
On Exploration, Exploitation and Learning in Adaptive Importance Sampling

Xiaoyu Lu¹ Tom Rainforth¹ Yuan Zhou¹ Yee Whye Teh¹ Frank Wood¹
Hongseok Yang² Jan-Willem van de Meent³

¹University of Oxford; ²KAIST, South Korea; ³Northeastern University
{xiaoyu.lu, rainforth, y.w.teh}@stats.ox.ac.uk, yuan.zhou@cs.ox.ac.uk,
fwood@robots.ox.ac.uk, hongseok.yang@kaist.ac.kr,
j.vandemeent@northeastern.edu

1 Introduction

Monte Carlo methods form the bedrock upon which significant sections of probabilistic machine learning and computational statistics rest. An important Monte Carlo technique which forms the basis for many others is importance sampling (IS). Let $\pi(x) = f(x)/Z$ be a target density which can be evaluated pointwise up to an unknown normalising constant Z , and let $q(x)$ be a proposal distribution from which samples can be drawn and which can be evaluated pointwise. IS works by drawing a sequence of samples x_1, x_2, \dots from the proposal $q(x)$, and using these to estimate both Z and target statistics $\mathbb{E}_\pi[\phi(x)]$ for some test function $\phi(x)$. Let $w(x_t) = f(x_t)/q(x_t)$ be the importance weight of x_t . Then,

$$Z = \int f(x)dx = \mathbb{E}_q[w(X)] \approx \frac{1}{T} \sum_{t=1}^T w(x_t),$$
$$\mathbb{E}_\pi[\phi(x)] = \frac{\int f(x)\phi(x)dx}{\int f(x)dx} = \frac{\mathbb{E}_q[w(X)\phi(X)]}{\mathbb{E}_q[w(X)]} \approx \frac{\sum_{t=1}^T w(x_t)\phi(x_t)}{\sum_{t=1}^T w(x_t)}. \quad (1)$$

Note that the estimate for Z is unbiased, but that for the target statistics is biased but consistent.

The efficiency of IS is governed by the choice of proposal $q(x)$, with the intuition that the closer q is to π the better. Adaptive IS (AdaIS) techniques [1, 2, 3, 4] attempt to improve the efficiency of IS by adapting the proposal to be closer to the target, producing a sequence of proposals q_1, q_2, \dots . The IS estimates (1) still apply with $q(x_t)$ replaced by $q_t(x_t)$. The basic idea is that previous samples x_s along with evaluations $f(x_s)$ give information about the distribution of probability masses in the target, and the proposal should place more mass where the target has more mass.

Viewed in this way, AdaIS is in effect an online learning problem, that of learning the target density π through a sequence of queries of it. As opposed to the typical setup of density estimation where each query is an iid sample from the target, here each query involves drawing a sample x_t from the current proposal q_t , and evaluating $f(x_t)$, the target density up to an unknown normalising constant. At each iteration, the proposal q_t is both our current estimate of the target, as well as our tool for querying the target.

Note that this exposes a trade off between exploration and exploitation. We would like our proposal q_t to be as close as possible to the target, so that our IS estimate is as good as possible (exploit). At the same time, q_t directs where queries of the target are made, and probability mass needs to be spread over the sample space where we have high uncertainty of the target so that we may query and reduce our uncertainty to improve our estimate of the target for the future (explore).

2 A Bandits Approach to Adaptive Important Sampling

In this paper we take the first steps towards developing an AdaIS method which optimally trades off exploration versus exploitation. Methods which address the exploration-exploitation trade off has been well-studied in online learning, most successfully under the banner of bandit algorithms [5],

32 with the Upper Confidence Bound (UCB) methods [6] and Thompson sampling methods [7] being
 33 popular approaches. In (basic) UCB, a finite number of arms are present, and at each iteration an arm
 34 is chosen to be pulled, which then returns a random reward. The aim to maximise rewards in the long
 35 run, by finding the arm with highest average reward. UCB operates by maintaining an estimate of the
 36 average reward for each arm, and arms are picked according to the estimates plus optimism boosts
 37 which are larger for arms where our estimates are less certain to encourage exploration. [6] showed
 38 that UCB optimally trades-off exploration and exploitation by showing that the cumulative regret
 39 (relative to an oracle which knows which arm is optimal) grows logarithmically in the number of
 40 iterations, which is the best growth rate achievable [8].

41 In this section we will develop an AdaIS method which has a similar flavour to UCB. We assume
 42 that we have a partition of our sample space into K disjoint subsets (corresponding to bandit arms),
 43 and that we have a tractable base distribution $g_a(x)$ for each subset indexed by $a \in \{1, \dots, K\}$. We
 44 consider proposal distributions of the form

$$q_t(x) = \sum_{a=1}^K q_{at} g_a(x) \quad (2)$$

45 where $\sum_a q_{at} = 1$ and the probability masses q_{at} of the subsets are to be adapted in the scheme. In
 46 the next section we will extend the approach to a hierarchical partition of the sample space, where
 47 subsets are recursively split where necessary.

48 A measure of the (in)efficiency of the proposals q_t is also required so that we can define what makes
 49 an optimal proposal among the class above, and the regret of using a proposal relative to the optimal.
 50 We will use the KL divergence $\text{KL}(\pi||q_t)$ as this measure, which has been shown by Chatterjee and
 51 Diaconis [9] to be the correct measure of the inefficiency of an IS with proposal q_t . For proposals as
 52 given by (2), this is,

$$\text{KL}(\pi||q_t) = \sum_a \int_{\mathcal{X}_a} \pi(x) \log \frac{\pi(x)}{g_a(x)} dx - \sum_a \pi_a \log q_{at} \quad (3)$$

53 where $\pi_a = \int_{\mathcal{X}_a} \pi(x) dx = Z_a/Z$ is the target probability mass of subset \mathcal{X}_a and Z_a is the corre-
 54 sponding subset partition function. The optimal parameters q_{at} are seen to be $q_{at}^* = \pi_a$, and the regret
 55 $R(q_t)$ of proposal q_t is,

$$R(q_t) = \text{KL}(\pi||q_t) - \text{KL}(\pi||q^*) = \sum_a \pi_a \log \frac{\pi_a}{q_{at}}, \quad (4)$$

56 which is just the KL divergence between the two finite vectors of probability masses.

57 At iteration t , a query of the target starts with a sample x_t from q_t , which involves first picking a
 58 subset A_t according to the probability masses q_{at} , then drawing a sample $x_t \in \mathcal{X}_{A_t}$ from $g_{A_t}(x)$.
 59 This then informs our knowledge of the target probability mass of subset \mathcal{X}_{A_t} . For each subset a we
 60 have $\pi_a = Z_a / \sum_b Z_b$ where

$$Z_a = \int_{\mathcal{X}_a} f(x) dx = \mathbb{E}_{g_a} \left[\frac{f(x)}{g_a(x)} \right]. \quad (5)$$

61 And so, if $A_t = a$, the importance weight $Y_{at} := f(x_t)/g_a(x_t)$ of the proposal $g_a(x)$ for subset a
 62 gives an unbiased estimate of Z_a .

63 Naively, for the next iteration we can now estimate each Z_a using

$$\hat{Z}_{a,t+1} := \frac{\sum_{l \leq t: A_l = a} Y_{al}}{N_{a,t+1}} \quad (6)$$

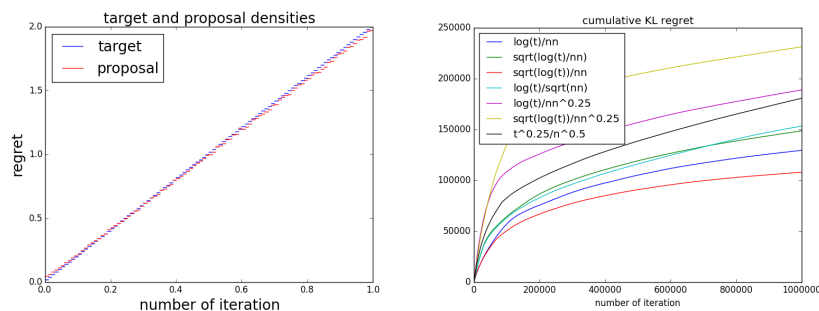
64 where $N_{a,t+1} = \#\{l : l \leq t, A_l = a\}$ is the number of times subset a has been chosen up to time
 65 t , and use the estimate $q_{a,t+1} = \hat{Z}_{a,t+1} / \sum_b \hat{Z}_{b,t+1}$. A problem with this naive scheme is that if
 66 by chance our estimate \hat{Z}_{at} for some subset a is too small, the resulting low estimated proposal
 67 probability will result in low probability for the subset to be picked in future, and hence the bad
 68 estimate may not be corrected. This is a symptom of under-exploration. As in UCB, we will consider
 69 encouraging exploration by using an optimism boost:

$$q_{a,t+1} = \frac{\hat{Z}_{a,t+1} + \sigma_{a,t+1}}{\sum_{a=1}^K (\hat{Z}_{a,t+1} + \sigma_{a,t+1})} \quad (7)$$

70 where σ_{at} should be decreasing with N_{at} but grows with t . The intuition is that if we have not
 71 explored the subset a sufficiently, σ_{at} is relatively large, which compensates and boosts q_{at} , allowing
 72 us to have higher chance to explore subset a and correct the under estimate. The growth with t is to
 73 ensure sufficient exploration of all subsets over time.

74 Note that the estimate $\hat{Z}_{a,t+1}$ above is in fact biased, since $N_{a,t+1}$ is random and correlated with
 75 previous values of Y_{al} for $l \leq t^1$. If it were unbiased, it is possible to bound the regret (4) by using
 76 concentration inequalities to bound the probability of large deviations. It is our ongoing work to fix
 77 this.

78 We demonstrate our method on a simple problem, and empirically evaluate a number forms of
 79 the optimism boost. Our sample space is the unit interval $[0, 1]$, and we partitioned it evenly into
 80 100 subintervals. We picked the target density $\pi(x)$ and the subproposals $g_a(x)$ such that Y_{at} has
 81 distribution $2a\text{Bernoulli}(1/100)/101$ for interval $a = \{1, \dots, 100\}$. We made this choice so that
 82 Y_{at} has a large but controllable variance, so that the resulting problem of adaptation is hard enough
 83 for exploration to be important. Figure 1a shows that our algorithm is able to adapt the proposals
 84 so that they converge to the optimum by balancing between exploitation and exploration well. Note
 85 that low probability intervals have their mass over-estimated by the proposal due to exploration. We
 86 compare the cumulative regret for different forms of optimism boosts σ_{at} in Figure 1b. For each
 87 form we have optimised over a constant multiplier to minimise cumulative regret. We observe that
 88 optimism boosts with an inverse relationship with N_{at} and slow growth with t work well, and seem
 89 to achieve sublinear cumulative regret.



(a) Target and proposal probabilities at final iteration, for optimism boost $\log(t)/\sqrt{N_{at}}$. (b) Cumulative regrets as functions of iteration.

Figure 1: Results for our AdaIS method, with 100 subsets, averaged over 10 runs.

90 3 Hierarchical Partition

91 Instead of fixing the number of partitions of the sample space to be K , we extend the approach
 92 to a binary hierarchical partition of the sample space using the *Polyá tree* idea from Bayesian
 93 nonparametrics [10][11] and the UCT algorithm[12]. The idea is that the subsets are recursively split
 94 when the relative masses are high in the subset, enabling a finer partitioning over the sample space
 95 which has higher density. We now explain in detail of how the algorithm works.

96 For simplicity, consider the sample space \mathcal{X} to be the interval $[0, 1]$, the idea is to conduct a hierarchical
 97 binary partition where the interval is recursively evenly split into left or right subspace, inducing a
 98 tree structure. At each iteration, sampling from the root \mathcal{X} down the tree according to the probability
 99 of going left or right at each node, which are computed using samples already drawn, until a leaf
 100 (a subspace) is visited. Draw a new sample in that subspace with a default proposal (e.g. uniform),
 101 and keep track of whether the new sample falls into the left or the right of the leaf node, as well as
 102 its importance weights. Relevant statistics are then updated along the path to the root to update the
 103 sequence of probabilities of going left/right. Now one can decide whether to further partition the
 104 leaf node according to some splitting criterion. For example, when the tree reaches some truncation
 105 level, *i.e.*, the depth of the tree is fixed; or more intuitively, we partition the leaf node further when it
 106 reduces the KL -divergence without compromising too much computational costs[9].

¹We thank Tom Rainforth for pointing this out to us.

107 We demonstrate the algorithm on a simple example, consider the target density to be $\pi(x) \propto$
 108 $\exp(10(x - 1))\mathbf{1}_{x \in (0,1)}$, and we set the truncation level to be 10. The result can be seen in figure
 109 2. It can be seen that more intervals have been partitioned over higher mass regions, and the
 110 algorithm is able to adapt the proposal distributions according to the relative masses it has already
 111 obtained. We also applied the algorithm to the classic 2D banana shaped problem, with density
 112 $f(x_1, x_2) \propto \exp\{-0.5(0.03x_1^2 + (x_2 + 0.03(x_1^2 - 100))^2)\}$, where we used the *KL*-divergence
 113 splitting criterion. It can be seen from figure 3 that we successfully recover the density within
 114 100,000 samples, and the the region with low density does not get splitted whereas the region with
 115 high density has finer grid for more accurate estimates. Moreover, the region with high density stops
 116 splitting infinitely to balance with the computational expense.

117 4 Conclusion

118 In this work, we have addressed the issue of exploration-exploitation in adaptive importance sampling,
 119 and proposed a novel approach through the lens of multi-armed bandit problems, borrowing the ideas
 120 of upper confidence bounds. We extend our method to the hierarchical case, where the sample space
 121 is recursively split in high density regions, and demonstrated experimentally that our method gives
 122 promising performance with little computational costs. In ongoing work we are investigating unbiased
 123 estimates of Z_a and a finite-time theoretical analysis to understand the growth of the cumulative
 124 regret. So far as we know, the concept of trading off exploration versus exploitation has not been
 125 widely considered in adaptive Monte Carlo, except for [13][14]. [13] considers a number of Monte
 126 Carlo estimates, where the idea is to pick the one with lowest variance with the goal of picking the
 127 optimal estimate as much as possible. [14] considers the problem of stratified sampling, where each
 128 strata is viewed as an arm and the mean in each strata is estimated. Both cases are closer to the bandit
 129 problem setting whereas our importance sampler does not choose the arm but learns a distribution
 130 over arms instead.

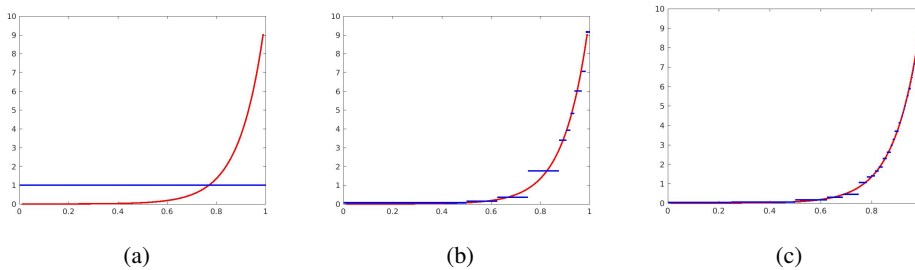


Figure 2: Results using the hierarchical algorithm: target and proposal density, the target density are plotted in red, and the adaptive proposal probabilities are plotted in blue, at iteration 1,(a); 50,(b); and 200,(c) respectively.

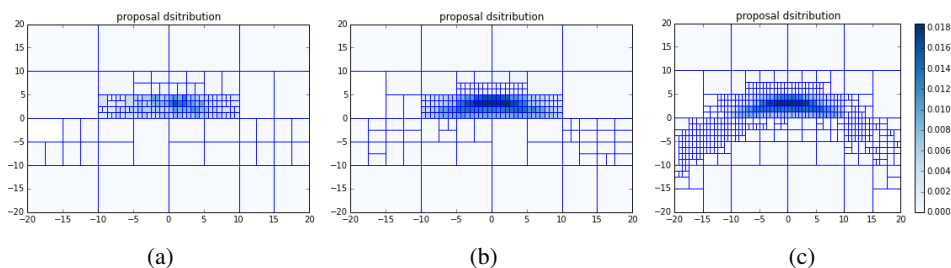


Figure 3: Banana shaped example, learned proposal distributions and partitions at iteration 1000,(a); 10000,(b); 100000,(c) respectively. Darker color indicated high density.

131 **References**

- 132 [1] Jun S Liu. *Monte Carlo strategies in scientific computing*. Springer Science & Business Media,
133 2008.
- 134 [2] Olivier Cappé, Arnaud Guillin, Jean-Michel Marin, and Christian P Robert. Population monte
135 carlo. *Journal of Computational and Graphical Statistics*, 13(4):907–929, 2004.
- 136 [3] Jean Cornuet, JEAN-MICHEL MARIN, Antonietta Mira, and Christian P Robert. Adaptive
137 multiple importance sampling. *Scandinavian Journal of Statistics*, 39(4):798–812, 2012.
- 138 [4] Olivier Cappé, Randal Douc, Arnaud Guillin, Jean-Michel Marin, and Christian P Robert.
139 Adaptive importance sampling in general mixture classes. *Statistics and Computing*, 18(4):447–
140 459, 2008.
- 141 [5] Donald A Berry and Bert Fristedt. *Bandit problems: sequential allocation of experiments*
142 (*Monographs on statistics and applied probability*). Springer, 1985.
- 143 [6] Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed
144 bandit problem. *Machine learning*, 47(2-3):235–256, 2002.
- 145 [7] Shipra Agrawal and Navin Goyal. Analysis of thompson sampling for the multi-armed bandit
146 problem. In *COLT*, pages 39–1, 2012.
- 147 [8] Tze Leung Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Ad-*
148 *vances in applied mathematics*, 6(1):4–22, 1985.
- 149 [9] Sourav Chatterjee and Persi Diaconis. The sample size required in importance sampling. *arXiv*
150 *preprint arXiv:1511.01437*, 2015.
- 151 [10] R Daniel Mauldin, William D Sudderth, and SC Williams. Polya trees and random distributions.
152 *The Annals of Statistics*, pages 1203–1221, 1992.
- 153 [11] R Daniel Mauldin and SC Williams. Reinforced random walks and random distributions.
154 *Proceedings of the American Mathematical Society*, 110(1):251–258, 1990.
- 155 [12] Levente Kocsis and Csaba Szepesvári. Bandit based monte-carlo planning. In *ECML*, volume 6,
156 pages 282–293. Springer, 2006.
- 157 [13] James Neufeld, András György, Dale Schuurmans, and Csaba Szepesvári. Adaptive monte
158 carlo via bandit allocation. *arXiv preprint arXiv:1405.3318*, 2014.
- 159 [14] Alexandra Carpentier and Rémi Munos. Finite time analysis of stratified sampling for monte
160 carlo. In *Advances in Neural Information Processing Systems*, pages 1278–1286, 2011.