



DEPARTMENT OF
STATISTICS

Title TBD

A THESIS SUBMITTED IN PARTIAL FULFILLMENT OF THE REQUIREMENTS FOR
transfer of status FROM PRS TO DPHIL

NOVEMBER 2018

ADAM FOSTER

UNIVERSITY COLLEGE

DEPARTMENT OF STATISTICS

UNIVERSITY OF OXFORD

Contents

Acknowledgements	3
Preface	4
1 Introduction and literature review	5
1.1 TODOs	5
1.2 Optimal experiment design	5
1.2.1 Literature review	5
2 Probabilistic modelling and inference	9
3 Probabilistic programming	10
4 Optimal experiment design	11
5 Future directions	12
5.1 EIG	12
5.1.1 EIG estimation on simpler models	12
5.1.2 EIG estimation on complex models	12
5.1.3 Theory of EIG estimators	12
5.1.4 EIG gradients	12
5.1.5 EIG optimisation	12
5.1.6 Model misspecification	13
5.1.7 Sequential design and active learning	13
5.1.8 Optional stopping	13
5.1.9 Dynamic models	13
5.2 Beyond EIG	13
5.2.1 Causal inference	13
5.2.2 Power	13
5.2.3 Cost	13
5.2.4 Non-greedy	13
5.2.5 Other criteria	14

5.2.6	Experiment design for model criticism	14
	Bibliography	15

Acknowledgements

I have been fortunate, in the first year of my DPhil, to have worked with a number of brilliant people. First and foremost, I would like to acknowledge the support, advice and patience of my supervisor, Yee Whye Teh. Also within Oxford, both Benjamin Bloem-Reddy and Tom Rainforth have been incredibly generous with their time and expertise. Without them, I would have struggled to achieve much in my first year. I would also like to thank Emile Mathieu for his unwavering support and friendship. I am grateful to Noah Goodman for taking a gamble on an unknown DPhil student and supporting my application to Uber Technologies for an internship. Noah's enthusiasm and guidance made my internship a joy, and also set me on a new and exciting research path. I'd like to extend my thanks to Martin Jankowiak and Eli Bingham, as well as the whole pyro team, for their continued help and support.

Preface

This thesis aims to describe the work I have done so far in my DPhil and discuss the future directions I plan to take in my research. Of primary relevance to my current work and future plans is Chapter 4 on optimal experiment design. Chapter 2 mostly concerns a project I worked on with Benjamin Bloem-Reddy which resulted in a UAI paper and oral presentation [Bloem-Reddy et al., 2018]. Chapter 3 relates to a question posed to me by Yee Whye Teh, namely, ‘is probabilistic programming useful for Bayesian nonparametrics?’. Our workshop paper [Bloem-Reddy et al., 2017] and my open source contributions to the language `pyro` informed this chapter. Chapter 4, the heart of this thesis, constitutes a draft of a paper that I plan to submit to ICML 2019, and represents the culmination of my internship with Uber. As detailed in Chapter 5, it is this project that I have found most exciting in my DPhil so far and that my future work will broadly be an extension of.

Chapter 1

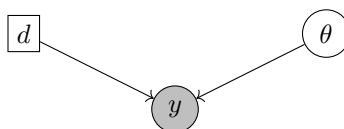
Introduction and literature review

1.1 TODOs

- Probabilistic modelling is important.
- Applied stats models
- Bayesian nonparametric models
- Probabilistic programming helps to do automation

1.2 Optimal experiment design

Consider the following graphical model



in which d represents a (non-random) control variate, the design of the experiment, θ represents a latent variable and y represents the observed outcome of the experiment.

1. introduce the problem 2. discuss EIG. Where does EIG come from? Interpretations of EIG 3. literature review– who uses OED, how have they approached the problem? 4. mutual information lit review

1.2.1 Literature review

TODO: Major revisions needed

Machine learning and statistics

There is long-standing interest in ‘classical’ statistical models and their design [Youssef,]. Consider a basic linear models with Gaussian noise. Optimal design here can be expressed in terms of the eigen-spectrum of XX^T (see [Chaloner, 1984]). For nonlinear models, see the section on Physics. What about

GLMs? Likely can solve the problem analytically again. These are great baselines. People in the linear models case are often concerned with proving the equivalence of different kinds of optimality [Youssef,].

In machine learning, experiment design is closely related to two common techniques in, for example, image classification: data augmentation and active learning.

In data augmentation images are rotated, translated, etc to create more training data. We could theoretically optimise the augmentation but this seems wasteful since copying the labels to new images is very easy.

A much more interesting area is *active learning*. In this context, there are a large number of unlabeled images. Labeling is expensive. We select which images to label either up front, or (more typical in active learning) in a sequential manner. The key difference here is that we have a finite pool of unlabeled instances. We may be more interested in reducing uncertainty in the labels of these unlabeled images than in our posterior entropy.

The connection between active learning and Bayesian optimal design was explored in [Golovin et al., 2010]. In this paper, they start from a place where the outcome of a test is deterministic (think of the 12 men on an island problem). In the noiseless setting, the sequential design can be encoded as a decision tree and the problem is called the Optimal Decision Tree problem. This problem is known to be NP-hard. The OED criterion is introduced later to account for noisy observations and the fact that true parameters need not be known exactly even after all tests have been run.

A particular active learning example can be found in [Nowak, 2009]. We have \mathcal{H} a hypothesis space (read parameter space) and \mathcal{X} a query space (read design space). The goal is to determine the true $h^* \in \mathcal{H}$. Each query outputs a label in $\{-1, 1\}$ corrupted with Bernoulli noise (independent between queries). The algorithm broadly works by targeting $x \in \mathcal{X}$ where the expected posterior label is near 0 (random guess). The convergence rate of $\mathbb{P}(\hat{h}_i \neq h^*) \rightarrow 0$ is studied (shown exponential). The importance of having access to unlabeled data is exploited by [Dasgupta, 2006].

Psychology

For an overview of optimal experiment design in probabilistic programming, [Ouyang et al., 2016] from Noah’s group is a good place to start. Experiment design is necessary to distinguish competing theories. We should select models with the highest *expected information gain*, written formally as

$$U(d) = \mathbb{E}_{(Y, \Theta) \sim p(y, \theta|d)} \left\{ \log \frac{p(\Theta|Y, d)}{p(\Theta)} \right\} \quad (1.1)$$

This equation has been studied by mathematical statisticians since the 50’s [Lindley, 1956]. We can naively evaluate $U(d)$ in a PPL via *nested inference*.

A canonical experiment discussed in this paper is the 5-4 experiment for category learning [Medin and Schaffer, 1978]. The experiment aimed to distinguish two competing models of category learning: the *exemplar model* (learn categories by comparing new items to all previous items) and the *prototype model* (learn categories by remembering a prototypical example). There are two models, so $\Theta = \{m_1, m_2\}$. During the experiment, participants are presented with a sequence of objects. In the training phase, they are also told the

correct label after guessing. In testing they have no feedback. The objects varied in four dimensions: colour, shape, size and count; we can consider the space of objects to be $\{0, 1\}^4$. Each object has a label A or B . The true labeling mechanism was limited by Medin and Schaffer to be linearly separable. There are 9 inputs in the final testing set and Medin and Schaffer restricted there to be 5 A s and 4 B s. The objects 0000 and 1111 have to be present. Under these restrictions there are 933 possible experiments up to permutation. So \mathcal{D} is a finite set of size 933. \mathcal{Y} is the ‘test’ responses, ie. the subjects responses when they are not given feedback. Thus $\mathcal{Y} = \{A, B\}^9$. The kind of participant numbers seen were 10-30.

In [Vincent and Rainforth, 2017], the canonical experiment is as follows. We want to model how humans discount future rewards relative to present ones (via utility indifference pricing framework). A single experiment takes the following form: ‘Would you prefer $\pounds A_1$ at time t_1 or $\pounds A_2$ at time t_2 ’? The parameter of interest is the discount factor. Formally, $\mathcal{D} = [0, \infty)^2$, $\mathcal{Y} = \{1, 2\}$ and $\Theta = [0, \infty)$. These three spaces fit together as follows. We first chose \mathcal{D} the space of possible designs. We subsequently chose \mathcal{Y} the space of possible outcomes. We posited a probabilistic model for Y in terms of parameters θ . Focus on non-nested estimation for finite \mathcal{Y} and sequential design. Sequential design means different participants will be asked different questions based on their previous answers.

Bioinformatics

In [Vanlier et al., 2012], the authors consider experiment design from the perspective that, with little data, many different parameter settings adequately describe the data. Canonical model. Biochemical network modeled as an ODE.

$$\begin{aligned}\dot{x} &= f(x, u, p) \\ \dot{y} &= g(x, q) + \xi \\ x(0) &= x_0\end{aligned}$$

u is the input, x, y are time varying (uncontrolled) with x latent and y observed, p, q, x_0 are parameters θ required to simulate the model and do not depend on t , ξ represents measurement noise. We treat ξ as iid Gaussian. The paradigm chosen here is expected variance reduction, as opposed to information gain. (Possibly wrong if people still do that.) The variance is in the posterior predictive density.

Physics

In [van Den Berg et al., 2003], we begin by discussing ‘classical’ experiment design procedures which assume linear dependence between model and outcome $y = G_{m_0}d$. One can solve this linear equation by least squares, $\hat{d} = G^T(GG^T)^{-1}y$. Define $L = G^T(GG^T)^{-1}$, possibly adding regularization as necessary. Basically, you want to maximize the max eigenvalue of G , which is essentially a gradient. The larger gradient, the more informative the experiment. In a linear setting, the gradient does not depend on the true parameter value.

Now consider linear noise but a nonlinear function between parameters and outcomes. For example,

the authors took

$$R_p = \left(\frac{1}{2} [1 + \tan^{-1} i] - 4c^2 \sin^2 i \right) \frac{\Delta\alpha}{\alpha} \quad (1.2)$$

where $\alpha = (\alpha_1 + \alpha_2)/2$, $\Delta\alpha = (\alpha_2 - \alpha_1)$. The parameter we want to optimize is α_2 .

This is a relatively simple and comprehensible case.

Chapter 2

Probabilistic modelling and inference

Chapter 3

Probabilistic programming

Chapter 4

Optimal experiment design

Chapter 5

Future directions

5.1 EIG

5.1.1 EIG estimation on simpler models

Current project. Focus is on linear-type models (see Kruschke) that are used in applied stats. We can implement some semi-whitebox methods here. The aim is to use these models and EIG estimators in active learning loops. So we want sub-second estimation. This leads to methods based on relatively structured guides.

5.1.2 EIG estimation on complex models

Very related, but taking a more black box approach. Assume that the model is too complex to build a structured guide, but that the experiment is very expensive. So we can spend more time on EIG estimation. Deep learning approaches, like Donsker-Varadhan, might look more attractive.

5.1.3 Theory of EIG estimators

Are estimators statistically consistent? Can we estimate, bound or approximate the error, or the relative error across different d ?

5.1.4 EIG gradients

How best to estimate the gradient $\partial_d \text{EIG}$? Can we obtain bounds? What would Rainforth gradient estimation look like? Can we optimise EIG in a GAN-like fashion – iterative updates of q and d .

5.1.5 EIG optimisation

Are there special features of EIG that we can exploit when using Bayes opt, or something else, to do EIG optimisation?

5.1.6 Model misspecification

How best to deal with model misspecification in experiment design. A uniform increase in y entropy does not change design... what would be the right paradigm for this?

5.1.7 Sequential design and active learning

Further considerations for using EIG estimation/optimisation in a live active learning loop.

5.1.8 Optional stopping

Suppose we use posterior entropy as an optional stopping criterion, and use EIG for sequential experiment design. How would this impact final conclusions that we are able to make about data?

5.1.9 Dynamic models

Experiment design for systems that change as a result of the experimentation. Things like the atmosphere or a pond.

5.2 Beyond EIG

5.2.1 Causal inference

What if we design an experiment for causal structure learning? *And* information? How do these fields intersect? Speak to Robin Evans.

5.2.2 Power

This is a theoretical question. How does the Bayesian notion of EIG intersect with frequentist notions of experiment design, in particular, statistical power?

5.2.3 Cost

Designing experiments for information, but with a cost associated with each experiment. Sequential case may be more interesting than one shot (which seems simple).

5.2.4 Non-greedy

Related to above. Solving the non-greedy experimental design problem brings in elements from POMDPs and RL. Should we use EIG here? Should we use RL reward functions? Are they in some sense (approximately) the same? Could greedy EIG optimisation arise as a good approximation to the RL task?

5.2.5 Other criteria

In active learning, they have criteria about the expected misclassification, and some other criteria. Can we connect these? In classical experiment design they have all these mysterious criteria like D -optimality and so on.

5.2.6 Experiment design for model criticism

Rather than assuming the model to be true and looking to gain information within the model, suppose instead that we have an empirical distribution and seek a new experiment to best expose flaws in the whole model. For instance, when comparing the posterior predictive and empirical distributions (possibly conditional on an input).

Bibliography

- [Bloem-Reddy et al., 2018] Bloem-Reddy, B., Foster, A., Mathieu, E., and Teh, Y. W. (2018). Sampling and inference for beta neutral-to-the-left models of sparse networks. In *Uncertainty in Artificial Intelligence*.
- [Bloem-Reddy et al., 2017] Bloem-Reddy, B., Mathieu, E., Foster, A., Rainforth, T., Teh, Y. W., Ge, H., Lomelí, M., and Ghahramani, Z. (2017). Sampling and inference for discrete random probability measures in probabilistic programs. In *NIPS Workshop on Advances in Approximate Bayesian Inference*.
- [Chaloner, 1984] Chaloner, K. (1984). Optimal bayesian experimental design for linear models. *The Annals of Statistics*, pages 283–300.
- [Dasgupta, 2006] Dasgupta, S. (2006). Coarse sample complexity bounds for active learning. In *Advances in neural information processing systems*, pages 235–242.
- [Golovin et al., 2010] Golovin, D., Krause, A., and Ray, D. (2010). Near-optimal bayesian active learning with noisy observations. In *Advances in Neural Information Processing Systems*, pages 766–774.
- [Lindley, 1956] Lindley, D. V. (1956). On a measure of the information provided by an experiment. *The Annals of Mathematical Statistics*, pages 986–1005.
- [Medin and Schaffer, 1978] Medin, D. L. and Schaffer, M. M. (1978). Context theory of classification learning. *Psychological review*, 85(3):207.
- [Nowak, 2009] Nowak, R. (2009). Noisy generalized binary search. In *Advances in neural information processing systems*, pages 1366–1374.
- [Ouyang et al., 2016] Ouyang, L., Tessler, M. H., Ly, D., and Goodman, N. (2016). Practical optimal experiment design with probabilistic programs. *arXiv preprint arXiv:1608.05046*.
- [van Den Berg et al., 2003] van Den Berg, J., Curtis, A., and Trampert, J. (2003). Optimal nonlinear bayesian experimental design: an application to amplitude versus offset experiments. *Geophysical Journal International*, 155(2):411–421.
- [Vanlier et al., 2012] Vanlier, J., Tiemann, C. A., Hilbers, P. A., and van Riel, N. A. (2012). A bayesian approach to targeted experiment design. *Bioinformatics*, 28(8):1136–1142.

[Vincent and Rainforth, 2017] Vincent, B. T. and Rainforth, T. (2017). The darc toolbox: automated, flexible, and efficient delayed and risky choice experiments using bayesian adaptive design.

[Youssef,] Youssef, N. A. A review on optimal experimental design.