

Active Learning

Laura, Peter, Simon

General introduction

Notation:

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

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.

We want to maximize the expected information gain of the input \mathbf{x} :

$$\mathcal{U}(\mathbf{x}) = H[p(\theta|\mathcal{D})] - \mathbb{E}_{P(\mathbf{y}|\mathbf{x}, \mathcal{D})} H[P(\theta|\mathcal{D}, \mathbf{x}, \mathbf{y})], \quad (1)$$

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

$$\mathcal{U}(\mathbf{x}) = I[\theta, \mathbf{y} | \mathcal{D}, \mathbf{x}] \quad (2)$$

$$= H[P(\mathbf{y} | \mathbf{x}, \mathcal{D})] - \mathbb{E}_{P(\theta | \mathcal{D})} H[P(\mathbf{y} | \mathbf{x}, \theta)], \quad (3)$$

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

NB: Maybe start even earlier, with the most important points of the Houlsby introduction.

We assume a prior $\pi(w)$ where \mathbf{w} is a vector of parameters that describe our model. As the model we use a sigmoid function

$$\sigma(w_0, w_1, x) = \frac{1}{1 + \exp[-(w_0 - w_1 x)]} \quad (4)$$

NB: This might not be the best way to write the sigmoid. We collect data N data

points by presenting a stimulus $x \in \mathbb{R}$ and observing a binary response $y \in \{0, 1\}$:

$$D = \{(x_1, y_1), \dots, (x_N, y_N)\} \equiv (X^N, Y^N) \quad (5)$$

NB: We drop the N , if it is not needed to dissociated the steps. The likelihood of the parameters \mathbf{w} given the data D is given by

$$P(Y^N | \mathbf{w}, X^N) = \prod_{i=1}^N P(y_i | \mathbf{w}, x_i) \quad (6)$$

$$= \prod_{i=1}^N \sigma(w_0 + w_1 x)^{y_i} (1 - \sigma(w_0 + w_1 x))^{1-y_i} \quad (7)$$

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

$$P(\mathbf{w} | X, Y) = \frac{P(Y | \mathbf{w}, X) \pi(\mathbf{w})}{P(Y | X)} \quad (8)$$

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

$$\int P(Y | \mathbf{w}', X) \pi(\mathbf{w}') d\mathbf{w}' = \iint P(Y | w_0, w_1, X) \pi(w_0, w_1) dw'_0 dw'_1 \quad (9)$$

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 \mathbf{w} . By *current* we mean that we use the data we have discovered in the n steps until now. To make this clear

we write X_N, Y_N instead of X, Y :

$$H[P(\mathbf{w}|X_N Y_N)] = - \int P(\mathbf{w}'|X_N, Y_N) \log[P(\mathbf{w}'|X_N, Y_N)] d\mathbf{w}'. \quad (10)$$

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:

$$K(x) = H[P(\mathbf{w}|X_N, Y_N, x, y = 0)] P(y = 0|X_N, Y_N, x) \quad (11)$$

$$+ H[P(\mathbf{w}|X_N, Y_N, x, y = 1)] P(y = 1|X_N, Y_N, x). \quad (12)$$

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

$$P(y = 0/1|X_N, Y_N, x) = \int P(y = 0/1|\mathbf{w}', x) P(\mathbf{w}'|X_N, Y_N) d\mathbf{w}' \quad (13)$$

24 1 Approximations

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

28 1.1 Restricting the tested stimuli

29 We need to discretize the x anyways. We choose values.

30 **1.2 Approximating the posterior**

31 Maybe MC methods?

32 **1.3 Using a different utility function**

33 BALD learning?