

1 Bayesian experimental design for model selection: variational and classification approaches

1.1 Introduction

Bayesian experimental design for model selection is an important and well-studied problem (Cavagnaro et al., 2010; Vanlier et al., 2014; Hainy et al., 2018). In this essay, we tackle two questions that are relevant to this problem. First, how do recently proposed variational methods for experimental design (Foster et al., 2019, 2020) translate into the model selection context? Second, how do these methods intersect with recently proposed classification-driven approaches to experimental design for model selection (Hainy et al., 2018)?

We begin by elucidating the key features of the model selection problem—it turns out that we can characterise the set-up as a semi-implicit model with a discrete latent variable of interest. The posterior or Barber–Agakov approach of Foster et al. (2019) involves training an amortised inference network from data simulated from the model. We find that, for model selection, this network is exactly a (neural) classifier that predicts the true model that synthesised an observation from that synthetic experimental observation. The marginal + likelihood method of Foster et al. (2019) also translates into the model selection case. This method involves variational density estimation of experimental outcomes for each possible model. In other words, it involves approximating the model evidence of the data for each possible model. Finally, we examine how the stochastic gradient design approach of Foster et al. (2020) applies here. This approach can build off the back of the Barber–Agakov bound, so it also utilises a classifier. The key difference here is that we differentiate the classifier output with respect to its input to learn the design at the same time as the classifier network parameters. This bears some similarities with adversarial approaches to neural network robustness (Carlini et al., 2019). Finally, we compare and contrast the variational approach with other classification driven approaches in the literature.

1.2 Characterising the problem

We denote experimental designs by ξ and experimental observations as y . Suppose there are K competing models $\{m_1, \dots, m_k\}$ and we have a prior distribution $p(m)$ on which model we think is likely to be correct. Given the choice of model, there are other model parameters $\psi \sim p(\psi|m)$. Conditional on the model, and on its parameters, we have a likelihood for the experiment $p(y|m, \psi, \xi)$ which we assume is known in closed form.

One important feature of the model selection problem is that we do *not* have a likelihood that directly relates the design ξ , observation y and the latent variable of interest m . Instead, we have to account for the auxiliary latent variable ψ . Indeed, we actually have $p(y|m, \xi) = \int_{\Psi} p(y|m, \psi, \xi) p(\psi|m) d\psi$. This case, where we have a closed form likelihood but for a larger set of variable, is referred to as a *semi-implicit* model.

In this essay, we focus on experimental design with the expected information gain (EIG) criterion, also called mutual information utility, that aims to reduce Shannon entropy in our beliefs about m . The EIG-optimal design is specifically,

$$\xi^* = \arg \max_{\xi} \mathbb{E}_{p(m)p(\psi|m)p(y|m, \psi, \xi)} \left[\log \frac{p(m|y, \xi)}{p(m)} \right]. \quad (1)$$

Finding ξ^* amounts to estimating the EIG objective function and optimising over the space of possible designs.

If we have already observed some data $\mathcal{D} = \{(\xi_1, y_1), \dots, (\xi_T, y_T)\}$, then we fit model-specific posteriors for the auxiliary variable ψ for each model $p(\psi|m, \mathcal{D})$, and we compute the posterior over models $p(m|\mathcal{D}) \propto p(m)p(\mathcal{D}|m)$. Thus, we update our priors $p(m)$ and $p(\psi|m)$ on the basis of past data.

1.3 The variational approach

1.3.1 Posterior lower bound

Foster et al. (2019) considered variational estimation of the EIG. Their general strategy was to optimise variational upper or lower bounds on the EIG. Their simplest bound was the posterior lower bound (also called the Barber–Agakov bound after Barber and Agakov (2003)). With the variables we have in this model, the bound would be expressed as

$$\mathbb{E}_{p(m)p(\psi|m)p(y|m,\psi,\xi)} \left[\log \frac{p(m|y,\xi)}{p(m)} \right] \geq \mathbb{E}_{p(m)p(\psi|m)p(y|m,\psi,\xi)} \left[\log \frac{q_\phi(m|y)}{p(m)} \right]. \quad (2)$$

The new term $q_\phi(m|y,\xi)$ was generically referred to as the *amortised approximate posterior* with variational parameters ϕ . It is an approximate posterior distribution on the latent variable m of interest. The amortisation here refers to the fact that we learn a function from y to a distribution over m (for different ξ , we would train separate functions). For the model selection approach, then, q_ϕ is a function from y to a distribution over the discrete model indicator m . First, since m is discrete, the choice of variational family is moot, because every distribution over m can be finitely represented. Second, q_ϕ has a very simple interpretation. It is a classifier that attempts to predict, on the basis of input y , which model of m_1, \dots, m_k generated that data, specifically trying to estimate the posterior probability $p(m|y,\xi)$ over the k different possibilities for m . Importantly though, rather than just attempting to predict the correct model that was responsible for generating the data y , it is essential that we have a *probabilistic* classifier that assigns probabilities to each possible model. For this probabilistic classifier, the issue of calibration becomes central, as we hope that our classifier probabilities will approach $p(m|y,\xi)$ during training.

We have established that q_ϕ is simply a probabilistic classifier for the model selection case. How should this classifier be trained? In general, Foster et al. (2019) proposed training q_ϕ by stochastic gradient methods (Robbins and Monro, 1951; Kingma and Ba, 2014) to maximise the lower bound with respect to ϕ

$$\phi^* = \arg \max_{\phi} \mathbb{E}_{p(m)p(\psi|m)p(y|m,\psi,\xi)} \left[\log \frac{q_\phi(m|y)}{p(m)} \right] \quad (3)$$

In model selection, training ϕ simply means training the parameters of the classifier. Maximising the posterior lower bound is equivalent to simply maximising the expected log likelihood under q , i.e.

$$\phi^* = \arg \max_{\phi} \mathbb{E}_{p(m)p(\psi|m)p(y|m,\psi,\xi)} [\log q_\phi(m|y)]. \quad (4)$$

This is true because $p(m)$ has no dependence on ϕ . So, we see that training q_ϕ to maximise the variational posterior lower bound amounts to maximum likelihood training of a neural classifier when we are in the setting of model selection. (Care may be needed to ensure the classifier produces good *probabilistic uncertainty*, as well as getting good predictions, as these probabilities are central to our method.)

In fact, we have an enhanced setting in which we can draw an infinite amount of training data by simulating from $p(m)p(\psi|m)p(y|m,\psi,\xi)$. To do this, we sample a random model m from its prior, then a random set of parameters $\psi \sim p(\psi|m)$ for the chosen model, and then simulate an experimental outcome under design ξ . Importantly, we do not need to draw a fixed training or test set, and we never need to show the classifier the same examples twice, we instead draw new batches on the fly. One particularly important consequence of this is that the spectre of *over-fitting* is much reduced in our case, as there is no fixed training set to overfit to.

We now see another important point—the negative log-likelihood loss of the classifier is essentially an estimate of the EIG, up to a constant. Suppose we have completed training and reached parameters $\hat{\phi}$. Then the EIG estimate is

$$\text{EIG}(\xi) \approx \mathbb{E}_{p(m)p(\psi|m)p(y|m,\psi,\xi)} \left[\log \frac{q_{\hat{\phi}}(m|y)}{p(m)} \right] = \mathbb{E}_{p(m)p(\psi|m)p(y|m,\psi,\xi)} [\log q_{\hat{\phi}}(m|y)] + H[p(m)] \quad (5)$$

and we can estimate the expectation with new, independent batches simulated from the model.

In summary, the posterior lower bound method for model selection amounts to training a classifier on (infinite) simulated data to predict m from y . The optimal design ξ^* will be approximated by the classifier which has the best (lowest) validation loss, which is a good approximation of having the highest EIG.

1.3.2 Marginal + likelihood estimator

The posterior lower bound is not the only way to estimate the EIG proposed by Foster et al. (2019). Both the marginal and the VNMC methods require an explicit likelihood, so they are not suitable for the semi-implicit model selection scenario. The marginal + likelihood estimator is

$$\text{EIG}(\xi) \approx \mathbb{E}_{p(m)p(\psi|m)p(y|m,\psi,\xi)} \left[\log \frac{q_\ell(y|m,\xi)}{q_p(y|\xi)} \right]. \quad (6)$$

This estimator translates, with some simplification, into the model selection setting. The ‘approximate likelihood’ $q_\ell(y|m,\xi)$ in the model selection setting is an approximation of the model evidence $q_\ell(y|m,\xi) \approx p(y|m,\xi)$. For model selection when m is discrete, we do not need to separately estimate q_p and q_ℓ , we can instead sum over m to obtain

$$q_p(y|\xi) = \sum_m p(m)q_\ell(y|m,\xi). \quad (7)$$

As shown in Appendix A.4 of Foster et al. (2019), the estimator actually becomes a lower bound

$$\text{EIG}(\xi) \geq \mathbb{E}_{p(m)p(\psi|m)p(y|m,\psi,\xi)} \left[\log \frac{q_\ell(y|m,\xi)}{\sum_{m'} p(m')q_\ell(y|m',\xi)} \right] \quad (8)$$

on the EIG in this case, which is not generally the case for the marginal + likelihood method. (In fact, this lower bound is itself a special case of the likelihood-free ACE lower bound introduced in Foster et al. (2020). Indeed, if we take the prior as the variational posterior and let $L \rightarrow \infty$ in the LF-ACE bound, we recover this lower bound.)

This lower bound also has a nice interpretation in the model selection scenario. The best design will be the one where the lower bound is largest, which happens, loosely speaking, when $q_\ell(y|m,\xi)$ is much larger than $\sum_m p(m)q_\ell(y|m,\xi)$. That means the approximate model evidence for the observation y under the correct model m is much larger than its evidence under other models. Thus, using the experiment with design ξ and observing y will allow us to easily discriminate between models.

To explicitly use this method, we need to choose trainable density estimators for $q_\ell(y|m,\xi;\phi)$ with parameters ϕ . The simplest method would be to have a distinct set of variational parameters for each value of m and ξ . Whilst it is possible to use a Gaussian density model, we could use more sophisticated methods such as normalising flows (Rezende and Mohamed, 2015). The training approach is similar to that for the posterior method. We use infinite simulated data, and maximise the variational lower bound using stochastic gradient optimisers.

The last two sections highlight a general feature of the variational methods of Foster et al. (2019)—we can either make variational approximations to densities over m or over y . Both lead to valid bounds.

1.4 Stochastic gradient optimisation of the design

So far, we have focused on variational estimation of the EIG. As shown in Foster et al. (2020), it is only a short jump from variational estimation of the EIG to stochastic gradient optimisation of the design using a variational lower bound on EIG. The benefit here, of course, is that we do not have to conduct a grid search, co-ordinate exchange or similar algorithm over the design space. What we require instead is a continuous design space and the ability to differentiate observations with respect to designs.

Whilst Foster et al. (2020) focused on explicit likelihood models, both the posterior (Barber–Agakov) lower bound and the LF-ACE bound are applicable to the semi-implicit model selection setting. There is just one thing to check, which is that we can compute a derivative $\partial y / \partial \xi$. In the semi-implicit case, this is often fine.

For example, if $p(y|m, \psi, \xi)$ takes the form $y = g(m, \psi, \xi, \epsilon)$ for a differentiable g and an independent noise random variable ϵ .

Assuming this is the case, we can train ξ by stochastic gradient using either the posterior bound or the simplified LF-ACE bound that was derived in equation (8). We focus on the posterior lower bound for simplicity. Recall that, for the posterior bound, we are training a classifier to predict m from y . We have

$$\text{EIG}(\xi) \geq \mathbb{E}_{p(m)p(\psi|m)p(y|m, \psi, \xi)} [\log q_\phi(m|y)] + H[p(m)] \quad (9)$$

where q_ϕ is the classifier. One thing that we skimmed over slightly in the previous section was that ϕ implicitly depends on ξ via the training data, and different ξ will have different classifiers with different optimal values of the classifier parameters ϕ .

In Foster et al. (2020), rather than training separate classifiers with different designs ξ , we update ξ and ϕ together in one stochastic gradient optimisation over the combined set of variables (ξ, ϕ) . To explicitly write down the ξ gradient here, let's assume that we do have $y = g(m, \psi, \xi, \epsilon)$, so we can write

$$\mathcal{L}(\xi, \phi) = \mathbb{E}_{p(m)p(\psi|m)p(\epsilon)} [\log q_\phi(m|g(m, \psi, \xi, \epsilon))] + H[p(m)]. \quad (10)$$

In this form, the ξ gradient can be simply calculated as

$$\frac{\partial \mathcal{L}}{\partial \xi} = \mathbb{E}_{p(m)p(\psi|m)p(\epsilon)} \left[\frac{\partial \log q_\phi}{\partial y} \bigg|_{m, g(m, \psi, \xi, \epsilon)} \frac{\partial g}{\partial \xi} \bigg|_{m, \psi, \xi, \epsilon} \right]. \quad (11)$$

The beauty of modern auto-diff frameworks, of course, means that we do not even need to calculate this explicitly ourselves.

For model selection, equation (11) has a natural interpretation. We want to increase the lower bound \mathcal{L} by moving to regions in which the classifier can confidently predict the correct model label m . This corresponds to moving y into regions in which $\log q_\phi(m|y)$ is larger *for the model that actually generated y* . In other words, we want the input to the classifier y to be pushed to regions where the classifier already finds it easy to classify correctly. That is, regions where deciding which model is correct is easier. We then exploit the differentiable relationship between ξ and y , and use this signal to ‘improve’ the input to the classifier by adjusting the design ξ to that such datasets y are more likely to be synthesised.

At the same time, we are constantly making gradient updates on the classifier parameters ϕ . This means that, as the distribution of (m, y) changes, the classifier can adjust accordingly.

If this sounds dubious, it is worth taking a step back. We are quite simply optimising the lower bound $\mathcal{L}(\xi, \phi)$ jointly with respect to ξ and ϕ , in the hopes that this global maximum may closely correspond to the EIG maximiser ξ^* . We actually have a guarantee that the value of \mathcal{L} at our final trained variables $\hat{\xi}, \hat{\phi}$ is a lower bound on $\text{EIG}(\hat{\xi})$, i.e. the true value of $\hat{\xi}$ cannot be worse than the value we estimate for it.

Whilst the method is approximate, because we cannot quantify the discrepancy between \mathcal{L} and the true EIG, it is highly scalable to very large design spaces. Other bounds presented in Foster et al. (2020) have the added benefit that they become equal to the EIG in a limit, providing some assurances that the global maximum of \mathcal{L} is a good design. Foster et al. (2020) also introduced the evaluation method of establishing *lower and upper* bounds on chosen designs. This numerically bounds the discrepancy between the training objective \mathcal{L} and the true EIG objective. Sadly, the upper bounds are only valid for explicit likelihood models; they don't work in the semi-implicit model selection case.

Finally, all of the above discussion carries over if we were to use the lower bound of equation (8) instead of the posterior bound.

1.5 Comparing with other classification approaches

We have established that the variational posterior approach of Foster et al. (2019) instructs us to learn a classifier to predict m from y and use the log probabilities $q_\phi(m|y)$ to estimate EIG. Other authors have considered supervised classification as a means to perform Bayesian experimental design for model selection.

Here, we focus on Hainy et al. (2018), which is “the first approach using supervised learning methods for optimal Bayesian design.” This method trains a classifier that predicts m using y , with separate classifiers for different ξ . They focus on training decision trees and random forest classifiers (Breiman, 2001). Since random forests are not generally trained by stochastic gradient methods, this means that they fall back on simulating fixed training and test datasets of samples $(m_j, y_j)_{j=1}^J$ from $p(m)p(\psi|m)p(y|m, \psi, \xi)$. The training dataset is used to train the classifier model, whilst the test dataset gives unbiased estimates of the posterior loss. There is a danger that the classifier may overfit to the training set in this case. Compare this with the training of stochastic gradient classifiers in our previous sections—here we can draw fresh training batches on the fly, and avoid overfitting to a training set.

Decision trees and random forests do provide estimates of the class probabilities $q(m|y)$, but they are relatively noisy. For this reason, Hainy et al. (2018) focus on the 0–1 loss to evaluate designs. In the language of classification, therefore, they choose the design which gives the best *test accuracy*. Again, this is different to the variational approach which fits a neural classifier that automatically provides smooth probability estimates $q_\phi(m|y)$. The latter case was applied to estimate the information gain, which we showed is equivalent to choosing the design which gives the best *test loss*, assuming a negative log-likelihood loss function.

The trade-offs between these methods are clear when we consider optimising over a large design space. For the variational method, we have to train a number of neural networks to convergence. For the classification approach of Hainy et al. (2018), we train a number of random forest classifiers—this may be significantly more computationally efficient. Hainy et al. (2018) propose embedding their 0–1 loss estimation within a co-ordinate exchange algorithm (Meyer and Nachtsheim, 1995) to optimise over designs. The variational method, on the other hand, can naturally be embedded in a unified stochastic gradient optimisation to find the optimal design through stochastic gradient optimisation. The former may be more effective when the design space is not continuous, the latter can work well in a high-dimensional design space that is difficult to search using discrete methods.

2 Bayesian active learning by disagreement and Bayesian experimental design

2.1 Introduction

The purpose of this essay is to highlight the connection between the Bayesian Active Learning by Disagreement (BALD) score as estimated by Gal et al. (2017) and the Prior Contrastive Estimation (PCE) bound of Foster et al. (2020). There is a deep connection between Bayesian experimental design and Bayesian active learning. A significant touchpoint is the use of the mutual information score (Lindley, 1956)

$$I(\xi) = \mathbb{E}_{p(\theta)p(y|\theta,\xi)} [H[p(\theta)] - H[p(\theta|y, \xi)]] . \quad (12)$$

to acquire new information in a Bayesian model with parameters θ where, y is the as yet unobserved outcome, and ξ is the design to be chosen.

2.2 Bayesian Active Learning by Disagreement

One of the computational challenges inherent in estimating equation (12) directly is that it involves repeated estimation of posterior distributions $p(\theta|y, \xi)$ for different simulated observations y . To remove this particular bottleneck, Houlby et al. (2011) introduced a rewriting of the mutual information score using Bayes rule

$$I(\xi) = H[p(y|\xi)] - \mathbb{E}_{p(\theta)} [H[p(y|\theta, \xi)]] . \quad (13)$$

Whilst this is exactly equal to the original mutual information score, the new way of expressing I removes the requirement to estimate posterior distributions over θ . They termed equation (13) the Bayesian Active Learning by Disagreement (BALD) score.

Unfortunately, the story does not end with the BALD score because it still typically involves some intractable computations that must be estimated. For example, Houlby et al. (2011) focused on approximations for Gaussian Process models (Williams and Rasmussen, 2006).

The more recent work by Gal et al. (2017) estimated the BALD score in the context of Bayesian deep learning classifiers. In such a model, θ represents the parameters of a classification model, and $p(y|\theta, \xi)$ is a probability distribution over classes $y \in \{c_1, \dots, c_k\}$. Computing $p(y|\theta, \xi)$ involves a forward pass through the classifier with input ξ and parameters θ , the network generally ends in a softmax activation to produce a normalised distribution. To sample different values of θ , Gal et al. (2017) employed Monte Carlo Dropout (Gal and Ghahramani, 2016). Given independent samples $\theta_1, \dots, \theta_M$ from $p(\theta)$, they proposed the following Deep BALD (DBALD) estimator of $I(\xi)$

$$I(\xi) \approx \hat{I}_{\text{DBALD}}(\xi) = H \left[\frac{1}{M} \sum_{i=1}^M p(y|\theta_i, \xi) \right] - \frac{1}{M} \sum_{i=1}^M H[p(y|\theta_i, \xi)] \quad (14)$$

where $H[P(y)] = -\sum_c P(y=c) \log P(y=c)$.

Notation For comparison with the original paper, we used θ in place of ω , ξ in place of \mathbf{x} , M in place of T and $p(\theta)$ is used in place of $q_\theta^*(\omega)$.

2.3 Prior Contrastive Estimation

In the context of stochastic gradient optimisation of Bayesian experimental designs, Foster et al. (2020) also considered the mutual information score $I(\xi)$ and the rearrangement equation (13). They proved the following Prior Contrastive Estimation (PCE) lower bound on $I(\xi)$

$$I(\xi) \geq \mathbb{E}_{p(\theta_0)p(y|\theta_0,\xi)p(\theta_1)\dots p(\theta_L)} \left[\log \frac{p(y|\theta_0, \xi)}{\frac{1}{L+1} \sum_{\ell=0}^L p(y|\theta_\ell, \xi)} \right] \quad (15)$$

and used this bound to optimise ξ by stochastic gradient. One approach to estimate this bound using finite samples is the estimator

$$\hat{I}_{\text{PCE-naive}}(\xi) = \frac{1}{M} \sum_{m=1}^M \log \frac{p(y_m | \theta_{m0}, \xi)}{\frac{1}{L+1} \sum_{\ell=0}^L p(y_m | \theta_{m\ell}, \xi)}. \quad (16)$$

where $y_m, \theta_{m0} \sim p(y, \theta | \xi)$ and $\theta_{m\ell} \sim p(\theta)$ for $\ell \geq 1$. However, we can also re-use samples more efficiently to give the estimator

$$\hat{I}_{\text{PCE}}(\xi) = \frac{1}{M} \sum_{m=1}^M \log \frac{p(y_m | \theta_m, \xi)}{\frac{1}{M} \sum_{\ell=1}^M p(y_m | \theta_\ell, \xi)}. \quad (17)$$

where $y_m, \theta_m \sim p(y, \theta | \xi)$. (To check the expectation of this version matches the PCE bound with $L = M - 1$, we simply move the \mathbb{E} sign inside of the summation.) Finally, Foster et al. (2020) discussed a speed-up that is possible when y is a discrete random variable taking values in $\{c_1, \dots, c_k\}$. In this case, we can integrate out y by summing over it, rather than by drawing random samples of y . This method, called Rao-Blackwellisation, results in the estimator

$$\hat{I}_{\text{PCE-RB}}(\xi) = \frac{1}{M} \sum_{m=1}^M \sum_c p(y = c | \theta_m, \xi) \log \frac{p(y = c | \theta_m, \xi)}{\frac{1}{M} \sum_{\ell=1}^M p(y = c | \theta_\ell, \xi)}. \quad (18)$$

2.4 PCE and DBALD equivalence

We have looked at two parallel ways of approximating $I(\xi)$. The interesting result is that *the Rao-Blackwellised PCE estimator and the DBALD estimator are the same*. We can see this by direct calculation

$$\hat{I}_{\text{PCE-RB}}(\xi) = \frac{1}{M} \sum_{m=1}^M \sum_c p(y = c | \theta_m, \xi) \log \frac{p(y = c | \theta_m, \xi)}{\frac{1}{M} \sum_{\ell=1}^M p(y = c | \theta_\ell, \xi)} \quad (19)$$

$$= \frac{1}{M} \sum_{m=1}^M \sum_c p(y = c | \theta_m, \xi) \log p(y = c | \theta_m, \xi) - \frac{1}{M} \sum_{m=1}^M \sum_c p(y = c | \theta_m, \xi) \log \left(\frac{1}{M} \sum_{\ell=1}^M p(y = c | \theta_\ell, \xi) \right) \quad (20)$$

$$= -\frac{1}{M} \sum_{m=1}^M H[p(y | \theta_m, \xi)] - \frac{1}{M} \sum_{m=1}^M \sum_c p(y = c | \theta_m, \xi) \log \left(\frac{1}{M} \sum_{\ell=1}^M p(y = c | \theta_\ell, \xi) \right) \quad (21)$$

$$= -\frac{1}{M} \sum_{m=1}^M H[p(y | \theta_m, \xi)] - \sum_c \left(\frac{1}{M} \sum_{m=1}^M p(y = c | \theta_m, \xi) \right) \log \left(\frac{1}{M} \sum_{\ell=1}^M p(y = c | \theta_\ell, \xi) \right) \quad (22)$$

$$= -\frac{1}{M} \sum_{m=1}^M H[p(y | \theta_m, \xi)] + H \left[\frac{1}{M} \sum_{m=1}^M p(y | \theta_m, \xi) \right] \quad (23)$$

$$= \hat{I}_{\text{DBALD}}(\xi). \quad (24)$$

A major consequence of this result is that *the expectation of the DBALD score is a lower bound on the true mutual information score*. We also note that this estimator has been used by Vincent and Rainforth (2017) in the context of Bayesian experimental design, although they did not show that it was a stochastic lower bound.

2.5 New diagnostic for the DBALD score

One advantage of making this connection is that we can bring certain diagnostics that were applied by Foster et al. (2020) over to the active learning setting. In particular, Foster et al. (2020) paired their PCE lower

bound with a complementary *upper bound* on $I(\xi)$. This provides a very useful diagnostic tool to tune the number of samples M used to compute the DBALD score. If the lower bound and upper bound are very close, we know that the difference between the DBALD score and the true mutual information must also be small. On the other hand, if the upper and lower bounds are far apart, then the DBALD score might not yet be close to the true mutual information.

One upper bound upper by Foster et al. (2020) was the Nested Monte Carlo (NMC) (Vincent and Rainforth, 2017) estimator. For the discrete y case with Rao-Blackwellisation, the estimator is

$$\hat{I}_{\text{NMC-RB}}(\xi) = -\frac{1}{M} \sum_{m=1}^M \sum_c p(y=c|\theta_m, \xi) \log \left(\frac{1}{M-1} \sum_{\ell \neq m} p(y=c|\theta_\ell, \xi) \right) - \frac{1}{M} \sum_{m=1}^M H[p(y|\theta_m, \xi)] \quad (25)$$

$$= \frac{1}{M} \sum_{m=1}^M H \left[p(y|\theta_m, \xi), \frac{1}{M-1} \sum_{\ell \neq m} p(y|\theta_\ell, \xi) \right] - \frac{1}{M} \sum_{m=1}^M H[p(y|\theta_m, \xi)] \quad (26)$$

where $H[p, q]$ is the cross-entropy. The expectation of this mutual information estimator is always an upper bound on $I(\xi)$. So both the DBALD score and the NMC-RB estimator converge to $I(\xi)$ as $M \rightarrow \infty$, but from opposite directions. We suggest NMC-RB as a diagnostic for the parameter M .

2.6 BALD estimators for regression

The connection to PCE may also be helpful when considering regression models. The standard parametrisation of a Bayesian neural network for regression is for the output of the network with parameters θ and input ξ to be the predictive mean μ and standard deviations σ of a Gaussian $y|\theta, \xi \sim N(\mu(\theta, \xi), \sigma(\theta, \xi)^2)$. (It is normal for y , μ and σ to be vector-valued and for the Gaussian to have a diagonal covariance matrix.)

For the DBALD estimator for a regression model, the entropy of a Gaussian is known in closed form, so $H[p(y|\theta_i, \xi)] = \frac{1}{2} \log(2\pi e \sigma(\theta_i, \xi)^2)$. However, the entropy of a mixture of Gaussians $H\left[\frac{1}{M} \sum_{i=1}^M p(y|\theta_i, \xi)\right]$ cannot be computed analytically. Instead, we could estimate this mixture of Gaussians entropy using Monte Carlo by sampling $i \in \{1, \dots, M\}$ uniformly, sampling y from $p(y|\theta_i, \xi)$ and calculating the log-density at y .

Despite the fact that we are using an analytic entropy for one term, and a Monte Carlo estimate for the other, it's easy to see that this new estimator is a *partially Rao-Blackwellised* PCE estimator. (This can be proved starting from equation (17).) That means all the existing facts, such as the estimator being a stochastic lower bound on $I(\xi)$, carry over naturally to the regression case.

3 Deep Adaptive Design and Bayesian reinforcement learning

3.1 Introduction

The purpose of this essay is to discuss the connections between the recently proposed Deep Adaptive Design (DAD) (Foster et al., 2021) method and the field of Bayesian reinforcement learning (Ghavamzadeh et al., 2016). That such a connection exists is hinted at by a high-level appraisal of the DAD method—it solves a sequential decision making problem to optimise a certain objective function, decision optimality is dependent on a *state* which is the experimental data already gathered, and the automated decision maker is a design *policy* network. We begin by showing how the sequential Bayesian experimental design problem solved by DAD can be viewed as a Bayes Adaptive Markov Decision Process (BAMDP) (Ross et al., 2007; Guez et al., 2012), making this connection formally precise. We also isolate some of the key differences between the problem DAD is solving and a conventional Bayesian RL problem, noting that the reward in DAD is intractable. Much of the effort of DAD is in establishing a differentiable surrogate for the true objective. The differentiability of the surrogate reward is also a key feature of the DAD problem, which facilitates the direct policy optimisation approach taken to train the policy that is rarely applicable in standard RL problems. We also highlight other features of the DAD method, such as its avoidance of explicitly estimating any posterior distributions, i.e. the avoidance of explicit belief state estimation.

Having studied DAD in some detail, we consider possible extensions of the method that make use of the RL connection. First, there are rather natural extensions of DAD to more general objective functions that incorporate design costs, terminal decisions and other functionals of the posterior distribution. Second, more standard approaches to (Bayesian) RL, such as Q-learning (Watkins and Dayan, 1992; Dearden et al., 1998) can be applicable to the sequential Bayesian experimental design problem. They may be particularly useful for long- or infinite-horizon problems.

3.2 Background on Bayesian Reinforcement Learning

3.2.1 Markov Decision Processes

The Markov Decision Process (MDP) (Bellman, 1957; Duff, 2002) is a highly successful mathematical framework for sequential decision problems in a known environment. Formally, a MDP consists of a state space S , an action space A , a transition model \mathcal{P} , a reward distribution R , a discount factor $0 \leq \gamma \leq 1$ and a time horizon T which may be infinite. An agent operates in the MDP by moving between different states in discrete time. For example, if the agent is in state s_t at time t and chooses to play action a_t , then the next state s_{t+1} will be sampled randomly according to the transition model $s_{t+1} \sim \mathcal{P}(s|s_t, a_t)$. Since the distribution over the next state depends only on s_t and a_t , the transitions are Markovian. Finally, by making the transition $s_t \xrightarrow{a_t} s_{t+1}$, the agent receives a random reward $r_t \sim R(r|s_t, a_t, s_{t+1}) \in \mathbb{R}$. The agent’s objective is to maximise the discounted sum of rewards $\sum_{t=0}^T \gamma^t r_t$. Given the Markovian nature of the problem, it is sufficient to choose actions according to some *policy* π , where $a_t = \pi(s_t)$. The optimality condition for a policy is

$$\pi^* = \arg \max_{\pi} \mathcal{J}(\pi), \quad (27)$$

where

$$\mathcal{J}(\pi) = \mathbb{E}_{s_0 \sim p(s_0) \prod_{t=0}^T a_t = \pi(s_t), s_{t+1} \sim \mathcal{P}(s|s_t, a_t), r_t \sim R(r|s_t, a_t, s_{t+1})} \left[\sum_{t=0}^T \gamma^t r_t \right]. \quad (28)$$

In a classical MDP, we assume that \mathcal{P} and R are known during the planning phase, when the agent devises their policy π . Of particular utility in planning a policy is the value function, defined as

$$V^\pi(s) = \mathbb{E}_{s' \sim \mathcal{P}(\cdot|s, \pi(s)), r \sim R(r|s, \pi(s), s')} [r + \gamma V^\pi(s')] \quad (29)$$

and the Q -function

$$Q^\pi(s, a) = \mathbb{E}_{s' \sim \mathcal{P}(\cdot|s, a), r \sim R(r|s, a, s')} [r + \gamma V^\pi(s')]. \quad (30)$$

These equations are valid when $T = \infty$, for finite time horizon we also have to take account of time t in state evaluations.

3.2.2 Bayes Adaptive Markov Decision Processes

The BAMDP (Duff, 2002; Ross et al., 2007; Guez et al., 2012; Ghavamzadeh et al., 2016) is one approach to generalising the MDP to deal with unknown transition models. In the BAMDP, the agent retains an explicit posterior distribution over the transition model called a belief state. This allows a formally elegant approach to behaviour under uncertainty which can trade off exploration (learning the transition model) and exploitation (executing actions that receive a high reward).

To set this up formally using the notation of Guez et al. (2012), we begin by considering an outer probabilistic model over the transition probabilities with prior $P(\mathcal{P})$. Given a history of states, actions and rewards $h_t = s_0 a_0 \dots r_{t-1} a_{t-1} s_t$, we can compute a posterior distribution on \mathcal{P} by

$$P(\mathcal{P}|h_t) \propto P(\mathcal{P})P(h_t|\mathcal{P}) = P(\mathcal{P}) \prod_{\tau=0}^t \mathcal{P}(s_{\tau+1}|s_{\tau}, a_{\tau}). \quad (31)$$

To bring this back into the MDP formulation, we consider an augmented state space S^+ which consists of entire histories, and which encapsulates both the current state and our beliefs about the transition model. Transitions in the augmented state space S^+ are given by integrating over the current beliefs on \mathcal{P}

$$\mathcal{P}^+(h_{t+1}|h_t, a_t) = \int P(\mathcal{P}|h_t) \mathcal{P}(s_{t+1}|s_t, a_t) d\mathcal{P}. \quad (32)$$

It is also possible for BAMDPs to incorporate unknown reward distributions (see e.g. Zintgraf et al. (2019)), where an outer model over reward distributions is updated on the basis of h_t in the same manner as for the transition probabilities. Specifically, if we have a prior $P(R)$ over reward distributions, then the reward function for playing action a_t in augmented state h_t is

$$R^+(r|h_t, a_t, h_{t+1}) = \int P(R|h_{t+1}) R(r|s_t, a_t, s_{t+1}) dR. \quad (33)$$

Combining these gives a new MDP with state space S^+ of histories, unchanged action space A , augmented transition model \mathcal{P}^+ , augmented reward distribution R^+ , discount factor γ and time horizon T . Optimal action in this new MDP gives the optimal trade-off between exploration and exploitation.

3.3 The Bayesian RL formulation of DAD

In DAD (Foster et al., 2021), we choose a sequence of designs ξ_1, \dots, ξ_T with a view to maximising the expected information gained about a latent parameter of interest θ . To place DAD in a Bayesian RL setting, we begin by associating the design ξ_t chosen before observing an outcome with the action a_{t-1} . The difference in time labels is necessary because ξ_t is chosen before y_t is observed. Since the observation distribution $p(y|\xi, \theta)$ depends on the unknown θ , we are not in a MDP, but rather a BAMDP. As in the previous section, it seems sensible to consider the state space for DAD as the space of histories $h_t = \xi_1 y_1 \dots \xi_t y_t$. Uncertainty over the transition model in DAD is captured by uncertainty in θ . Specifically, we have the following transition distribution for history states

$$p(h_{t+1}|h_t, \xi_{t+1}) = \int p(\theta|h_t) p(y_{t+1}|\xi_{t+1}, \theta) d\theta \quad (34)$$

which is the analogue of equation (32), but now expressed in the notation of experimental design. Unlike the standard reinforcement learning setting, there are no external rewards in DAD. Instead, rewards are defined in terms of information gathered about θ . Specifically, we can take the reward distribution on augmented

states $R^+(r|h_t, a_t, h_{t+1})$ to be a deterministic function of h_{t+1} that represents the information gained about θ by moving from h_t to h_{t+1} . This is given by the reduction in entropy

$$R^+(h_t, a_t, h_{t+1}) = H[p(\theta|h_t)] - H[p(\theta|h_{t+1})]. \quad (35)$$

To complete the BAMDP specification, we take $\gamma = 1$ and we use a time horizon of T . This gives the objective function for policies

$$\mathcal{J}(\pi) = \mathbb{E} \left[\sum_{t=1}^T r_t \right] = \mathbb{E}_{p(\theta)p(h_T|\theta, \pi)} \left[\sum_{t=1}^T H[p(\theta|h_{t-1})] - H[p(\theta|h_t)] \right]. \quad (36)$$

To connect this with the objective that is used in DAD, we apply Theorem 1 of Foster et al. (2021), which tells us that

$$\mathcal{J}(\pi) = \mathbb{E}_{p(\theta)p(h_T|\theta, \pi)} \left[\sum_{t=1}^T H[p(\theta|h_{t-1})] - H[p(\theta|h_t)] \right] \stackrel{\text{Theorem 1}}{=} \mathcal{I}_T(\pi) \quad (37)$$

where

$$\mathcal{I}_T(\pi) = \mathbb{E}_{p(\theta)p(h_T|\theta, \pi)} \left[\log \frac{p(h_T|\theta, \pi)}{\mathbb{E}_{p(\theta')} [p(h_T|\theta', \pi)]} \right]. \quad (38)$$

In summary, we can cast the problem that DAD solves as a BAMDP. We identify designs with actions, experimental histories with augmented states, we use the probabilistic model to give a natural transition distribution on these states, we introduce non-random rewards that are one-step information gains, we set $\gamma = 1$ and generally assume a finite number of experiment iterations T .

3.4 What makes the experimental design problem distinctive?

Having established a theoretical connection between sequential Bayesian experimental design and Bayesian RL, one might naturally ask whether there is any reason to develop specialist algorithms for experimental design when general purpose Bayesian RL algorithms are applicable. First, we focus on the reward structure of the Bayesian experimental design problem. The rewards $r_t = H[p(\theta|h_{t-1})] - H[p(\theta|h_t)]$ are generally intractable, requiring Bayesian inference on θ . Rather than attempting to estimate this reward, DAD proposes the sPCE lower bound on the total expected information gain under policy π , namely

$$\mathcal{I}_T(\pi) \geq \mathcal{L}_T(\pi, L) = \mathbb{E}_{p(\theta_0)p(h_T|\theta_0, \pi)p(\theta_{1:L})} \left[\log \frac{p(h_T|\theta_0, \pi)}{\frac{1}{L+1} \sum_{\ell=0}^L p(h_T|\theta_\ell, \pi)} \right]. \quad (39)$$

Interestingly, there is a way to interpret the sPCE objective within the RL framework. First, we use *root sampling* to sample θ_0 and h_T together. We also fix the contrasts $\theta_{1:L}$. Finally, we use the surrogate rewards

$$\tilde{r}_t = \log \frac{p(h_t|\theta_0, \pi)}{\frac{1}{L+1} \sum_{\ell=0}^L p(h_t|\theta_\ell, \pi)} - \log \frac{p(h_{t-1}|\theta_0, \pi)}{\frac{1}{L+1} \sum_{\ell=0}^L p(h_{t-1}|\theta_\ell, \pi)}. \quad (40)$$

Since these rewards depend on θ_0 , we can treat them as randomised rewards if we are only conditioning on h_t .

One important feature of these rewards is that, whilst intractable, the surrogate $\mathcal{L}_T(\pi, L)$ is differentiable with respect to the designs $(\xi_t)_{t=1}^T$ and observations $(y_t)_{t=1}^T$. In the simplest form of DAD, we further assume a differentiable relationship between y_t and ξ_t that is encapsulated by a reparametrisable way to sample $p(y|\theta, \xi)$. Concretely, for example, we might have $y|\theta, \xi = \mu(\theta, \xi) + \sigma(\theta, \xi)\varepsilon$ where $\varepsilon \sim N(0, 1)$ and μ and σ are differentiable functions. The result of these assumptions is that we can directly differentiate the surrogate objective $\mathcal{L}_T(\pi, L)$ with respect to the parameters ϕ of the policy network π_ϕ that generates the designs $(\xi_t)_{t=1}^T$ according to the formula $\xi_t = \pi_\phi(h_{t-1})$. DAD optimises the policy π_ϕ directly by gradient descent on $\mathcal{L}_T(\pi, L)$.

Thus, DAD can be characterised in RL language as a direct policy optimisation method. Whilst direct policy optimisation methods (Lorberbom et al., 2019; Howell et al., 2021) are used in RL, they are far

from the norm, with methodologies such as Q-learning (Watkins and Dayan, 1992) and actor-critic (Konda and Tsitsiklis, 2000) being more dominant. This may be because RL does not typically assume that the reward function is differentiable—for example, rewards from a real environment rarely come with gradient information. It may also be because discrete action problems are more the focus.

DAD also contrasts with many approaches to *Bayesian* RL in that it avoids the estimation of the posteriors $p(\theta|h_t, \pi)$. In Bayesian RL, these posterior distributions are referred to as *belief states*. Many methods for tackling Bayesian RL problems utilise the estimation of belief states (Ghavamzadeh et al., 2016; Igl et al., 2018; Zintgraf et al., 2019). DAD instead relies on an approach that is closer to the method of root sampling (Guez et al., 2012). This is also one difference between DAD and the previous approach to non-greedy sequential Bayesian experimental design of Huan and Marzouk (2016).

3.5 New objective functions for DAD

Seeing DAD in the framework of Bayesian RL naturally invites the question of whether the general DAD methodology can be applied to objective functions (rewards) that are not information gains. The preceding discussion suggests that, using root sampling so a dependence on θ is possible, we could consider rewards of the form

$$r_t^{\text{general}} = R(\theta, h_t, \epsilon_t) \quad (41)$$

where R is a known differentiable function and ϵ_t is an independent noise random variable. Clearly, the information gain reward r_t fits this pattern, being a function of h_t only. Combining the differentiable reward function with the reparametrisation assumption would mean that the general reward

$$\mathcal{J}^{\text{general}}(\pi) = \mathbb{E}_{p(\theta)p(h_T)p(\epsilon_{1:T})} \left[\sum_{t=1}^T r_t^{\text{general}} \right] \quad (42)$$

can be optimised with respect to π by direct policy gradients. In the experimental design context, this opens the door to two relatively simple extensions of DAD. First, we can assign a (differentiable) cost to each design. Suppose we augment the original expected information gain objective with the negative sum of the costs of the designs. Using λ to trade off cost and information, we arrive at

$$\mathcal{J}^{\text{costed}}(\pi) = \mathcal{I}_T(\pi) - \lambda \mathbb{E} \left[\sum_{t=1}^T C(\xi_t) \right] \quad (43)$$

which we can tackle using an approach that is essentially the same as DAD. Second, we can consider different measures of the quality of the final posterior distribution. For instance, with a one-dimensional θ , we might be more interested in reducing posterior *variance* than posterior entropy. We could take the reward function

$$r_t^{\text{variance}} = \text{Var}_{p(\theta|h_{t-1})}[\theta] - \text{Var}_{p(\theta|h_t)}[\theta]. \quad (44)$$

Whilst there are certain reasons why the entropy approach is considered more theoretically well-justified (Lindley, 1956), using a different functional of the posterior distribution as a reward signal does fit relatively naturally into the DAD framework. The remaining piece of the puzzle would be whether that functional could be estimated efficiently as DAD estimates the information gain using sPCE. For the variance, we have

$$\mathbb{E}_{p(\theta)p(h_T|\theta,\pi)} \left[\sum_{t=1}^T r_t^{\text{variance}} \right] \geq \text{Var}_{p(\theta)}[\theta] - \mathbb{E}_{p(\theta)p(h_T|\theta,\pi)} [(\theta - f_{\phi'}(h_T))^2] \quad (45)$$

where $f_{\phi'}$ is a learnable function. Note the similarity with the Barber–Agakov bound (Barber and Agakov, 2003; Foster et al., 2019, 2020).

3.6 RL algorithms for Bayesian experimental design

To conclude, making the formal connection between sequential Bayesian experimental design opens up the possibility of using the vast literature on Bayesian RL and control theory to improve our ability to plan

sequential experiments. Whilst the direct policy optimisation approach of DAD works remarkably well, understanding the connection to RL should aid us when this training method begins to break down. The application of existing Bayesian RL algorithms to experimental design is an exciting area for new research that is well within reach.

A case of potential difficulty for DAD, where such insights may be useful, is in long-horizon experiments. In order to plan effectively for long experiments, DAD simulates thousands of possible experimental trajectories. However, the efficiency of this simulation is likely to drop as T increases. DAD is extremely data hungry—it resimulates completely new trajectories at each gradient step. This avoids any problems of the training data becoming out-of-date, but it increases the training cost.

It is also conceivable that, in some settings, it is impossible to plan for all future eventualities. The RL analogy would be a strongly stochastic environment in which a game is selected at random from a long list at the start of play. The agent, therefore, has to first discover which game it is playing, and then to play it successfully. If all planning is conducted up-front, then the RL agent has to learn how to play every single game well before starting on the real environment. The alternative is to introduce some real data and retrain the policy as we go. In the RL setting, that would mean discovering which game is being played before knowing how to play the games, which could be achieved with a much simpler policy. Once this discovery is made with good confidence, we can retrain to learn to play that specific game. In the experimental design setting, we are often in the ‘unknown game’ setting. This is because, until we have observed some data, it is almost impossible to know which later experiments will be optimal to run. The DAD approach is to simulate different possibilities and learn to ‘play’ well across the board. The retraining alternative would be a hybrid approach between the standard greedy method and DAD in which some real data is used to retrain the policy as we progress.

References

- David Barber and Felix Agakov. The IM algorithm: a variational approach to information maximization. *Advances in Neural Information Processing Systems*, 16:201–208, 2003.
- Richard Bellman. A markovian decision process. *Journal of mathematics and mechanics*, 6(5):679–684, 1957.
- Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- Nicholas Carlini, Anish Athalye, Nicolas Papernot, Wieland Brendel, Jonas Rauber, Dimitris Tsipras, Ian Goodfellow, Aleksander Madry, and Alexey Kurakin. On evaluating adversarial robustness. *arXiv preprint arXiv:1902.06705*, 2019.
- Daniel R Cavagnaro, Jay I Myung, Mark A Pitt, and Janne V Kujala. Adaptive design optimization: A mutual information-based approach to model discrimination in cognitive science. *Neural computation*, 22(4):887–905, 2010.
- Richard Dearden, Nir Friedman, and Stuart Russell. Bayesian q-learning. In *Aaai/iaai*, pages 761–768, 1998.
- Michael O’Gordon Duff. *Optimal Learning: Computational procedures for Bayes-adaptive Markov decision processes*. University of Massachusetts Amherst, 2002.
- Adam Foster, Martin Jankowiak, Elias Bingham, Paul Horsfall, Yee Whye Teh, Thomas Rainforth, and Noah Goodman. Variational Bayesian Optimal Experimental Design. In *Advances in Neural Information Processing Systems 32*, pages 14036–14047. Curran Associates, Inc., 2019.
- Adam Foster, Martin Jankowiak, Matthew O’Meara, Yee Whye Teh, and Tom Rainforth. A unified stochastic gradient approach to designing bayesian-optimal experiments. In *International Conference on Artificial Intelligence and Statistics*, pages 2959–2969. PMLR, 2020.
- Adam Foster, Desi R Ivanova, Ilyas Malik, and Tom Rainforth. Deep adaptive design: Amortizing sequential bayesian experimental design. *arXiv preprint arXiv:2103.02438*, 2021.
- Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059. PMLR, 2016.
- Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image data. In *International Conference on Machine Learning*, pages 1183–1192. PMLR, 2017.
- Mohammad Ghavamzadeh, Shie Mannor, Joelle Pineau, and Aviv Tamar. Bayesian reinforcement learning: A survey. *arXiv preprint arXiv:1609.04436*, 2016.
- Arthur Guez, David Silver, and Peter Dayan. Efficient bayes-adaptive reinforcement learning using sample-based search. In *Advances in neural information processing systems*, pages 1025–1033, 2012.
- Markus Hainy, David J Price, Olivier Restif, and Christopher Drovandi. Optimal bayesian design for model discrimination via classification. *arXiv preprint arXiv:1809.05301*, 2018.
- Neil Houlsby, Ferenc Huszár, Zoubin Ghahramani, and Máté Lengyel. Bayesian active learning for classification and preference learning. *arXiv preprint arXiv:1112.5745*, 2011.
- Taylor A Howell, Chunjiang Fu, and Zachary Manchester. Direct policy optimization using deterministic sampling and collocation. *IEEE Robotics and Automation Letters*, 6(3):5324–5331, 2021.
- Xun Huan and Youssef M Marzouk. Sequential bayesian optimal experimental design via approximate dynamic programming. *arXiv preprint arXiv:1604.08320*, 2016.
- Maximilian Igl, Luisa Zintgraf, Tuan Anh Le, Frank Wood, and Shimon Whiteson. Deep variational reinforcement learning for pomdps. In *International Conference on Machine Learning*, pages 2117–2126. PMLR, 2018.

- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Vijay R Konda and John N Tsitsiklis. Actor-critic algorithms. In *Advances in neural information processing systems*, pages 1008–1014, 2000.
- Dennis V Lindley. On a measure of the information provided by an experiment. *The Annals of Mathematical Statistics*, pages 986–1005, 1956.
- Guy Lorberbom, Chris J Maddison, Nicolas Heess, Tamir Hazan, and Daniel Tarlow. Direct policy gradients: Direct optimization of policies in discrete action spaces. *arXiv preprint arXiv:1906.06062*, 2019.
- Ruth K Meyer and Christopher J Nachtsheim. The coordinate-exchange algorithm for constructing exact optimal experimental designs. *Technometrics*, 37(1):60–69, 1995.
- Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International conference on machine learning*, pages 1530–1538. PMLR, 2015.
- Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951.
- Stephane Ross, Brahim Chaib-draa, and Joelle Pineau. Bayes-Adaptive POMDPs. In *NIPS*, pages 1225–1232, 2007.
- Joep Vanlier, Christian A Tiemann, Peter AJ Hilbers, and Natal AW van Riel. Optimal experiment design for model selection in biochemical networks. *BMC systems biology*, 8(1):1–16, 2014.
- Benjamin T Vincent and Tom Rainforth. The DARC toolbox: automated, flexible, and efficient delayed and risky choice experiments using bayesian adaptive design. *Retrieved from psyarxiv.com/yehjb*, 2017.
- Christopher JCH Watkins and Peter Dayan. Q-learning. *Machine learning*, 8(3-4):279–292, 1992.
- Christopher K Williams and Carl Edward Rasmussen. *Gaussian processes for machine learning*, volume 2. MIT press Cambridge, MA, 2006.
- Luisa Zintgraf, Kyriacos Shiarlis, Maximilian Igl, Sebastian Schulze, Yarin Gal, Katja Hofmann, and Shimon Whiteson. Varibad: A very good method for bayes-adaptive deep rl via meta-learning. *arXiv preprint arXiv:1910.08348*, 2019.