

Bayesian Active Learning

Assaf Dvora

May 3, 2022

1 Introduction

Bayesian Active Learning by Disagreement (BALD) [1] is an information theoretic approach for active learning designed for the Gaussian Process Classifier. The following sections describes this approach.

2 Bayesian Information Theoretic Active Learning

We consider a fully discriminative model where the goal of active learning is to discover the dependencies of some variable $y \in \mathcal{Y}$ on an input variable $\mathbf{x} \in \mathcal{X}$. The key idea in active learning is that the learner chooses the input queries $\mathbf{x}_i \in \mathcal{X}$ and observes the system's responses y_i , rather than passively receiving (\mathbf{x}_i, y_i) pairs.

Within a Bayesian framework we assume existence of some latent parameters, $\boldsymbol{\theta}$ (e.g. GP latent function), that controls the dependence between inputs and outputs through the conditional distribution $p(y|\mathbf{x}, \boldsymbol{\theta})$. Having observed data $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, a posterior distribution over the latent parameters is inferred, $p(\boldsymbol{\theta}|\mathcal{D})$. The central goal of information theoretic active learning is to reduce the number of possible hypotheses maximally fast. I.e., minimizing the uncertainty about the parameters using Shannon's entropy. Data points \mathcal{D}' are selected that satisfy

$$\arg \min_{\mathcal{D}'} H[\boldsymbol{\theta}|\mathcal{D}'] = - \int p(\boldsymbol{\theta}|\mathcal{D}') \log p(\boldsymbol{\theta}|\mathcal{D}') \quad (1)$$

Since solving this problem is NP-hard, a greedy policy is often used. Therefore the objective is to seek that data point \mathbf{x} that maximizes the decrease in posterior entropy [1]

$$\arg \max_{\mathbf{x}^*} \{H[\boldsymbol{\theta}|\mathcal{D}] - H[\boldsymbol{\theta}|y, \mathbf{x}, \mathcal{D}]\} \quad (2)$$

In practice, y is unknown and therefore the second term in Eq. (2) is replaced by the expected posterior entropy

$$\arg \max_{\mathbf{x}^*} \{H[\boldsymbol{\theta}|\mathcal{D}] - \mathbb{E}_{y \sim p(y|\mathbf{x}, \mathcal{D})} [H[\boldsymbol{\theta}|y, \mathbf{x}, \mathcal{D}]]\} \quad (3)$$

This solution however arises a computational difficulty: if $N_{\mathbf{x}}$ data points are under consideration, and N_y responses may be seen, then $\mathcal{O}(N_{\mathbf{x}}N_y)$ posterior updates are required.

An important insight arises if we note that the objective in Eq. (2) is equal to the mutual information between the latent parameters $\boldsymbol{\theta}$ and the unknown responses y^* , $I[\boldsymbol{\theta}, y^* | \mathbf{x}^*, \mathcal{D}]$. Using this insight it is simple to show that the objective can be rearranged to compute the entropies in the y space (see Appendix A.)

$$\arg \max_{\mathbf{x}^*} \{H[y | \mathbf{x}, \mathcal{D}] - \mathbb{E}_{\boldsymbol{\theta} \sim p(\boldsymbol{\theta} | \mathcal{D})} [H[y | \boldsymbol{\theta}, \mathbf{x}]]\} \quad (4)$$

Eq. (4) overcomes the computational difficulty described in Eq. (3). The latent parameter $\boldsymbol{\theta}$ is now conditioned only on \mathcal{D} , so only $\mathcal{O}(1)$ posterior updates are required. Eq. (4) also provides us with an interesting intuition about the objective; we seek the \mathbf{x} for which the model is marginally most uncertain about y (high $H[y | \mathbf{x}, \mathcal{D}]$), but for which, given individual settings of the parameters, y is confident (low $\mathbb{E}_{\boldsymbol{\theta} \sim p(\boldsymbol{\theta} | \mathcal{D})} [H[y | \boldsymbol{\theta}, \mathbf{x}]]$). Further analysis of the objective function is given in Section (5). We note that the argument in (4) is non-negative as we show in Appendix B.

3 Gaussian Process Classifier (GPC)

In this section we introduce the Gaussian Process Classifier (GPC) [2]. The probabilistic model underlying GPC is as follows

$$\begin{aligned} f(\mathbf{x}) &\sim \mathcal{GP}(0, k(\cdot, \cdot)) \\ y | \mathbf{x}, f &\sim \text{Bernoulli}(\Phi(f(\mathbf{x}))) \end{aligned}$$

where Φ is the Gaussian CDF. The latent parameter, now called $f(\mathbf{x})$, is a function $\mathcal{X} \rightarrow \mathbb{R}$, and is assigned a GP prior with zero mean and covariance function (or kernel).

Inference in the GPC model is non-Gaussian and intractable. Throughout this section we will assume that a Gaussian approximation of the posterior (e.g. Laplace approximation) is used. The posterior predictive distribution of $f(\mathbf{x})$ is then given by

$$p(f(\mathbf{x}) | \mathbf{x}, \mathcal{D}) = \mathcal{N}(f(\mathbf{x}) | \mu_{\mathbf{x}}, \sigma_{\mathbf{x}}^2) \quad (5)$$

The posterior predictive of y is then given by [2]

$$p(y | \mathbf{x}, \mathcal{D}) = \Phi \left(\frac{\mu_{\mathbf{x}}}{\sqrt{1 + \sigma_{\mathbf{x}}^2}} \right) \quad (6)$$

4 GPC Active Learning

In the GPC case, the objective function given in Eq. (4), takes the following form

$$\arg \max_{\mathbf{x}^*} \{H[y | \mathbf{x}, \mathcal{D}] - \mathbb{E}_{f_{\mathbf{x}} \sim p(f_{\mathbf{x}} | \mathbf{x}, \mathcal{D})} [H[y | f_{\mathbf{x}}]]\} \quad (7)$$

The first term in (7) which is the entropy of the posterior predictive of y can be handled analytically:

$$H[y|\mathbf{x}, \mathcal{D}] = -p(y|\mathbf{x}, \mathcal{D}) \log p(y|\mathbf{x}, \mathcal{D}) - (1 - p(y|\mathbf{x}, \mathcal{D})) \log(1 - p(y|\mathbf{x}, \mathcal{D}))$$

The second term in (7) involves integration over f space of the entropy of $y|f_{\mathbf{x}}$

$$\mathbb{E}_{f_{\mathbf{x}}} [H[y|f_{\mathbf{x}}]] = \int H[y|f_{\mathbf{x}}] \mathcal{N}(f_{\mathbf{x}}|\mu_{\mathbf{x}}, \sigma_{\mathbf{x}}^2) df_{\mathbf{x}} \quad (8)$$

where the entropy of $y|f_{\mathbf{x}}$ is given by

$$H[y|f_{\mathbf{x}}] = -p(y|f_{\mathbf{x}}) \log p(y|f_{\mathbf{x}}) - (1 - p(y|f_{\mathbf{x}})) \log(1 - p(y|f_{\mathbf{x}})) \quad (9)$$

$$= -\Phi(f(\mathbf{x})) \log \Phi(f(\mathbf{x})) - (1 - \Phi(f(\mathbf{x}))) \log(1 - \Phi(f(\mathbf{x}))) \quad (10)$$

5 Numerical Integration

To compute the objective function (7), one must compute the expectation $\mathbb{E}_{f_{\mathbf{x}} \sim p(f_{\mathbf{x}}|\mathbf{x}, \mathcal{D})} \{H[y|f_{\mathbf{x}}]\}$. Using the strong law of large number (SLLN) one can approximate the expectation by summation over samples from the posterior. Alternatively, numerical integration can be used to solve the integral

$$\begin{aligned} \mathbb{E}_{f_{\mathbf{x}}} [H[y|f_{\mathbf{x}}]] &= \int H[y|f_{\mathbf{x}}] \mathcal{N}(f_{\mathbf{x}}|\mu_{\mathbf{x}}, \sigma_{\mathbf{x}}^2) df_{\mathbf{x}} = \\ &= \frac{1}{\sqrt{2\pi\sigma_{\mathbf{x}}^2}} \int H[y|f_{\mathbf{x}}] \exp \left[-\frac{1}{2} \left(\frac{f_{\mathbf{x}} - \mu_{\mathbf{x}}}{\sigma_{\mathbf{x}}} \right)^2 \right] df_{\mathbf{x}} \end{aligned} \quad (11)$$

By using change of variables we obtain

$$z = \frac{f_{\mathbf{x}} - \mu_{\mathbf{x}}}{\sigma_{\mathbf{x}}}, \quad \frac{dz}{df_{\mathbf{x}}} = \frac{1}{\sigma_{\mathbf{x}}} \quad (12)$$

and (11) becomes

$$\mathbb{E}_{f_{\mathbf{x}}} [H[y|f_{\mathbf{x}}]] = \frac{1}{\sqrt{2\pi}} \int H[y|\sigma_{\mathbf{x}}z + \mu_{\mathbf{x}}] \exp[-\frac{z^2}{2}] dz \quad (13)$$

The integral can now be approximated using numerical integration

$$\mathbb{E}_{f_{\mathbf{x}}} [H[y|f_{\mathbf{x}}]] \approx \frac{1}{\sqrt{2\pi}} \sum_{i=1}^n H[y|\sigma_{\mathbf{x}}z_i + \mu_{\mathbf{x}}] \exp[-\frac{z_i^2}{2}] \Delta z \quad (14)$$

Appendix A.

Using Bayes rule we can express the objective in Eq. (3) in the y space

$$\begin{aligned}
& \arg \max_{\mathbf{x}} \{H[\boldsymbol{\theta}|\mathcal{D}] - \mathbb{E}_{y \sim p(y|\mathbf{x}, \mathcal{D})}[H[\boldsymbol{\theta}|y, \mathbf{x}, \mathcal{D}]]\} = \\
& \arg \min_{\mathbf{x}} \mathbb{E}_{y \sim p(y|\mathbf{x}, \mathcal{D})}[H[\boldsymbol{\theta}|y, \mathbf{x}, \mathcal{D}]] = \\
& \arg \min_{\mathbf{x}} \mathbb{E}_{y \sim p(y|\mathbf{x}, \mathcal{D})}[\mathbb{E}_{\boldsymbol{\theta} \sim p(\boldsymbol{\theta}|y, \mathbf{x}, \mathcal{D})}[-\log p(\boldsymbol{\theta}|y, \mathbf{x}, \mathcal{D})]] = \\
& \arg \max_{\mathbf{x}} \mathbb{E}_{y \sim p(y|\mathbf{x}, \mathcal{D})}[\mathbb{E}_{\boldsymbol{\theta} \sim p(\boldsymbol{\theta}|y, \mathbf{x}, \mathcal{D})}[\log p(y|\boldsymbol{\theta}, \mathbf{x}) + \log p(\boldsymbol{\theta}|\mathcal{D}) - p(y|\mathbf{x}, \mathcal{D})]] = \\
& \arg \max_{\mathbf{x}} \{H[y|\mathbf{x}, \mathcal{D}]
\end{aligned}$$

Appendix B.

In this section we show that the argument in the objective function (4) is non-negative.

Denote with $h(\cdot)$ the convex function $h(p) = p \log p, p > 0$ and denote with $\psi_{k,\mathbf{x}}(\cdot)$ the random variable $\psi_{k,\mathbf{x}}(\boldsymbol{\theta}) = p(y = k|\boldsymbol{\theta}, \mathbf{x})$, using Jensen inequality we obtain

$$\begin{aligned}
\mathbb{E}_{\boldsymbol{\theta} \sim p(\boldsymbol{\theta}|\mathcal{D})}[H[y|\boldsymbol{\theta}, \mathbf{x}]] &= \mathbb{E}_{\boldsymbol{\theta} \sim p(\boldsymbol{\theta}|\mathcal{D})}\left\{-\sum_{k=1}^K h[\psi_{k,\mathbf{x}}(\boldsymbol{\theta})]\right\} = \\
&= -\sum_{k=1}^K \mathbb{E}_{\boldsymbol{\theta} \sim p(\boldsymbol{\theta}|\mathcal{D})}\{h[\psi_{k,\mathbf{x}}(\boldsymbol{\theta})]\} = \\
&= -\sum_{k=1}^K \int h[\psi_{k,\mathbf{x}}(\boldsymbol{\theta})]p(\boldsymbol{\theta}|\mathcal{D})d\boldsymbol{\theta} \leqslant \\
&\leqslant -\sum_{k=1}^K h\left[\int \psi_{k,\mathbf{x}}(\boldsymbol{\theta})p(\boldsymbol{\theta}|\mathcal{D})d\boldsymbol{\theta}\right] = \\
&= -\sum_{k=1}^K h[p(y = k|\mathbf{x}, \mathcal{D})] = H[y|\mathbf{x}, \mathcal{D}]
\end{aligned}$$

Thus, $H[y|\mathbf{x}, \mathcal{D}] \geqslant \mathbb{E}_{\boldsymbol{\theta} \sim p(\boldsymbol{\theta}|\mathcal{D})}[H[y|\boldsymbol{\theta}, \mathbf{x}]]$.

References

- [1] Houlshby, Neil, et al. "Bayesian active learning for classification and preference learning." arXiv preprint arXiv:1112.5745 (2011).
- [2] Rasmussen, C. and Williams, C. (2005). Gaussian Processes for Machine Learning. The MIT Press.