

# Active Learning

Laura, Peter, Simon

## General introduction

Notation:

Data  $\mathcal{D}$ . New data  $\mathcal{D}^*$ . Parameters  $\theta$ . Input  $\mathbf{x}$ . Output  $\mathbf{y}$ . Model  $\mathcal{M}$ .

**NB: We might not mention the  $\mathcal{M}$  explicitly in all the probabilities.**

We are interested in the posterior probability of the model parameters given the data

$$P(\theta|\mathcal{D}, \mathcal{M}) = \frac{P(\mathcal{D}|\theta, \mathcal{M})\pi(\theta, \mathcal{M})}{P(\mathcal{D}, \mathcal{M})} \quad (1)$$

The posterior reflects the parameters for which the model best describes the data and the underlying uncertainty. It is typically desirable to have little uncertainty in the posterior. The amount of uncertainty can be measure by the entropy

$$H[P(\theta|\mathcal{D})] = -\mathbb{E}_{P(\theta|\mathcal{D})} \log P(\theta|\mathcal{D}) \quad (2)$$

In the psychophysics experiment that we conduct, the data is comprised of stimulus-response pairs. We can choose the stimulus  $\mathbf{x}$  and observe the response  $\mathbf{y}$ . Say we have collected some stimulus-response pairs already (represented with  $\mathcal{D}$ ). The goal is to choose the next stimulus  $\mathbf{x}$  such that the uncertainty in the posterior is

decreased. In other words, we would like to choose  $\mathbf{x}$  such that the corresponding decrease in entropy

$$H[p(\theta|\mathcal{D})] - H[p(\theta|\mathcal{D}, \mathbf{x}, \mathbf{y})] \quad (3)$$

is maximal. But we don't know the answer  $\mathbf{y}$  that we are going to get. We can only maximize the decrease in expected posterior entropy:

$$\mathcal{U}(\mathbf{x}) = H[p(\theta|\mathcal{D})] - \mathbb{E}_{P(\mathbf{y}|\mathbf{x}, \mathcal{D})} H[p(\theta|\mathcal{D}, \mathbf{x}, \mathbf{y})], \quad (4)$$

which corresponds to minimizing the second term. This is called posterior entropy minimization. We can get an alternative formulation by noting that

$$\mathcal{U}(\mathbf{x}) = I[\theta, \mathbf{y}|\mathcal{D}, \mathbf{x}] \quad (5)$$

$$= H[P(\mathbf{y}|\mathbf{x}, \mathcal{D})] - \mathbb{E}_{P(\theta|\mathcal{D})} H[P(\mathbf{y}|\mathbf{x}, \theta)], \quad (6)$$

8 where  $I$  is the mutual information, which is symmetric in its arguments. Writing  
9 it this way allows for a different interpretation of the utility function  $\mathcal{U}(\mathbf{x})$ . Now  
10  $H[P(\mathbf{y}|\mathbf{x}, \mathcal{D})]$  should be large, which makes sense, because we should choose an  
11 input  $\mathbf{x}$  for which we don't know yet what the output  $\mathbf{y}$  will be. Furthermore  
12  $\mathbb{E}_{P(\theta|\mathcal{D})} H[P(\mathbf{y}|\mathbf{x}, \theta)]$  should be small, because we don't want to choose an input  $\mathbf{x}$   
13 for which the output  $\mathbf{y}$  is very uncertain. **NB: Copied from Houlby thesis:** In other  
14 words, we seek the input  $\mathbf{x}$  for which the parameters under the posterior make con-  
15 fident predictions (term 2), but these predictions are highly diverse. That is, the  
16 parameters disagree about the output  $\mathbf{y}$ , hence this formulation is named Bayesian  
17 Active Learning by Disagreement (BALD).

The *predictive distribution* is the probability  $P(\mathcal{D}^*|\mathcal{D},\mathcal{M})$  of observing a new data point  $\mathcal{D}^*$  given the old data  $\mathcal{D}$  and a model  $\mathcal{M}$ .

$$\begin{aligned} P(\mathcal{D}^*|\mathcal{D},\mathcal{M}) &= \int d\theta P(\mathcal{D}^*,\theta|\mathcal{D},\mathcal{M}) \\ &= \int d\theta P(\mathcal{D}^*|\mathcal{D},\mathcal{M},\theta)P(\theta|\mathcal{D},\mathcal{M}) \\ &= \int d\theta P(\mathcal{D}^*|\mathcal{M},\theta)P(\theta|\mathcal{D},\mathcal{M}), \end{aligned}$$

where in the last step we used  $P(\mathcal{D}^*|\mathcal{D},\mathcal{M},\theta) = P(\mathcal{D}^*|\mathcal{M},\theta)$ , because the new data should depend only on the model and the parameters and not on the collected data. That is, we assume that the model captures all the structure in the data. This assumption is typical for Bayesian inference.

## 1 Experiment

We show two gratings to participants. The frequency of the grating is fixed. The orientation of both gratings is the same but varied in each trial to avoid afterimages. One grating is always of the same contrast level, but the side is chosen at random. We vary the contrast of the other grating. We characterize the difference in contrast between the grating that is shown on the left and the grating that is shown on the right with  $x$ . For negative  $x$  the grating on the left is of higher contrast, for positive  $x$  the stimulus on the right is of higher contrast. If we denote the fixed baseline contrast with  $x_b$ , then the value of  $x$  is in the range  $[-(1-x_b), (1-x_b)]$ .

We chose the presented  $x$  according to different strategies and record the answers (left  $L$ , or right  $R$ ) of the participants when they decide on which side the grating with higher contrast is shown.

## 34 1.1 Choosing the presented stimulus

We assume a prior  $\pi(\theta)$  where  $\theta = \{w_0, w_1, \lambda\}$  is a set of parameters that describe our model. As the model  $\mathcal{M}$  we use a sigmoid function

$$\sigma(\theta, x) = \lambda/2 + \frac{1 - \lambda/2}{1 + \exp[-w_1(x - w_0)]}, \quad (7)$$

where  $\lambda$  is the lapse rate. The lapse rate accounts for wrong answers that are not because the task was too difficult, but because the participant hit mistakenly hits the wrong button. NB: We drop the model  $\mathcal{M}$  in all the subsequent probabilities. Whenever we use  $\theta$  we mean the parameters together with the sigmoid model. We collect data  $N$  data points by presenting a stimulus  $x \in [-1.0, 1.0]$  and observing a binary response  $y \in \{0, 1\}$ :

$$\mathcal{D} = \{(x_1, y_1), \dots, (x_N, y_N)\} \equiv (X^N, Y^N) \quad (8)$$

NB: We drop the  $N$ , if it is not needed to dissociate the steps. The likelihood of the parameters  $\theta$  given the data  $D$  is given by

$$P(Y^N | \theta, X^N) = \prod_{i=1}^N P(y_i | \theta, x_i) \quad (9)$$

$$= \prod_{i=1}^N \sigma(\theta, x)^{y_i} (1 - \sigma(\theta, x))^{1-y_i} \quad (10)$$

NB: The stimuli  $x$  are not considered part of the data, because we have control over it. The posterior probability of the parameters  $\mathbf{w}$  is

$$P(\theta | X, Y) = \frac{P(Y | \theta, X) \pi(\theta)}{P(Y | X)} \quad (11)$$

The denominator, i.e. the marginal likelihood, is computed by taking the integral over all hypotheses:

$$P(Y|X) = \int P(Y|\theta', X) \pi(\theta') d\theta' \quad (12)$$

The goal is to get a posterior  $P(\mathbf{w}|X, Y)$  that is of low uncertainty. We use entropy as a measure of the current uncertainty of our estimation of  $\theta$ . By *current* we mean that we use the data  $\mathcal{D}$  we have discovered in the  $N$  steps until now. The new data points are labeled  $x, y$ .

$$H[P(\theta|\mathcal{D})] = - \int P(\theta'|\mathcal{D}) \log[P(\theta'|\mathcal{D})] d\theta'. \quad (13)$$

In principle we would now like to choose our next stimulus  $x$  such that it minimizes the resulting entropy  $H(\mathbf{w}|X_N, Y_N, x, y)$ , but we do not know what  $y$  is going to be. So we want to find the  $x$  that minimizes the mean:

$$H[P(\mathbf{w}|\mathcal{D}, x, y = 0)] P(y = 0|\mathcal{D}, x) \quad (14)$$

$$+ H[P(\mathbf{w}|\mathcal{D}, x, y = 1)] P(y = 1|\mathcal{D}, x). \quad (15)$$

Here  $P(y = 0/1|\mathcal{D}, x)$  is called the *predictive distribution*. Determining them again requires an integral over the hypotheses:

$$P(y = 0/1|\mathcal{D}, x) = \int P(y = 0/1|\theta, x) P(\theta|\mathcal{D}) d\theta \quad (16)$$

Instead we should use BALD learning and find the  $x$  that maximizes:

$$\mathcal{U}(x) = H[P(y|\mathcal{D}, x)] - \mathbb{E}_{P(\theta|\mathcal{D})} H[P(y|\theta, x)]. \quad (17)$$

## 36 1.2 Humans

37 Here the contrast difference  $x$  is chosen by a human. They can select every possible  
38 value for  $x$ . To help them they

## 39 2 Approximations

40 We need to determine  $K(x)$  for all  $x$  values that we consider as worthwhile new  
41 stimuli. This can be many values and doing the involved integrals over posteriors is  
42 costly. There are several ways to deal with this problem.

### 43 2.1 Restricting the tested stimuli

44 We need to discretize the  $x$  anyways. We choose values . . . .

### 45 2.2 Approximating the posterior

46 We often determine the mean of a function over the posterior distribution. This is  
47 computationally expensive, in particular for large parameter spaces. If we take sam-  
48 ples of the posterior and approximate the integrals by smaller sums over the samples,  
49 we can save computation time. To get good samples from the posterior distribution  
50 we use the Metropolis-Hastings algorithm as an implementation of Markov Chain  
51 Monte Carlo integration. As a proposal distribution  $Q(x; x')$  we choose a multi vari-  
52 ate normal distribution which determines the random walk that samples from the

53 posterior.

## 54 2.3 Focused active learning

55 Instead of maximizing the utility function  $\mathcal{U}(x)$  over  $x$  we can maximize it only with  
56 respect to a subset of parameters. In the scenario of the sigmoid we might be  
57 interested in choosing the  $x$  where the entropy is maximized

## 58 3 Paradox of ladder stimulus presentation

59 Assume a stimulus range from  $-a$  to  $a$ . Imagine we present stimuli in the following  
60 way:

- 61 • We present the first at the center ( $x_1 = 0$ ) and record the answer  $y_1$
- For the  $n$ -th stimulus ( $n > 1$ ) we take

$$x_n = x_{n-1} + (-1)^{y_{n-1}} \frac{a}{2^{n-1}} \quad (18)$$

62 For  $N \rightarrow \infty$ , this leads to a distribution of data points that seems to be best fitted  
63 by a step function. When we collect data this way, we know that it will always lead  
64 to a distribution that looks this way and we know that we should not conclude that  
65 the underlying psychometric function should have a steep step. If, however, someone  
66 would present us this data and claim that it was obtained without using the ladder  
67 algorithm above, we might be tempted to conclude that the underlying psychometric  
68 curve is very steep. This is paradoxical, because the likelihood principle claims that  
69 the data is all we need to determine our parameter. It should not matter how the  
70 data was obtained.

The solution of this paradox is that we should not conclude a step function in neither of the two cases. Instead the posterior suggests that all sigmoids that cross  $y^* = \frac{\# \text{ of response 1 around } x_\infty}{\# \text{ of response 0 around } x_\infty}$  at the stimulus value  $x_\infty$  to which the ladder method converges are equally probable. NB: 'around' is not properly defined here, but it can't be  $x_\infty$  precisely, because it is never reached. It should be some  $\varepsilon$ -ball. With our parameterization this is the set of sigmoids for which

$$w_0 = -x_\infty w_1 + c, \tag{19}$$

71 where  $c > 0$  if  $y^* > 1/2$ ,  $c = 0$  if  $y^* = 1/2$  and  $c < 0$  if  $y^* < 1/2$ . So for  $N \rightarrow \infty$  the  
72 posterior  $P(w_0, w_1 | \mathcal{D})$  turns into a line that is characterized by Eq. (19).  
73 Marginalizing over the threshold value  $w_0$  leads to a flat posterior for the slope  
74  $P(w_1 | \mathcal{D})$ . This shows that all slopes are equally likely and Bayesian inference does  
75 not conclude a step function.