \sum_{k=1}^{K} P'(y(k)) \log \left(\frac{P'(y(k))}{\widetilde{P}_{c(P_{0})}(y(k))}\right),\tag{4}$$

according to Csiszár (2008, §4), which summarizes Csiszár (1991). The distribution  $P_{c(P_0)}$  satisfying those conditions is called the  $c(P_0)$ -adversarial distribution.

**Lemma 1.** The  $c(P_0)$ -adversarial distribution  $P_{c(P_0)}$  satisfies

$$P_{c(P_{0})}(y(k)) = \frac{\widetilde{P}_{c(P_{0})}(y(k))}{\sum_{k'=1}^{K} \widetilde{P}_{c(P_{0})}(y(k'))}$$
(5)

for all k = 1, ..., K.

*Proof.* By definition,  $P_{c(P_0)}$  satisfies equation (4), the solution of which is equation (5), as proved using Lagrange multipliers (Jaynes, 2003, §12.3).

## 4 The adversarial distribution as the moderated distribution

#### 4.1 The degree of caution as the degree of moderation

This section equates the degree  $\mu$  to which a raw distribution  $P_0$  is moderated (§2) with the degree  $c(P_0)$  of caution toward  $P_0$  as opposed to  $P_1$  (§3). In short,  $\mu = c(P_0)$ . That  $P_1$  would be the distribution for decision making under complete uncertainty, as in Section 2.

Then, for any  $c(P_0)$  between 0 and 1, a  $c(P_0)$ -adversarial distribution is a special case of a  $c(P_0)$ -moderated distribution. That is formally stated with  $\mu$  in place of  $c(P_0)$ :

**Theorem 1.** Consider any  $\mu \in [0, 1]$ . Every  $\mu$ -adversarial distribution  $P_{\mu}$  is a  $\mu$ -moderated distribution and satisfies

$$P_{\mu}(y(k)) = \frac{\operatorname{logit}^{-1}((1-\mu)\operatorname{logit} P_{0}(y(k)) + \mu\operatorname{logit} P_{1}(y(k)))}{\sum_{k'=1}^{K} \operatorname{logit}^{-1}((1-\mu)\operatorname{logit} P_{0}(y(k')) + \mu\operatorname{logit} P_{1}(y(k')))}$$
(6)

for all k = 1, ..., K.

Proof. Plugging  $\mu$  into the  $c(P_0)$  of equation (3) yields  $\tilde{P}_{\mu}$ . By Lemma 1, every  $\mu$ -adversarial distribution  $P_{\mu}$  satisfies equation (5) with  $\mu$  substituted for  $c(P_0)$ . The substitution of  $\tilde{P}_{\mu}$  then yields equation (6). Since  $P_{\mu}$  weakly converges to  $P_0$  as  $\mu \to 0$  and to  $P_1$  as  $\mu \to 1$ ,  $P_{\mu}$  is a  $\mu$ -moderated distribution.

**Example 3.** In the notation of Example 1, equation (6) is

$$P_{\mu}(y(k)|(x,y),x_{t}) = \frac{\log t^{-1}((1-\mu)\log t P_{0}(y(k)|(x,y),x_{t}) + \mu \log t P_{1}(y(k)))}{\sum_{k'=1}^{K} \log t^{-1}((1-\mu)\log t P_{0}(y(k')|(x,y),x_{t}) + \mu \log t P_{1}(y(k')))}$$
(7)

rather than equation (1). That expression will be illustrated in Section 5.  $\blacktriangle$ 

#### 4.2 Convergence as the sample size increases

As the sample size n increases, the  $\mu$ -adversarial distribution converges to the truth, provided that the raw distribution does so, except in degenerate cases such as  $\mu = 1$ :

**Corollary 1.** Suppose there is a category y such that  $\lim_{n\to\infty} P_0(y) = 1$  with probability 1 and  $0 < P_1(y) < 1$ . Then  $\lim_{n\to\infty} P_\mu(y) = 1$  with probability 1 for any  $\mu < 1$ , where  $P_\mu$  is a  $\mu$ -adversarial distribution.

*Proof.* Assume  $0 \le \mu < 1$ . Equation (6) holds by Theorem 1. Thus, since  $\lim_{n\to\infty} P_0(y) = 1$  with probability 1,

$$\lim_{n \to \infty} P_{\mu}(y) \propto \operatorname{logit}^{-1} \left( (1-\mu) \operatorname{logit} \lim_{n \to \infty} P_{0}(y) + \mu \operatorname{logit} P_{1}(y) \right)$$
$$= \operatorname{logit}^{-1} \left( (1-\mu) \lim_{n \to \infty} \log \frac{P_{0}(y)}{1-P_{0}(y)} + \mu \log \frac{P_{1}(y)}{1-P_{1}(y)} \right)$$
$$= \operatorname{logit}^{-1} \left( (1-\mu) \lim_{n \to \infty} \log \frac{P_{0}(y)}{1-P_{0}(y)} \right)$$
$$= \left( 1 + e^{-(1-\mu) \lim_{n \to \infty} \log \frac{P_{0}(y)}{1-P_{0}(y)}} \right)^{-1} = (1+0)^{-1} = 1$$

with probability 1. Since  $\lim_{n\to\infty} P_0(y) = 1$  with probability 1 and  $\sum_{k=1}^{K} P_0(y(k)) = 1$ , we have

 $\lim_{n\to\infty} P_0(y(k)) = 0$  with probability 1 for all k = 1, ..., K such that  $y(k) \neq y$ . Therefore, for each of those values of k, equation (6) gives

$$\begin{split} \lim_{n \to \infty} P_{\mu} \left( y \left( k \right) \right) &\propto \operatorname{logit}^{-1} \left( (1 - \mu) \operatorname{logit} \lim_{n \to \infty} P_{0} \left( y \left( k \right) \right) + \mu \operatorname{logit} P_{1} \left( y \left( k \right) \right) \right) \\ &= \operatorname{logit}^{-1} \left( (1 - \mu) \lim_{n \to \infty} \log \frac{P_{0} \left( y \left( k \right) \right)}{1 - P_{0} \left( y \left( k \right) \right)} + \mu \log \frac{P_{1} \left( y \left( k \right) \right)}{1 - P_{1} \left( y \left( k \right) \right)} \right) \\ &= \operatorname{logit}^{-1} \left( (1 - \mu) \lim_{n \to \infty} \log \frac{P_{0} \left( y \left( k \right) \right)}{1 - P_{0} \left( y \left( k \right) \right)} \right) = \operatorname{logit}^{-1} \left( - (1 - \mu) \lim_{n \to \infty} \log \frac{1 - P_{0} \left( y \left( k \right) \right)}{P_{0} \left( y \left( k \right) \right)} \right) \\ &= \left( 1 + e^{+(1 - \mu) \lim_{n \to \infty} \log \frac{1 - P_{0} \left( y \left( k \right) \right)}{P_{0} \left( y \left( k \right) \right)}} \right)^{-1} = 0 \end{split}$$

with probability 1. Those two expressions of proportionality together yield the claim since  $\sum_{k=1}^{K} P_{\mu}(y(k)) = 1$  according to equation (6).

That property is highly desirable since it means moderating the raw distribution does not interfere with its convergence as the sample size increases. Intuitively, moderation becomes less necessary as the sample becomes larger. That property is not shared by all moderated distributions; for example, it does not hold for the arithmetic mean of equation (1).

#### 4.3 Decision-theoretic moderation of distributions

For a parameter  $\theta$  of interest, consider a null hypothesis y(1) such as  $H_0: \theta = 0$  and a mutually exclusive alternative hypothesis y(2) such as  $H_1: \theta \neq 0$ . The  $\Delta$ -discounted posterior probability of y(1) is

$$\widetilde{P}_{\Delta}^{\star}(y(1)) = \left(1 + \left(\frac{P(y(1))}{P(y(2))}\right)^{-\Delta}\right)^{-1} = \left(1 + \left(\frac{P(y(1))}{1 - P(y(1))}\right)^{-\Delta}\right)^{-1},\tag{8}$$

which is derived from a decision-theoretic method of moderating posterior distributions, where  $\Delta \geq 1$  is the *degree of discounting* (Bickel, 2017, Example 1). The case of no discounting ( $\Delta = 1$ ) then results in  $\widetilde{P}^{\star}_{\Delta}(y(1)) = P(y(1))$ . Analogously, the  $\Delta$ -discounted posterior probability of y(2) is

$$\widetilde{P}_{\Delta}^{\star}(y(2)) = \left(1 + \left(\frac{P(y(2))}{P(y(1))}\right)^{-\Delta}\right)^{-1} = \left(1 + \left(\frac{P(y(2))}{1 - P(y(2))}\right)^{-\Delta}\right)^{-1}.$$
(9)

The maximum entropy argument of Section 3 leads to normalizing the discounted probabilities: for k = 1 and k = 2,

$$P^{\star}_{\Delta}(y(k)) = \frac{\widetilde{P}^{\star}_{\Delta}(y(k))}{\widetilde{P}^{\star}_{\Delta}(y(1)) + \widetilde{P}^{\star}_{\Delta}(y(2))}.$$
(10)

The probability mass function  $P^{\star}_{\Delta}$  on  $\{y(1), y(2)\}$  that satisfies equation (10) is called the  $\Delta$ discounted posterior distribution. It is a special case of an adversarial distribution: **Corollary 2.** The  $\Delta$ -discounted posterior distribution is the  $(1 - 1/\Delta)$ -adversarial distribution based on the moderating distribution given by  $P_1(y(1)) = P_1(y(2)) = 1/2$ .

*Proof.* By Theorem 1, the  $(1 - 1/\Delta)$ -adversarial distribution based on the moderating distribution given by  $P_1(y(1)) = P_1(y(2)) = 1/2$  satisfies

$$\begin{split} P_{1-\frac{1}{\Delta}}\left(y\left(k\right)\right) &= \frac{\log \mathrm{i} t^{-1}\left(\frac{1}{\Delta}\log \mathrm{i} P_{0}\left(y\left(k\right)\right) + \left(1-\frac{1}{\Delta}\right)\log \mathrm{i} \frac{1}{2}\right)}{\sum_{k'=1}^{2}\log \mathrm{i} t^{-1}\left(\frac{1}{\Delta}\log \mathrm{i} P_{0}\left(y\left(k'\right)\right) + \left(1-\frac{1}{\Delta}\right)\log \mathrm{i} \frac{1}{2}\right)}{1-\frac{1}{2}}\right) \\ &\propto \log \mathrm{i} t^{-1}\left(\frac{1}{\Delta}\log \frac{P\left(y\left(k\right)\right)}{1-P\left(y\left(k\right)\right)} + \left(1-\frac{1}{\Delta}\right)\log \frac{\frac{1}{2}}{1-\frac{1}{2}}\right)\right) \\ &= \log \mathrm{i} t^{-1}\left(\log\left(\frac{P\left(y\left(k\right)\right)}{1-P\left(y\left(k\right)\right)}\right)^{\frac{1}{\Delta}}\right) \\ &= \left(1+e^{-\log\left(\frac{P\left(y\left(k\right)\right)}{1-P\left(y\left(k\right)\right)}\right)^{-1}} = \left(1+e^{\log\left(\frac{P\left(y\left(k\right)\right)}{1-P\left(y\left(k\right)\right)}\right)^{-1}}\right)^{-1} \\ &= \left(1+\left(\frac{P\left(y\left(k\right)\right)}{1-P\left(y\left(k\right)\right)}\right)^{-\frac{1}{\Delta}}\right)^{-1} = \widetilde{P}_{\Delta}^{\star}\left(y\left(k\right)\right) \end{split}$$

for k = 1, 2, with the last step following from equations (8) and (9). Since  $P_{1-\frac{1}{\Delta}}(y(1)) + P_{1-\frac{1}{\Delta}}(y(2)) = 1$  is required by equation (6), it follows that

$$P_{1-\frac{1}{\Delta}}\left(y\left(k\right)\right) = \frac{\widetilde{P}_{\Delta}^{\star}\left(y\left(k\right)\right)}{\widetilde{P}_{\Delta}^{\star}\left(y\left(1\right)\right) + \widetilde{P}_{\Delta}^{\star}\left(y\left(2\right)\right)} = P_{\Delta}^{\star}\left(y\left(k\right)\right)$$

for k = 1, 2, with the last step resulting from equation (10).

That result says that in the cases considered, the degrees of moderation and discounting are related by  $\mu = 1 - 1/\Delta$  and  $\Delta = 1/(1 - \mu)$ .

### 5 Sentiment analysis improved by distribution moderation

#### 5.1 Natural language processing by deep learning

Sentiment analysis is an approach to computational linguistics that automatically extracts opinions from natural language communications between humans. Sentiment analysis has been applied not only to the social sciences and business analytics but also to healthcare (Satapathy et al., 2018). For example, it can be used to infer the quality of care from patients' descriptions of them on social media (Greaves et al., 2013) and can improve the determination of patients' moods from their participation in online networks (Beaunoyer et al., 2017).



Figure 1: Two string-sentiment pairs randomly selected from the training set.

The hierarchical structure of natural language suggests the use of deep learning via neural networks that have multiple hidden layers (Satapathy et al., 2018, §1.6). Such deep neural networks have outperformed other forms of machine learning for problems in natural language processing (Hasan and Farri, 2019). For instance, Lee et al. (2017) used a deep neural network to detect adverse drug events from social media such as Twitter feeds.

#### 5.2 A sentiment analysis based on a deep neural network

A simple form of sentiment analysis classifies text from a product review as expressing a positive or negative sentiment (Satapathy et al., 2018, §1.5). That problem has been addressed by applying deep neural networks to movie reviews (Radford et al., 2017). The description of that research in Wolfram Research, Inc. (2019b) is summarized in the rest of this subsection.

The sentiment of each movie review is either negative or positive. The training set and nontraining set respectively consist of 7462 and 3200 string-sentiment pairs, with each pair consisting of a string of text from a movie review and the sentiment corresponding to the string (Figure 1).

The value of a certain output state of a 27-layer neural network trained on that training set is called a *sentiment score* since it quantifies the sentiment of each string. Replacing the strings in the training set and in the non-training set with their sentiment scores results in a *sentiment training set* and a *sentiment non-training set*, as seen in Figure 2. A naïve Bayes classifier trained on the sentiment training set is surprisingly accurate according to the sentiment non-training set.

#### 5.3 Application of adversarial distributions

The sentiment training set and the sentiment non-training set of Section 5.2 illustrate the method proposed in Section 4. In the notation of Examples (1) and 3, the sentiment scores are the independent variables and the possible sentiments are y(1) = negative and y(2) = positive.





Figure 2: The two score-sentiment pairs from the sentiment training set that correspond to the two pairs displayed in Figure 1.

The pairs in the sentiment non-training set were assigned randomly to a validation set and a test set of equal size. Thus, the sentiment training set, the validation set, and the test set consist of 7462, 1600, and 1600 score-sentiment pairs, respectively. To assess the effect of the sample size on performance, the sentiment training set and validation set were randomly permuted and then reduced to their first 7462 $\phi$  and 1600 $\phi$  pairs for each  $\phi \in \{0.01, 0.04, 0.16, 1\}$ . After the random permutations, the score-sentiment pair of the *i*th movie review is  $(x_i, y_i)$ , where  $x_i$  is the sentiment score from the review of a movie and  $y_i \in \{\text{negative, positive}\}$  is the sentiment of the movie for  $i = 1, 2, \dots, 7462 + 3200$ .

Let *n* denote the number of text-sentiment pairs in (x, y), the data set actually used in training a classifier at a value of  $\phi$ . Depending on the classifier, either (x, y) is the sentiment training set of  $n = 7462\phi$  text-sentiment pairs, in which case the validation set of  $1600\phi$  pairs could be used to optimize the degree  $\mu$  of moderation, or (x, y) is the union of the sentiment training set and the validation set, in which case no data are available for optimizing  $\mu$ .  $\mathcal{T} = \{7462 + 1601, \ldots, 7462 + 3200\}$  is the set of indices of the test set  $((x_{7462+1601}, y_{7462+1601}), \ldots, (x_{7462+3200}, y_{7462+3200}))$ , and  $|\mathcal{T}| = 1600$  is the size of the test set, which is not affected by the value of  $\phi$ .

Each classifier considered below is identified by a value of the variable written as clssfr. The test-set mean log-loss and Brier-loss of each classifier are

$$\begin{aligned} \widehat{\ell}_{\log}\left(\bullet, \text{clssfr}\right) &= \frac{1}{|\mathcal{T}|} \sum_{t \in \mathcal{T}} -\log P_{\text{clssfr}}\left(y_t \mid (x, y), x_t\right);\\ \widehat{\ell}_{\text{Brier}}\left(\bullet, \text{clssfr}\right) &= \frac{1}{|\mathcal{T}|} \sum_{t \in \mathcal{T}} \sum_{k=1}^{2} \left(P_{\text{clssfr}}\left(y\left(k\right) \mid (x, y), x_t\right) - \chi_t\left(k\right)\right)^2\\ &= \frac{1}{|\mathcal{T}|} \sum_{t \in \mathcal{T}} \left(P_{\text{clssfr}}\left(\text{negative} \mid (x, y), x_t\right) - \chi_t\left(1\right)\right)^2 + \left(P_{\text{clssfr}}\left(\text{positive} \mid (x, y), x_t\right) - \chi_t\left(2\right)\right)^2 \end{aligned}$$

where  $\chi_t(k) = 1$  if  $y(k) = y_t$  and  $\chi_t(k) = 0$  if  $y(k) \neq y_t$ . Then  $\chi_t(1)$  indicates whether

 $y_t$  = negative, and  $\chi_t(2)$  whether  $y_t$  = positive. Both loss functions penalize reporting high probabilities of incorrect sentiments.

The mean losses of these classifiers are compared:

- 1. The data-independent classifier yielding the moderating distribution  $P_1$  that assigns 50% probability to each sentiment:  $P_1$  (negative) =  $P_1$  (positive) = 1/2. That uniform classifier is abbreviated as clssfr = "U".
- 2. Four classifiers using the Classify function in Wolfram Research, Inc. (2019a):
  - (a) The default logistic regression yields  $P_0^{\text{logistic}}(\bullet | (x, y), x_t)$  as the raw distribution  $P_0(\bullet | (x, y), x_t)$ . The classifier with the training set as (x, y), the data set actually used for training, is denoted by clssfr = "tL", whereas the classifier instead trained on the union of the sentiment training set and validation set as (x, y) is denoted by clssfr = "tvL".
  - (b) The default support vector machine yields  $P_0^{\text{SVM}}(\bullet|(x, y), x_t)$  as the raw distribution  $P_0(\bullet|(x, y), x_t)$ . There are two SVM classifiers, named according to the convention used for logistic: clssfr = "tS" and clssfr = "tvS".
- 3. The eight combined classifiers yielding these versions of the  $\mu$ -moderated distributions  $P_{\mu,\text{arithmetic}}^{\text{logistic}}(\bullet | (x, y), x_t)$ ,  $P_{\mu,\text{arithmetic}}^{\text{SVM}}(\bullet | (x, y), x_t)$ , and  $P_{\mu}^{\text{SVM}}(\bullet | (x, y), x_t)$ , formed by combining either  $P_0^{\text{logistic}}(\bullet | (x, y), x_t)$  or  $P_0^{\text{SVM}}(\bullet | (x, y), x_t)$  with  $P_1$ , either according to equation (1) for  $P_{\mu,\text{arithmetic}}^{\text{logistic}}$  and  $P_{\mu,\text{arithmetic}}^{\text{SVM}}(\bullet | (x, y), x_t)$  with these degrees of moderation:
  - (a) In the simpler case,  $\mu = 1/2$  with training logistic or SVM on the union of the sentiment training set and the validation set. In that way, "U" is combined with "tvL" or "tvS". The resulting classifier's name begins with "1/2" and ends with "L" or "S".
  - (b) Alternatively, μ, denoted in this case by μ̂, is fit to minimize the mean log or Brier loss of the validation set after training logistic or SVM on the training data alone. In other words, "U" is combined with "tL" or "tS". The resulting classifier's name begins with "\*" and ends with "L" or "S".

Each of those 13 classifiers is specified by a value for clssfr according to the above symbols in quotation marks, as follows:

clssfr 
$$\in \{ "U", "tL", "tvL", "tS", "tvS", "*oL", "1/2oL", "*oS", "1/2oS", "1/2aL", "*aL", "*aS", "1/2aS" \}.$$
  
(11)

The test-set mean losses of the best-performing classifiers are displayed in Figure 3. When 1% of the sentiment training set and 1% of the validation set is used ( $\phi = 1\%$ ), the  $\hat{\mu}$ -moderated SVM distributions ("\*oS", "\*aS") perform better than the SVM distributions ("tS", "tvS"), and the  $\hat{\mu}$ -moderated logistic distributions ("\*oL", "\*aL") perform better than the logistic distributions ("tL", "tvL"). When  $\phi \in \{4\%, 16\%, 100\%\}$ , the same pattern holds for SVM, but the moderation of probabilities does not show a clear advantage for logistic. It thus appears that the original probabilities reported by logistic require no moderation since they adequately reflect the uncertainty about the sentiment. By contrast, the original probabilities reported by SVM, inadequately reflecting the uncertainty about the sentiment, are improved by moderation.

Figure 3 does not indicate a clear advantage of the adversarial distribution ("\*oL", "<sup>1</sup><sub>2</sub>oL", "\*oS", "<sup>1</sup><sub>2</sub>oS") over the arithmetic mean ("<sup>1</sup><sub>2</sub>aL", "\*aL", "\*aS", "<sup>1</sup><sub>2</sub>aS"). That suggests that the sentiment training set is too small for Corollary 1 to be relevant.

# Acknowledgments

This research was partially supported by the Natural Sciences and Engineering Research Council of Canada (RGPIN/356018-2009).

# References

- Augustin, T., Coolen, F., de Cooman, G., Troffaes, M. (Eds.), 2014. Introduction to Imprecise Probabilities. Wiley Series in Probability and Statistics. Wiley.
- Beaunoyer, E., Arsenault, M., Lomanowska, A. M., Guitton, M. J., 2017. Understanding online health information: Evaluation, tools, and strategies. Patient Education and Counseling 100 (2), 183 – 189.
- Bickel, D. R., 2015. Inference after checking multiple Bayesian models for data conflict and applications to mitigating the influence of rejected priors. International Journal of Approximate Reasoning 66, 53–72.
- Bickel, D. R., 2017. Departing from Bayesian inference toward minimaxity to the extent that the posterior distribution is unreliable, working paper, HAL-01673783. URL https://hal.archives-ouvertes.fr/hal-01673783



Figure 3: Test-set mean losses  $\hat{\ell}_{\log}(\bullet, \text{clssfr})$  and  $\hat{\ell}_{\text{Brier}}(\bullet, \text{clssfr})$  for sentiment classification, where each classifier is represented by a prefix and then a suffix, as follows. The suffixes indicate either the logistic regression classifier ("L") or the support vector machine classifier ("S"). The prefixes refer to training the classifier on the sentiment training set alone ("t"), training the classifier on the sentiment training set and the validation set ("tv"), or averaging the classifier with the uniform distribution using either 50% moderation and the validation set for training ("1/2") or a degree of moderation fitted to the validation set ("\*") and using either the arithmetic mean ("a") or the odds-based mean ("o"). (Each displayed classifier appears as a value of clssfr in expression (11).) Each of the first four plots corresponds to a different fraction  $\phi$  of the training and validation sets, and the last two plots zoom in on the very best performers when needed to distinguish them from 12

- Bickel, D. R., 2019. Reporting Bayes factors or probabilities to decision makers of unknown loss functions. Communications in Statistics - Theory and Methods 48, 2163–2174.
- Cooke, R. M., 1991. Experts in Uncertainty: Opinion and Subjective Probability in Science. Oxford University Press.
- Cox, D. R., 2001. Comment on 'Statistical modeling: The two cultures'. Statistical Science 16 (3), 216–218.
- Csiszár, I., 1991. Why least squares and maximum entropy? An axiomatic approach to inference for linear inverse problems. Ann. Stat. 19, 2032–2066.
- Csiszár, I., 2008. Axiomatic characterizations of information measures. Entropy 10, 261–273.
- Greaves, F., Ramirez-Cano, D., Millett, C., Darzi, A., Donaldson, L., 2013. Harnessing the cloud of patient experience: using social media to detect poor quality healthcare. BMJ Quality & Safety 22 (3), 251–255.
- Hasan, S. A., Farri, O., 2019. Clinical Natural Language Processing with Deep Learning. Springer International Publishing, Cham, pp. 147–171.
- Jaynes, E., 2003. Probability Theory: The Logic of Science. Cambridge University Press, Cambridge.
- Kittler, J., Hatef, M., Duin, R. P. W., Matas, J., March 1998. On combining classifiers. IEEE Transactions on Pattern Analysis and Machine Intelligence 20 (3), 226–239.
- Lee, K., Qadir, A., Hasan, S. A., Datla, V., Prakash, A., Liu, J., Farri, O., 2017. Adverse drug event detection in tweets with semi-supervised convolutional neural networks. In: Proceedings of the 26th International Conference on World Wide Web. International World Wide Web Conferences Steering Committee, Republic and Canton of Geneva, CHE, pp. 705–714. URL https://doi.org/10.1145/3038912.3052671
- Radford, A., Józefowicz, R., Sutskever, I., 2017. Learning to generate reviews and discovering sentiment. CoRR abs/1704.01444. URL http://arxiv.org/abs/1704.01444
- Satapathy, R., Cambria, E., Hussain, A., 2018. Sentiment Analysis in the Bio-Medical Domain: Techniques, Tools, and Applications. Socio-Affective Computing. Springer International Publish-

ing, New York.

 ${\rm URL}\ {\tt https://books.google.ca/books?id=gbZIDwAAQBAJ}$ 

Troffaes, M. C. M., 2007. Decision making under uncertainty using imprecise probabilities. International Journal of Approximate Reasoning 45 (1), 17–29.

Wolfram Research, Inc., 2019a. Mathematica, version 12.0.0.0. URL https://www.wolfram.com

Wolfram Research, Inc., 2019b. Sentiment Language Model Trained on Amazon Product Review Data.

URL http://bit.ly/30em3aN