

# DEEP BAYESIAN BANDITS SHOWDOWN

AN EMPIRICAL COMPARISON OF BAYESIAN DEEP NETWORKS FOR THOMPSON SAMPLING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

*Recent interest in decision making with deep neural networks, e.g. for reinforcement learning, has led to a significant body of work developing practical methods for trading off exploration and exploitation. Bayesian methods for deep learning are especially appealing for this purpose as their quantification of epistemic uncertainty helps inform strategies for exploration. However, these methods are rarely compared on benchmarks that evaluate their uncertainty, particularly for decision making, and their practical trade-offs seem poorly understood. In this paper, we compare a variety of well-established and recent methods under the lens of Thompson Sampling over a series of contextual bandit problems.*

## 1 INTRODUCTION

Recent advances in areas such as reinforcement learning have sparked renewed interest in sequential decision making with deep neural networks. These models have proven to be powerful and flexible function approximators, allowing one to learn mappings directly from complex state, e.g. pixels, to estimates of expected reward or value. While such models can be tremendously accurate on data they have been trained on, quantifying epistemic uncertainty remains challenging. However, having an understanding of what is not yet known or well understood is critical to some central tasks of machine learning, such as effective exploration for decision making.

Unfortunately, the exact Bayesian posterior is generally intractable for such highly parameterized models. As such, significant effort has been dedicated to approximate Bayesian methods for deep neural networks. These range from variational methods (Graves, 2011; Blundell et al., 2015; Kingma et al., 2015) to stochastic minibatch Markov Chain Monte Carlo (Neal, 1994; Welling & Teh, 2011; Li et al., 2013; Ahn et al., 2012; Mandt et al., 2016). Because the exact posterior is intractable, evaluating these approaches is challenging. Furthermore, these methods are rarely compared on benchmarks that measure the quality of their estimates of uncertainty for downstream tasks.

To address this challenge, we develop a benchmark for methods for exploration using deep neural networks. We compare a variety of well-established and recent methods under the lens of Thompson sampling for contextual bandits, a classical task in sequential decision making. We will open-source our code and implementations to help researchers develop and compare new approaches.

A fundamental challenge in sequential decision making is the exploration-exploitation dilemma: in order to maximize cumulative reward, agents need to trade-off between what is expected to be best at the moment, (i.e., exploitation), with potentially sub-optimal exploratory actions. Solving this trade-off in an efficient manner, i.e. to maximize cumulative reward, is a significant challenge as it requires epistemic uncertainty, a quantification of what is not yet known or understood. Furthermore, exploratory actions should be coordinated throughout the entire decision making process, i.e. deep exploration, rather than performed myopically at each state.

Exploration in the context of reinforcement learning is a highly active area of research. Simple strategies such as epsilon-greedy remain extremely competitive (Mnih et al., 2015; Schaul et al., 2016). However, recently a number of promising techniques have emerged involving encouraging exploration through carefully adding random noise to the parameters (Plappert et al., 2017; Meire Fortunato, 2017; Gal & Ghahramani, 2016) or bootstrap sampling (Osband et al., 2016) before making decisions. These methods rely explicitly or implicitly on Posterior or Thompson sampling for exploration.

Thompson Sampling is a classic algorithm (Thompson, 1933) which requires only that one can sample from the posterior distribution over plausible problem instances (for example, values or rewards). At

each round, it draws a sample and takes a greedy action under that sample. The posterior distribution is then updated after the result of the action is observed. Thompson sampling has been shown to be extremely effective in practice for bandit problems (Chapelle & Li, 2011; Granmo, 2010) and has been shown to exhibit appealing theoretical guarantees of efficiency (Agrawal & Goyal, 2012). It is especially appealing for deep neural networks as one rarely has access to the full posterior but can often approximately sample from it.

In this paper, we investigate how different posterior approximations affect the performance of Posterior Sampling from an empirical standpoint. For simplicity, we restrict ourselves to one of the most basic sequential decision making scenarios: that of contextual bandits. In this case, Posterior Sampling is usually known as Thompson Sampling.

As expected, no single algorithm beats all the others in every bandit problem. We observe, however, some general trends. We find that simple heuristics, including bootstrapping, dropout, and injecting random noise, do not significantly and systematically improve performance, while they may still sometimes help. This suggests that the intrinsic randomness of stochastic gradient descent is enough to explore in many cases. Other algorithms, like Variational Inference and Monte Carlo approaches, strongly couple their complex representation and uncertainty estimates. This proves problematic when decisions are made based on partial optimization of both, as in online scenarios. On the other hand, making decisions according to a Bayesian linear regression on the representation provided by the last layer of a deep network offers a robust and easy-to-tune approach.

In Section 2 we discuss Thompson Sampling, and present the contextual bandit problem. The different algorithmic principles to approximate the posterior distribution fed to Thompson Sampling are introduced in Section 3, while the linear case is described in Section 4. The main experimental results are presented in Section 5, and discussed in Section 6. Finally, Section 7 concludes.

## 2 DECISION-MAKING VIA THOMPSON SAMPLING

The contextual bandit problem works as follows. At time  $t = 1, \dots, n$  a new context  $X_t \in \mathbf{R}^d$  arrives and is presented to algorithm  $\mathcal{A}$ . The algorithm—based on its internal model and  $X_t$ —selects one of the  $k$  available actions,  $a_t$ . Some reward  $r_t = r_t(X_t, a_t)$  is then generated and returned to the algorithm, that may update its internal model with the new data. At the end of the process, the reward for the algorithm is given by  $r = \sum_{t=1}^n r_t$ , and cumulative regret is defined as  $R_{\mathcal{A}} = \mathbf{E}[r^* - r]$ , where  $r^*$  is the cumulative reward of the optimal policy (i.e., the policy that always selects the action with highest expected reward given the context). The goal is to minimize  $R_{\mathcal{A}}$ .

The main research question we address in this paper is how approximated model posteriors affect the performance of decision making via Thompson Sampling (Algorithm 1) in contextual bandits. We study a variety of algorithmic approaches to approximate a posterior distribution, together with different empirical and synthetic data problems that highlight several aspects of decision making. We consider distributions  $\pi$  over the space of parameters that completely define a problem instance  $\theta \in \Theta$ . For example,  $\theta$  could encode the reward distributions of a set of arms in the multi-armed bandit scenario, or—more generally—all the parameters of an MDP in reinforcement learning.

In the following sections we rely on the idea that, if we had access to the actual posterior  $\pi_t$  given the observed data at all times  $t$ , then choosing actions using Thompson Sampling would lead to near-optimal cumulative regret or, more informally, to good performance. It is important to remark that in some problems this is not necessarily the case; for example, when actions that have no chance of being optimal still convey useful information about other actions. Thompson Sampling (or UCB approaches) would never select such actions, even if they are worth their cost (Russo & Van Roy, 2014). In addition, Thompson Sampling does *not* take into account the time horizon where the process ends, and if known, exploration efforts could be tuned accordingly (Russo et al., 2017). Nonetheless, under the assumption that very accurate posterior approximations lead to efficient decisions, the question is: what happens when the approximations are not so accurate? There are two main possibilities. In some cases, the mismatch in posteriors may not hurt in terms of decision making, and we will still end up with good decisions. Unfortunately, in other cases, this mismatch together with its induced feedback loop will degenerate in a significant loss of performance. We would like to understand the main aspects that determine which way it goes.

**Algorithm 1** Thompson Sampling

---

```

1: Input: Prior distribution over models,  $\pi_0 : \theta \in \Theta \rightarrow [0, 1]$ .
2: for time  $t = 0, \dots, N$  do
3:   Observe context  $X_t \in \mathbb{R}^d$ .
4:   Sample model  $\theta_t \sim \pi_t$ .
5:   Compute  $a_t = \text{BestAction}(X_t, \theta_t)$ .
6:   Select action  $a_t$  and observe reward  $r_t$ .
7:   Update posterior distribution  $\pi_{t+1}$  with  $(X_t, a_t, r_t)$ .

```

---

This is an important practical question as, in large and complex systems, computational sacrifices and statistical assumptions are made to favor simplicity and tractability. But, what is their impact?

### 3 ALGORITHMS

In this section, we describe the main algorithmic design principles behind the specific implementations used to approximately sample from the posterior in our simulations of Section 5. In Figure 3, we visualize the posteriors of the nonlinear algorithms on a synthetic one dimensional problem.

**Linear Methods** We can apply well-known closed-form updates for Bayesian linear regression (Bishop, 2006) for exact posterior inference in linear models. We provide the specific formulas below, and note that they admit a computationally-efficient online version. We consider exact linear posteriors as a baseline; i.e., these formulas lead to the right posterior *only* when the data was generated according to  $Y = X^T \beta + \epsilon$  where  $\epsilon \sim \mathcal{N}(0, \sigma^2)$ . Importantly, we model the *joint* distribution of  $\beta$  and  $\sigma^2$  for each action. Sequentially estimating the noise level  $\sigma^2$  for each action allows the algorithm to adaptively improve its understanding of the volume of the hyperellipsoid of plausible  $\beta$ 's, leading, in general, to a more aggressive initial exploration phase (in both  $\beta$  and  $\sigma^2$ ).

The posterior at time  $t$  for action  $i$ , after observing  $X, Y$ , is  $\pi_t(\beta, \sigma^2) = \pi_t(\beta \mid \sigma^2) \pi_t(\sigma^2)$ , where we assume  $\sigma^2 \sim \text{IG}(a_t, b_t)$ , and  $\beta \mid \sigma^2 \sim \mathcal{N}(\mu_t, \sigma^2 \Sigma_t)$ , an Inverse Gamma and Gaussian distribution, respectively. Their parameters are given by

$$\Sigma_t = (X^T X + \Lambda_0)^{-1}, \quad \mu_t = \Sigma_t (\Lambda_0 \mu_0 + X^T Y), \quad (1)$$

$$a_t = a_0 + t/2, \quad b_t = b_0 + \frac{1}{2} (Y^T Y + \mu_0^T \Sigma_0 \mu_0 - \mu_t^T \Sigma_t^{-1} \mu_t). \quad (2)$$

We set the prior hyperparameters to  $\mu_0 = 0$ , and  $\Lambda_0 = \lambda \text{Id}$ , while  $a_0 = b_0 = \alpha > 1$ .

We consider two approximations to (1) motivated by function approximators where  $d$  is large. While posterior distributions or confidence ellipsoids should capture dependencies across parameters as observed above (say, a dense  $\Sigma_t$ ), in practice, computing the correlations across all pairs of parameters is too expensive, and diagonal covariance approximations are common. For linear models it may still be feasible to exactly compute (1), whereas in the case of Bayesian neural networks, unfortunately, this may no longer be possible. Accordingly, we study two linear approximations where  $\Sigma_t$  is diagonal. Our goal is to understand the impact of such approximations in the simplest case, to properly set our expectations for the loss in performance of equivalent approximations in more complex approaches, like mean-field variational inference or Stochastic Gradient Langevin Dynamics.

Assume for simplicity the noise standard deviation is known. In Figure 2a, for  $d = 2$ , we see the posterior distribution  $\beta_t \sim \mathcal{N}(\mu_t, \Sigma_t)$  of a linear model based on (1), in red, together with two diagonal approximations. Each approximation tries to minimize a different objective. In blue, the *PrecisionDiag* posterior approximation finds the diagonal  $\hat{\Sigma} \in \mathbb{R}^{d \times d}$  minimizing  $\text{KL}(\mathcal{N}(\mu_t, \hat{\Sigma}) \parallel \mathcal{N}(\mu_t, \Sigma_t))$ , like in mean-field variational inference. In particular,  $\hat{\Sigma} = \text{Diag}(\Sigma_t^{-1})^{-1}$ . On the other hand, in green, the *Diag* posterior approximation finds the diagonal matrix  $\bar{\Sigma}$  minimizing  $\text{KL}(\mathcal{N}(\mu_t, \Sigma_t) \parallel \mathcal{N}(\mu_t, \bar{\Sigma}))$  instead, similarly to expectation-propagation algorithms. In this case, the solution is simply  $\bar{\Sigma} = \text{Diag}(\Sigma_t)$ .

**Neural Linear** The main problem linear algorithms face is their lack of representational power, which they complement with accurate uncertainty estimates. A natural attempt at getting the best of both worlds consists in performing a Bayesian linear regression on top of the representation of the last layer of a neural network. The predicted value  $v_i$  for each action  $a_i$  is given by  $v_i = \beta_i^T z_x$ ,

where  $z_x$  is the output of the last hidden layer of the network for context  $x$ . While linear methods directly try to regress values  $v$  on  $x$ , we can independently train a deep net to learn a representation  $z$ , and then use a Bayesian linear regression to regress  $v$  on  $z$ , obtain uncertainty estimates on the  $\beta$ 's, and make decisions accordingly via Thompson Sampling. Note that we do not explicitly consider the weights of the linear output layer of the network to make decisions; further, the network is *only* used to find good representations  $z$ . In addition, we can update the network and the linear regression at different time-scales. It makes sense to keep an exact linear regression (as in (1) and (2)) at all times, adding each new data point as soon as it arrives. However, we only update the network after a number of points have been collected. In our experiments, after updating the network, we perform a forward pass on all the training data to obtain their latest representation, which is then fed to the Bayesian regression. In practice this may be too expensive, and the last representation computed on each point could be used instead with online updates on the regression. We call this algorithm *Neural Linear*.

**Neural Greedy** We refer to the algorithm that simply trains a neural network and acts greedily (i.e., takes the action whose predicted score for the current context is highest) as **RMS**, as we train it using the RMSProp optimizer. This is our non-linear baseline, and we tested several versions of it (based on whether the training step was decayed, reset to its initial value for each re-training or not, and how long the network was trained for). We also tried the  $\epsilon$ -greedy version of the algorithm, where a random action was selected with probability  $\epsilon$  for some decaying schedule of  $\epsilon$ . Somewhat surprisingly, this did not help.

**Variational Inference** Variational approaches approximate the posterior by finding a distribution within a tractable family that minimizes the KL divergence to the posterior (Hinton & Van Camp, 1993). Typically (and in our experiments), the posterior is approximated by a mean-field or factorized distribution. Recent advances in variational inference have scaled these approaches to estimate the posterior of neural networks with millions of parameters (Blundell et al., 2015). A common criticism of variational inference is that it underestimates uncertainty (e.g., (Bishop, 2006)), which could lead to under-exploration.

Dropout (Srivastava et al., 2014) can be seen as optimizing a variational objective (Kingma et al., 2015; Gal & Ghahramani, 2016), which leads to a straight-forward posterior approximation and application to Thompson sampling.

**Monte Carlo** Monte Carlo sampling remains one of the simplest and reliable tools in the Bayesian toolbox. Rather than parameterizing the full posterior, Monte Carlo methods estimate the posterior through drawing samples. This is naturally appealing for highly parameterized deep neural networks for which the posterior is intractable in general and even simple approximations such as multivariate Gaussian are too expensive (i.e. require computing and inverting a covariance matrix over all parameters). Among Monte Carlo methods, Hamiltonian Monte Carlo (Neal, 1994) (HMC) is often regarded as a gold standard algorithm for neural networks as it takes advantage of gradient information and momentum to more effectively draw samples. However, it remains unfeasible for larger datasets as it involves a Metropolis accept-reject step that requires computing the log likelihood over the whole data set. A variety of methods have been developed to approximate HMC using mini-batch stochastic gradients. These Stochastic Gradient Langevin Dynamics (SGLD) methods (Neal, 1994; Welling & Teh, 2011) add Gaussian noise to the model gradients during stochastic gradient updates in such a manner that each update results in an approximate sample from the posterior. Different strategies have been developed for augmenting the gradients and noise according to a preconditioning matrix. Li et al. (2013) show that a preconditioner based on the RMSprop algorithm performs well on deep neural networks. Patterson & Teh (2013) suggested using the Fisher information matrix as a preconditioner in SGLD. Unfortunately theory requires that the approximations of SGLD hold only if the learning rate is asymptotically annealed to zero. Ahn et al. (2012) introduced Stochastic Gradient Fisher Scoring to elegantly remove this requirement by preconditioning according to the Fisher information (or a diagonal approximation thereof). Mandt et al. (2016) develop methods for approximately sampling from the posterior using a constant learning rate in stochastic gradient descent and develop a prescription for a stable version of SGFS. We evaluate the diagonal-SGFS and *constant-SGD* algorithms from Mandt et al. (2016) in this work. Specifically for constantSGD we use a constant learning rate for stochastic gradient descent, where the learning rate,  $\epsilon$  is given by  $\epsilon = 2 \frac{S}{N} \mathbf{B}\mathbf{B}^T$  where  $S$  is the batch size,  $N$  the number of data points and  $\mathbf{B}\mathbf{B}^T$  is an online average of the diagonal empirical Fisher information matrix. For Stochastic Gradient Fisher Scoring we use

the following stochastic gradient update for the model parameters  $\theta$  at step  $t$ :

$$\theta(t+1) = \theta(t) - \epsilon \mathbf{H} g(\theta(t)) + \sqrt{\epsilon} \mathbf{H} \mathbf{E} \nu, \quad \nu \sim \mathcal{N}(0, I) \quad (3)$$

where we take the noise covariance  $\mathbf{E}\mathbf{E}^T$  to also be  $\mathbf{B}\mathbf{B}^T$  and  $\mathbf{H} = \frac{2}{N}(\epsilon\mathbf{B}\mathbf{B}^T + \mathbf{E}\mathbf{E}^T)^{-1}$ .

**Bootstrap** A simple and empirical approach to approximate the sampling distribution of any estimator is the Bootstrap (Efron, 1982). The main idea is to simultaneously train  $q$  models, where each model  $i$  is based on a different dataset  $D_i$ . When all the data  $D$  is available in advance,  $D_i$  is typically created by sampling  $|D|$  elements from  $D$  at random with replacement. In our case, however, the data grows one example at a time. Accordingly, we set a parameter  $p \in (0, 1]$ , and append the new datapoint to each  $D_i$  independently at random with probability  $p$ . In order to emulate Thompson Sampling, we sample a model uniformly at random (i.e., with probability  $p_i = 1/q$ .) and take the action predicted to be best by the sampled model. We tested versions where  $q = 5, 10$  and  $p = 0.8, 1.0$ , with neural network models. Note that even when  $p = 1$  and the datasets are identical, the random initialization of each network, together with the randomness from SGD, leads to different predictions.

**Direct Noise Injection** Parameter-Noise (Plappert et al., 2017) is a recently proposed approach for exploration in deep RL that has shown promising results. The training updates for the network are unchanged, but when selecting actions, the network weights are perturbed with isotropic Gaussian noise. Crucially, the network uses layer normalization (Ba et al., 2016), which ensures that all weights are on the same scale. The magnitude of the Gaussian noise is adjusted so that the overall effect of the perturbations is similar in scale to  $\epsilon$ -greedy with a linearly decaying schedule (see (Plappert et al., 2017) for details). Because the perturbations are done on the model parameters, we might hope that the actions produced by the perturbations are more sensible than  $\epsilon$ -greedy.

**Bayesian Non-parametric** Gaussian processes (Rasmussen & Williams, 2005) are a gold-standard method for modeling distributions over non-linear continuous functions. It can be shown that, in the limit of infinite hidden units and under a Gaussian prior, a Bayesian neural network converges to a Gaussian process (Neal, 1994). As such, GPs would appear to be a natural baseline. Unfortunately, standard GPs computationally scale cubically in the number of observations, limiting their applicability to relatively small datasets. There are a wide variety of methods to approximate Gaussian processes using, for example, pseudo-observations (Snelson & Ghahramani) or variational inference (Titsias, 2009). However, these are outside the scope of this comparison.<sup>1</sup> Due to the scaling issue, we stop adding inputs to the GP after 1000 observations. This performed significantly better than randomly sampling inputs. Our implementation is a multi-task Gaussian process (Bonilla et al., 2008) with a linear and Matern 3/2 product kernel over the inputs and an exponentiated quadratic kernel over latent vectors for the different tasks. The hyperparameters of this model and the latent task vectors are optimized over the GP marginal likelihood. This allows the model to learn correlations between the outputs of the model. Specifically, the covariance function  $K(\cdot)$  of the GP is given by:

$$K(\mathbf{x}^k, \hat{\mathbf{x}}^l) = k_{\text{matern}}(\mathbf{x}^k, \hat{\mathbf{x}}^l) \cdot k_{\text{lin}}(\mathbf{x}^k, \hat{\mathbf{x}}^l) \cdot k_{\text{task}}(\mathbf{v}^k, \mathbf{v}^l) \quad (4)$$

$$k_{\text{matern}}(\mathbf{x}^k, \hat{\mathbf{x}}^l) = \alpha(1 + \sqrt{3}r_{\lambda_m}(\mathbf{x}, \hat{\mathbf{x}})) \exp(-\sqrt{3}r_{\lambda_m}(\mathbf{x}, \hat{\mathbf{x}})) \quad (5)$$

$$k_{\text{lin}}(\mathbf{x}^k, \hat{\mathbf{x}}^l) = \beta(\mathbf{x} \odot \lambda_l)(\hat{\mathbf{x}} \odot \lambda_l)^T \quad (6)$$

and the task kernel between tasks  $t$  and  $l$  are  $k_t(\mathbf{x}^k, \hat{\mathbf{x}}^l) = \exp(-r_{\lambda_m}(\mathbf{v}^k, \mathbf{v}^l))^2$  where  $\mathbf{v}^l$  indexes the latent vector for task  $l$  and  $r_{\lambda}(\mathbf{x}, \hat{\mathbf{x}}) = |(\mathbf{x} \odot \lambda) - (\hat{\mathbf{x}} \odot \lambda)|$ . The length-scales,  $\lambda_m$  and  $\lambda_l$ , and amplitude parameters  $\alpha, \beta$  are optimized via the log marginal likelihood.

There are a variety of other methods, including expectation propagation (Hernández-Lobato & Adams, 2015) and  $\alpha$ -divergence minimization of which variational inference is a special case (Hernández-Lobato et al., 2016). We leave evaluating these approaches to future work.

## 4 FEEDBACK LOOP IN THE LINEAR CASE

In this section, we illustrate some of the subtleties that arise when uncertainty estimates drive sequential decision-making using simple linear examples.

There is a fundamental difference between *static* and *dynamic* scenarios. In a static scenario, e.g. supervised learning, we are given a model family  $\Theta$  (like the set of linear models, trees, or

<sup>1</sup>We invite others to apply these and new methods upon the open source release of this benchmark.

neural networks with specific dimensions), a prior distribution  $\pi_0$  over  $\Theta$ , and some observed, and importantly assumed i.i.d., data  $\mathbf{D}$ . Our goal is to return an approximate posterior distribution:  $\tilde{\pi} \approx \pi = \mathbf{P}(\theta \mid \mathbf{D})$ . Let the distance  $d(\tilde{\pi}, \pi)$  represent the quality of our approximation.

On the other hand, in *dynamic* settings, our estimate at time  $t$ , say  $\tilde{\pi}_t$ , will be used via some mechanism  $\mathbf{M}$ , in this case Thompson sampling, to collect the next data-point, which is then appended to  $\mathbf{D}_t$ . In this case, the data-points in  $\mathbf{D}_t$  are no longer independent.  $\mathbf{D}_t$  will now determine two distributions: the posterior given the data that was actually observed,  $\pi_{t+1} = \mathbf{P}(\theta \mid \mathbf{D}_t)$ , and our estimate  $\tilde{\pi}_{t+1}$ . When the goal is to make good sequential decisions in terms of cumulative regret, the distance  $d(\tilde{\pi}_t, \pi_t)$  is in general no longer a definitive proxy for performance. For instance, a poorly-approximated decision boundary could lead an algorithm, based on  $\tilde{\pi}$ , to get stuck repeatedly selecting a single sub-optimal action  $a$ . After collecting lots of data for that action,  $\tilde{\pi}_t$  and  $\pi_t$  could start to agree (to their capacity) on the models that explain what was observed for  $a$ , while both would stick to something close to the prior regarding the other actions. At that point,  $d(\tilde{\pi}_t, \pi_t)$  may show relatively little disagreement, but the regret would already be terrible.

Let  $\pi_t^*$  be the posterior distribution  $\mathbf{P}(\theta \mid \mathbf{D}_t)$  under the Thompson Sampling assumption, that is, data was always collected according to  $\pi_j^*$  for  $j < t$ . We follow the idea that  $\tilde{\pi}_t$  being close to  $\pi_t^*$  for all  $t$  leads to strong performance. However, this concept is difficult to formalize: once different decisions are made, data for different actions is collected and it is hard to compare posterior distributions.

We illustrate the previous points with a simple example, see Figure 1. Data is generated according to a bandit with  $k = 6$  arms. For a given context  $X \sim \mathcal{N}(\mu, \Sigma)$ , the reward obtained by pulling arm  $i$  follows a linear model  $r_{i,X} = X^T \beta_i + \epsilon$  with  $\epsilon \sim \mathcal{N}(0, \sigma_i^2)$ . The posterior distribution over  $\beta_i \in \mathbb{R}^d$  can be exactly computed using the standard Bayesian linear regression formulas presented in Section 3. We set the contextual dimension  $d = 20$ , and the prior to be  $\beta \sim \mathcal{N}(0, \lambda \mathbf{I}_d)$ , for  $\lambda > 0$ .

In Figure 1, we show the posterior distribution for two dimensions of  $\beta_i$  for each arm  $i$  after  $n = 500$  pulls. In particular, in Figure 1a, two independent runs of Thompson Sampling with their posterior distribution are displayed in red and green. While strongly aligned, the estimates for some arms disagree (especially for arms that are best only for a *small* fraction of the contexts, like Arm 2 and 3, where fewer data-points are available). In Figure 1b, we also consider Thompson Sampling with an approximate posterior with diagonal covariance matrix, *Diag* in red, as defined in Section 3. Each algorithm collects its own data based on its current posterior (or approximation). In this case, the posterior disagreement after  $n = 500$  decisions is certainly stronger. However, as shown in Figure 1c, if we computed the approximate posterior with a diagonal covariance matrix based on the data collected by the *actual posterior*, the disagreement would be reduced as much as possible within the approximation capacity (i.e., it still cannot capture correlations in this case). Figure 1b shows then the effect of the feedback loop. We look next at the impact that this mismatch has on regret.

We illustrate with a similar example how inaccurate posteriors sometimes lead to quite different behaviors in terms of regret. In Figure 2a, we see the posterior distribution  $\beta \sim \mathcal{N}(\mu, \Sigma)$  of a linear model in green, together with the two diagonal linear approximations introduced in Section 3: the *Diag* (in red) and the *PrecisionDiag* (in blue) approximations, respectively. We now assume there are  $k$  linear arms,  $\beta_i \in \mathbb{R}^d$  for  $i = 1, \dots, k$ , and decisions are made according to the posteriors in Figure 2a. In Figures 2b and 2c we plot the regret of Thompson Sampling when there are  $k = 20$  arms, for both  $d = 15$  and  $d = 30$ . We see that, while the *PrecisionDiag* approximation does even outperform the actual posterior, the diagonal covariance approximation truly suffers poor regret when we increase the dimension  $d$ , as it is heavily penalized by *simultaneously* over-exploring in a large number of dimensions and repeatedly acting according to implausible models.

## 5 EMPIRICAL EVALUATION

In this section, we present the simulations and outcomes of several synthetic and real-world data bandit problems with each of the algorithms introduced in Section 3. In particular, we first explain how the simulations were set up and run, and the metrics we report. We then split the experiments according to how data was generated, and the underlying models fit by the algorithms from Section 3.

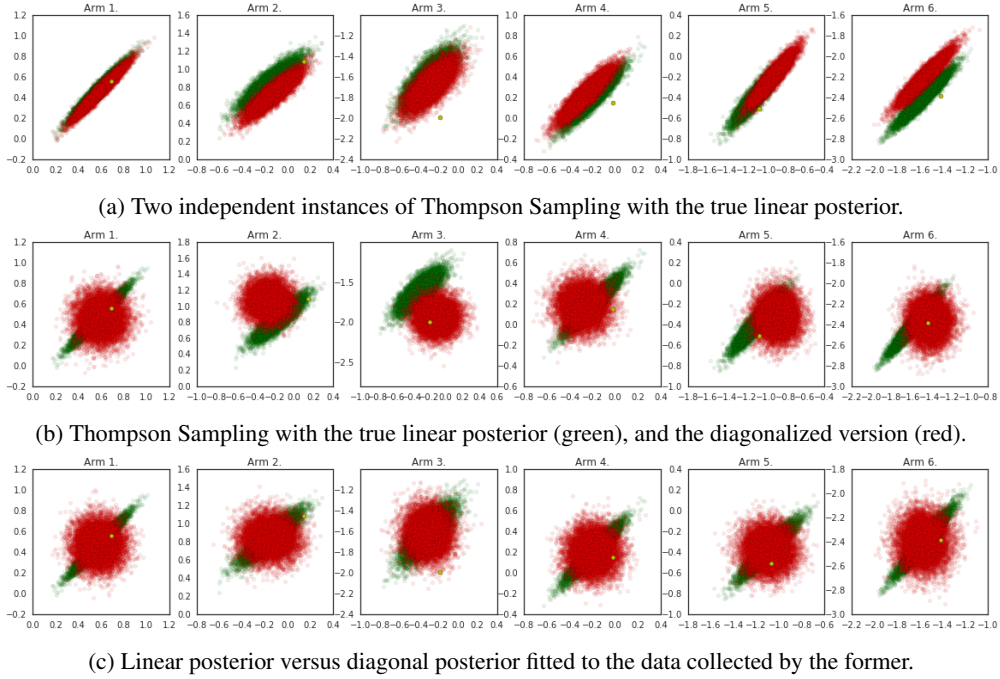


Figure 1: Visualizations of the posterior approximations in a linear example.

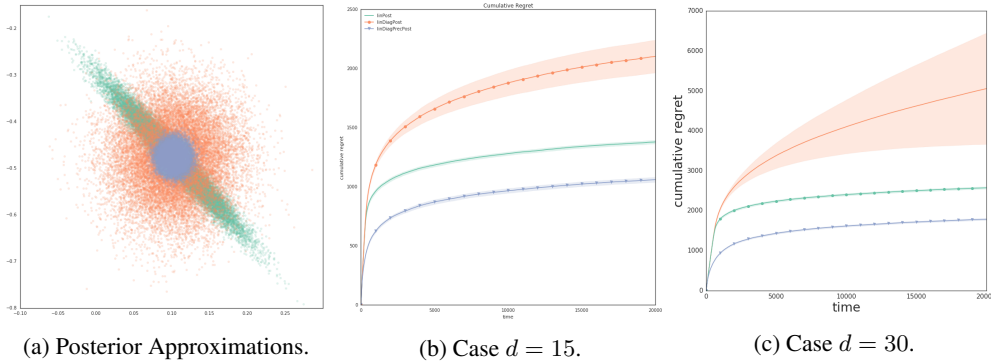


Figure 2: In The impact on regret of different approximated posteriors. We show (green) the actual linear posterior, (red) the diagonal posterior approximation and (blue) the precision approximation in 2a. In 2b and 2c we visualize the impact of the approximations on cumulative regret.

## 5.1 THE EXPERIMENTAL FRAMEWORK

We run the contextual bandit experiments as described at the beginning of Section 2, and discuss below some implementation details of both experiments and algorithms.

**Neural Network Architectures.** All algorithms based on neural networks as function approximators share the same architecture. In particular, we fit a simple fully-connected feedforward network with two hidden layers with 100 units each and softplus activations.<sup>2</sup> The input of the network has dimension  $d$  (same as the contexts), and there are  $k$  outputs, one per action. Note that for each training point  $(X_t, a_t, r_t)$  only one action was observed.

**Updating Models.** A key question is how often and for how long models are updated. Ideally, we would like to train after each new observation and for as long as possible. However, this may limit the applicability of our algorithms in online scenarios where decisions must be made immediately.

<sup>2</sup>We also tried relu activations, which led to similar results.

We update linear algorithms after each time-step by means of (1) and (2). For neural networks, the default behavior was to train for  $t_s = 20$  mini-batches every  $t_f = 20$  timesteps<sup>3</sup>. The size of each mini-batch was 512. We experimented with increasing values of  $t_s$ , and it proved essential for some algorithms like variational inference approaches.

**Metrics** We report two metrics: cumulative regret and simple regret. We approximate the latter as the mean cumulative regret in the last 500 time-steps, a proxy for the quality of the final policy (see further discussion on *pure* exploration settings, Bubeck et al. (2009)).

**Hyper-Parameter Tuning** Deep learning methods are known to be very sensitive to the selection of a wide variety of hyperparameters and many of the algorithms presented are no exception. Moreover, that choice is known to be highly dataset dependent. However, in the bandits scenario, we would not have access to each problem a-priori to perform tuning. Thus, for each algorithm, the same hyperparameters are used across all tasks; for some algorithms they were chosen through careful tuning over cumulative regret on a small subset of linear bandits (Dropout, SGFS, RMS, BBBN, pNoise) and some (ConstSGD, Neural Linear and all linear methods) required very little tuning.

**Buffer** After some experimentation, we decided not to use a data buffer as some evidence of catastrophic forgetting was observed. All observations are sampled with equal probability to be part of a mini-batch. In addition, as is standard in bandit algorithms, each action was initially selected  $s = 3$  times using round-robin independently of the context.

## 5.2 REAL-WORLD DATA PROBLEMS WITH NON-LINEAR MODELS

We test the algorithms in a broad range of bandit problems created using real-world data. In particular, we test on the Mushroom, Statlog, Coverttype, Financial, Jester, Adult, Census, and Song datasets. A detailed description of each dataset and its bandit problem is provided in the Appendix, Section A. They exhibit a broad range of properties: small and large sizes, one dominating action versus more homogeneous optimality, learnable or little signal, stochastic or deterministic rewards, etc.

For space reasons, the outcome of some datasets are directly presented in the Appendix. The Statlog, Coverttype, Adult, and Census datasets were originally tested in Elmachet et al. (2017). The results are displayed in Table 1 (cumulative regret) and Table 2 (simple regret).

## 5.3 REAL-WORLD DATA PROBLEMS WITH LINEAR MODELS

As the algorithms from Section 3 can be implemented for any model architecture, here we use linear models as a baseline comparison across algorithms. The datasets are the same as in the previous subsection. The results are displayed in Table 1 (cumulative regret) and Table 2 (simple regret).

## 5.4 THE WHEEL BANDIT

Some of the real-data problems presented above do not require significant exploration. We design an artificial problem where the need for exploration is smoothly parameterized. The *wheel* bandit is defined as follows (see Figure 5). Set  $d = 2$ , and  $\delta \in (0, 1)$ , the exploration parameter. Contexts are sampled uniformly at random in the unit circle,  $X \sim U(\mathbf{D})$ . There are  $k = 5$  possible actions. The first action  $a_1$  always offers reward  $r \sim \mathcal{N}(\mu_1, \sigma^2)$ , independently of the context. For contexts inside the blue circle, i.e.  $\|X\| \leq \delta$ , the other four actions are equal and sub-optimal, with  $r \sim \mathcal{N}(\mu_2, \sigma^2)$  for  $\mu_2 < \mu_1$ . When  $\|X\| > \delta$ , we are outside the blue circle, and one of the actions  $a_1, \dots, a_4$  is optimal depending on the sign of  $X_1$  and  $X_2$ . If  $X_1, X_2 > 0$ , action 2 is optimal. If  $X_1 > 0, X_2 < 0$ , action 3 is optimal, and so on. Non-optimal actions still deliver  $r \sim \mathcal{N}(\mu_2, \sigma^2)$ , while the optimal action provides  $r \sim \mathcal{N}(\mu_3, \sigma^2)$ , with  $\mu_3 \gg \mu_1$ . We set  $\mu_1 = 1.2, \mu_2 = 1.0, \mu_3 = 50.0$ , and  $\sigma = 0.01$ . Note that the probability of a context randomly falling in the high-reward region is  $1 - \delta^2$ .

The difficulty of the problem increases with  $\delta$ , and we expect algorithms to get stuck repeatedly selecting action  $a_1$ . The problem can be easily generalized for  $d > 2$ . Results are shown in Table 3.

<sup>3</sup>For reference, the standard strategy for Deep Q-Networks on Atari is to make one model update after every 4 actions performed (Mnih et al., 2015; Osband et al., 2016; Plappert et al., 2017; Meire Fortunato, 2017).



## 6 DISCUSSION

We discuss next our main findings for each class of algorithms.

**Linear Methods.** Linear methods offer a reasonable baseline, strong in many cases. While their representation power is certainly a limiting factor, their ability to compute informative uncertainty measures seems to payoff and balance their initial disadvantage. They do well in several datasets, and are able to react fast to surprising or extreme rewards (maybe as single points can have a heavy impact in fitted models, and their updates are deterministic and exact). Some datasets clearly need more complex non-linear representations, and linear methods are unable to efficiently solve those. In addition, linear methods obviously offer computational advantages.

**NeuralLinear.** The NeuralLinear algorithm sits near a sweet spot that is worth further studying. In general it seems to consistently beat the RMS neural network it is based on, suggesting its exploration mechanisms add concrete value. We believe its main strength relies on the fact that it is able to *simultaneously* learn a data representation that greatly simplifies the task at hand, and to accurately quantify the uncertainty over linear models that explain the observed rewards in terms of the proposed representation. While the former process may be noisier and heavily dependent on the amount of training steps that were taken and available data, the latter always offers the exact solution to its approximate parent problem. This, together with the partial success of linear methods with poor representations, may explain its promising results. In some sense, it knows what it knows. While conceptually simple, its deployment to large scale systems may involve some technical difficulties; mainly, to update the Bayesian estimates when the network is re-trained. We believe, however, standard solutions to similar problems (like running averages) could greatly mitigate these issues.

**Variational Inference.** Overall, BBB performed poorly, even when the model was linear. Variational methods are known to underestimate uncertainty (Bishop, 2006), so we expected that might lead to poor exploration, however, this does not explain the poor performance entirely. When the model is linear, LinDiagPrecPost computes the exact solution to the minimization problem BBB optimizes with stochastic variational inference. The poor performance of BBB relative to LinDiagPrecPost suggests that the problem is more than underexploration. Figure 4 shows that performance improves as we increase the number of training steps, suggesting that it is not sufficient to partially optimize the variational parameters when the uncertainty estimates directly affect the data being collected. Optimizing to convergence at every step significantly reduces regret, but at major computational cost.

**Monte Carlo.** ConstantSGD comes out as the winner on Statlog and Covertype, two problems that both appear to require non-linearity and exploration as evidenced by performance of the linear baseline approaches. The method is especially appealing as it doesn't require tuning learning rates or exploration parameters. SGFS, however, performs significantly worse on most benchmarks except the hard exploration variant of the wheel problem where it shines. The additional injected noise in SGFS may cause the model to explore more than other approaches and thus incur additional regret.

**Bootstrap.** The improvement, if any, provided by bootstrapping neural networks over the base model is surprisingly small. Moreover, it might not justify the important computational overhead of the method. The randomness from SGD may be enough for exploration in many bandits problems.

**Direct Noise Injection.** Parameter-Noise was based on the RMS2 implementation. Results suggest that it may only offer modest improvements over its base learner, while we found the algorithm hard to tune, and quite sensitive to the heuristic controlling the injected noise-level. In addition, developing an intuition for the heuristic is not straightforward as it lacks transparency, and may require previous and repeated access to the decision-making process.

On the other hand, we mainly experimented with two dropout versions: fixed  $p = 0.5$ , and  $p = 0.8$ . The latter consistently delivered better results, and it is the one we report. The algorithm required no extra-tuning, and sometimes offered small improvements. Dropout was used both for training and for decision-making; unfortunately, we did not add a baseline where dropout only applies during training. Consequently, it is not obvious how to disentangle the contribution of better training from that of better exploration. This remains as future work.

**Bayesian Non-parametrics.** Perhaps unsurprisingly, the Gaussian process performs remarkably well on problems with little data but struggles on larger problems. This may motivate the use of sparse GP-based approaches in the future.

## 7 CONCLUSIONS AND FUTURE WORK

In this work, we empirically studied the impact on performance of approximate model posteriors for decision making via Thompson Sampling in contextual bandits. The methods we found to perform best exactly measure uncertainty (possibly under the wrong model assumptions) on top of complex representations learned in parallel. More complicated approaches that learn the representation and its uncertainty together seemed to require heavier training, an important drawback in online scenarios, and exhibited stronger hyper-parameter dependence. Further exploring and developing the promising approaches is an exciting avenue for future work.

### ACKNOWLEDGMENTS

### REFERENCES

- Agrawal, Shipra and Goyal, Navin. Analysis of thompson sampling for the multi-armed bandit problem. In *Conference on Learning Theory*, 2012.
- Ahn, Sungjin, Balan, Anoop Korattikara, and Welling, Max. Bayesian posterior sampling via stochastic gradient fisher scoring. In *International Conference on Machine Learning*, 2012.
- Asuncion, Arthur and Newman, David. Uci machine learning repository, 2007.
- Ba, Jimmy Lei, Kiros, Jamie Ryan, and Hinton, Geoffrey E. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- Bertin-Mahieux, Thierry, Ellis, Daniel P.W., Whitman, Brian, and Lamere, Paul. The million song dataset. In *Proceedings of the 12th International Conference on Music Information Retrieval (ISMIR 2011)*, 2011.
- Bishop, Christopher M. *Pattern recognition and machine learning*. springer, 2006.
- Blundell, Charles, Cornebise, Julien, Kavukcuoglu, Koray, and Wierstra, Daan. Weight uncertainty in neural network. In *International Conference on Machine Learning*, pp. 1613–1622, 2015.
- Bonilla, Edwin V, Chai, Kian M., and Williams, Christopher. Multi-task gaussian process prediction. In *Advances in Neural Information Processing Systems*. 2008.
- Bubeck, Sébastien, Munos, Rémi, and Stoltz, Gilles. Pure exploration in multi-armed bandits problems. In *International conference on Algorithmic learning theory*, pp. 23–37. Springer, 2009.
- Chapelle, Olivier and Li, Lihong. An empirical evaluation of thompson sampling. In *Advances in Neural Information Processing Systems 24*. 2011.
- Efron, Bradley. *The jackknife, the bootstrap and other resampling plans*. SIAM, 1982.
- Elmachoub, Adam N, McNellis, Ryan, Oh, Sechan, and Petrik, Marek. A practical method for solving contextual bandit problems using decision trees. *arXiv preprint arXiv:1706.04687*, 2017.
- Gal, Yarin and Ghahramani, Zoubin. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *International conference on machine learning*, pp. 1050–1059, 2016.
- Goldberg, Ken, Roeder, Theresa, Gupta, Dhruv, and Perkins, Chris. Eigentaste: A constant time collaborative filtering algorithm. *Information Retrieval*, 4(2):133–151, 2001.
- Granmo, Ole-Christoffer. Solving two-armed bernoulli bandit problems using a bayesian learning automaton. *International Journal of Intelligent Computing and Cybernetics*, 3(2):207–234, 2010.
- Graves, Alex. Practical variational inference for neural networks. In *Advances in Neural Information Processing Systems*, pp. 2348–2356, 2011.
- Hernández-Lobato, José Miguel and Adams, Ryan P. Probabilistic backpropagation for scalable learning of bayesian neural networks. In *International Conference on Machine Learning*, 2015.

- Hernández-Lobato, José Miguel, Li, Yingzhen, Rowland, Mark, Bui, Thang D., Hernández-Lobato, Daniel, and Turner, Richard E. Black-box alpha divergence minimization. In *International Conference on Machine Learning*, 2016.
- Hinton, Geoffrey E and Van Camp, Drew. Keeping the neural networks simple by minimizing the description length of the weights. In *Proceedings of the sixth annual conference on Computational learning theory*, pp. 5–13. ACM, 1993.
- Kingma, Diederik P, Salimans, Tim, and Welling, Max. Variational dropout and the local reparameterization trick. In *Advances in Neural Information Processing Systems*, pp. 2575–2583, 2015.
- Kohavi, Ron. Scaling up the accuracy of naive-bayes classifiers: A decision-tree hybrid.
- Li, Ke, Swersly, Kevin, Ryan, and Zemel, Richard S. Efficient feature learning using perturb-and-map. *NIPS Workshop on Perturbations, Optimization, and Statistics*, 2013.
- Mandt, Stephan, Hoffman, Matthew D., and Blei, David M. A variational analysis of stochastic gradient algorithms. In *International Conference on Machine Learning*, 2016.
- Meire Fortunato, Mohammad Gheshlaghi Azar, Bilal Piot Jacob Menick Ian Osband Alex Graves Vlad Mnih Remi Munos Demis Hassabis Olivier Pietquin Charles Blundell Shane Legg. Noisy networks for exploration. *arXiv:1706.10295*, 2017.
- Mnih, Volodymyr, Kavukcuoglu, Koray, Silver, David, Rusu, Andrei A., Veness, Joel, Bellemare, Marc G., Graves, Alex, Riedmiller, Martin, Fidjeland, Andreas K., Ostrovski, Georg, Petersen, Stig, Beattie, Charles, Sadik, Amir, Antonoglou, Ioannis, King, Helen, Kumaran, Dharshan, Wierstra, Daan, Legg, Shane, and Hassabis, Demis. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, February 2015.
- Neal, Radford M. Bayesian learning for neural networks. *Dept. of Computer Science, University of Toronto*, 1994.
- Osband, Ian, Blundell, Charles, Pritzel, Alexander, and Van Roy, Benjamin. Deep exploration via bootstrapped dqn. In *Advances in Neural Information Processing Systems*, pp. 4026–4034, 2016.
- Patterson, Sam and Teh, Yee Whye. Stochastic gradient riemannian langevin dynamics on the probability simplex. In *Advances in Neural Information Processing Systems*. 2013.
- Plappert, Matthias, Houthoofd, Rein, Dhariwal, Prafulla, Sidor, Szymon, Chen, Richard Y., Chen, Xi, Asfour, Tamim, Abbeel, Pieter, and Andrychowicz, Marcin. Parameter space noise for exploration. *arXiv:1706.01905*, 2017.
- Rasmussen, Carl Edward and Williams, Christopher K. I. *Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning)*. The MIT Press, 2005.
- Riquelme, Carlos, Ghavamzadeh, Mohammad, and Lazaric, Alessandro. Active learning for accurate estimation of linear models. In Precup, Doina and Teh, Yee Whye (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 2931–2939, International Convention Centre, Sydney, Australia, 06–11 Aug 2017. PMLR. URL <http://proceedings.mlr.press/v70/riquelme17a.html>.
- Russo, Dan and Van Roy, Benjamin. Learning to optimize via information-directed sampling. In *Advances in Neural Information Processing Systems*, pp. 1583–1591, 2014.
- Russo, Daniel, Tse, David, and Van Roy, Benjamin. Time-sensitive bandit learning and satisficing thompson sampling. *arXiv preprint arXiv:1704.09028*, 2017.
- Schaul, Tom, Quan, John, Antonoglou, Ioannis, and Silver, David. Prioritized experience replay. In *International Conference on Learning Representations*, 2016.
- Schlimmer, Jeff. Mushroom records drawn from the audubon society field guide to north american mushrooms. *GH Lincoff (Pres)*, New York, 1981.

- Snelson, Edward and Ghahramani, Zoubin. Sparse gaussian processes using pseudo-inputs. In Weiss, Y., Schölkopf, P. B., and Platt, J. C. (eds.), *Advances in Neural Information Processing Systems*.
- Srivastava, Nitish, Hinton, Geoffrey E, Krizhevsky, Alex, Sutskever, Ilya, and Salakhutdinov, Ruslan. Dropout: a simple way to prevent neural networks from overfitting. *Journal of machine learning research*, 15(1):1929–1958, 2014.
- Thompson, William R. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3/4):285–294, 1933.
- Titsias, Michalis K. Variational learning of inducing variables in sparse gaussian processes. In *International Conference on Artificial Intelligence and Statistics*, 2009.
- Welling, Max and Teh, Yee Whye. Bayesian learning via stochastic gradient langevin dynamics. In *International Conference on Machine Learning*, 2011.

## APPENDIX

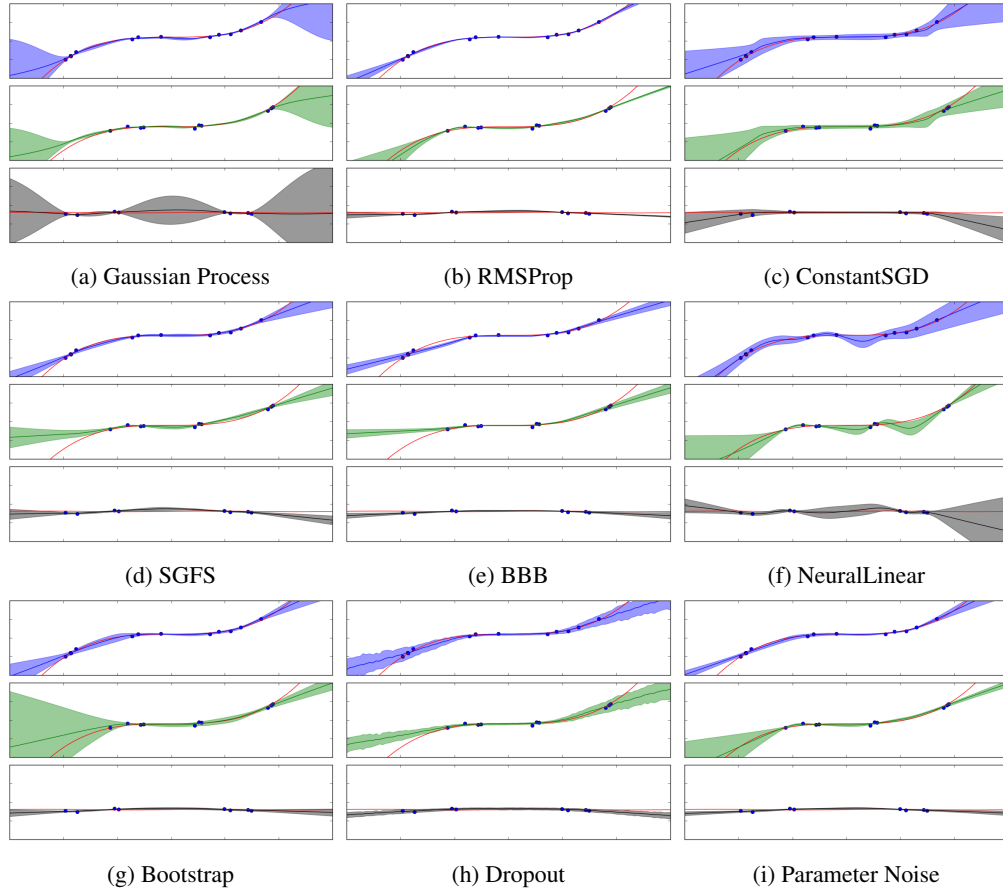


Figure 3: We qualitatively compare plots of the sample distribution from various methods, similarly to Hernández-Lobato et al. (2016). We plot the mean and standard deviation of 100 samples drawn from each method conditioned on a small set of observations with three outputs (two are from the same underlying function and thus strongly correlated while the third (bottom) is independent). The true underlying functions are plotted in red.

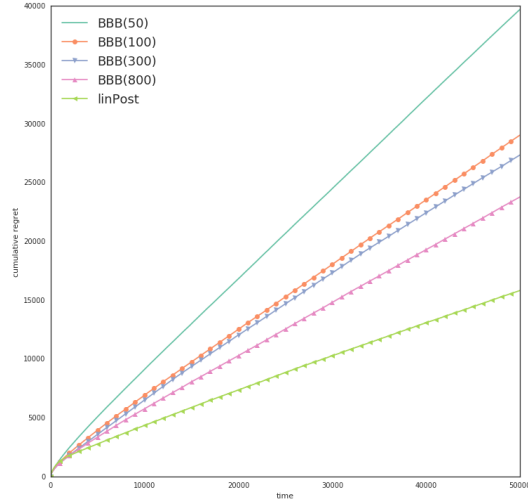


Figure 4: Performance of BBB algorithms with linear models in Mushroom Dataset for increasing number of training steps every 20 time-steps (in parenthesis next to the name).

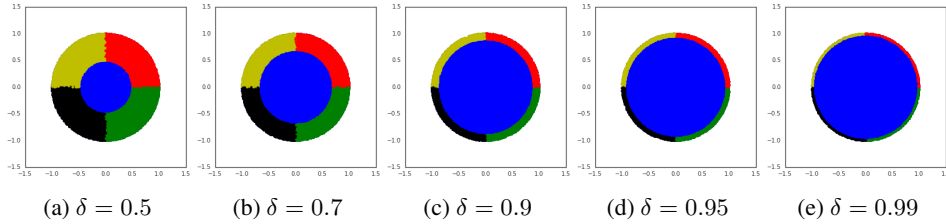


Figure 5: Chrome bandits for increasing values of  $\delta \in (0, 1)$ . Optimal action for yellow, red, black, green, and blue regions, are actions 1, 2, 3, 4, and 5, respectively.

Table 1: Cumulative regret incurred by the algorithms in Section 3 on the bandits described in Section A. The methods that can take advantage of nonlinear models use them. Results are all relative to the cumulative regret of NeuralLinear.

	Mushroom	Statlog	Covertime	Financial	Jester	Adult
pNoise2	205.41 $\pm$ 63.39	115.41 $\pm$ 20.02	110.55 $\pm$ 1.00	111.53 $\pm$ 7.48	99.62 $\pm$ 0.91	101.16 $\pm$ 0.69
linfullPost	120.93 $\pm$ 13.73	642.07 $\pm$ 1.74	129.61 $\pm$ 0.08	78.26 $\pm$ 0.79	82.40 $\pm$ 0.68	<b>93.42 <math>\pm</math> 0.04</b>
linfullDiagPrecPost	230.46 $\pm$ 33.83	642.71 $\pm$ 1.79	129.68 $\pm$ 0.10	56.36 $\pm$ 0.59	82.89 $\pm$ 0.67	93.66 $\pm$ 0.04
epsGreedyRMS1	314.69 $\pm$ 14.38	202.67 $\pm$ 9.98	103.91 $\pm$ 0.63	222.10 $\pm$ 7.56	100.71 $\pm$ 0.80	106.18 $\pm$ 0.33
epsGreedyRMS3	101.44 $\pm$ 4.73	195.07 $\pm$ 11.97	118.95 $\pm$ 0.81	243.08 $\pm$ 9.62	97.42 $\pm$ 1.04	99.04 $\pm$ 0.88
RMS2b	150.48 $\pm$ 45.68	171.73 $\pm$ 40.00	112.74 $\pm$ 3.18	109.53 $\pm$ 9.00	99.17 $\pm$ 1.09	101.77 $\pm$ 0.81
BBB	223.34 $\pm$ 44.00	2200.24 $\pm$ 0.11	228.25 $\pm$ 0.20	597.48 $\pm$ 31.48	92.13 $\pm$ 1.10	114.41 $\pm$ 0.29
linGreedy	623.08 $\pm$ 73.08	1112.96 $\pm$ 58.73	134.23 $\pm$ 0.86	<b>39.29 <math>\pm</math> 5.04</b>	79.24 $\pm$ 0.53	101.25 $\pm$ 0.83
linDiagPost	782.17 $\pm$ 9.99	4470.70 $\pm$ 2.36	363.90 $\pm$ 0.07	134.48 $\pm$ 1.00	80.04 $\pm$ 0.64	121.17 $\pm$ 0.02
pNoise	93.60 $\pm$ 6.98	185.85 $\pm$ 20.13	126.71 $\pm$ 2.86	245.77 $\pm$ 13.02	96.48 $\pm$ 1.04	98.68 $\pm$ 0.69
NeuralLinear	100.00 $\pm$ 4.48	100.00 $\pm$ 0.80	100.00 $\pm$ 0.19	100.00 $\pm$ 1.12	100.00 $\pm$ 0.63	100.00 $\pm$ 0.12
GP-TS	733.40 $\pm$ 34.61	231.55 $\pm$ 40.61	155.38 $\pm$ 1.17	61.34 $\pm$ 1.04	101.59 $\pm$ 1.13	111.92 $\pm$ 0.39
BBB3	414.44 $\pm$ 112.11	2200.15 $\pm$ 0.09	229.45 $\pm$ 1.12	801.61 $\pm$ 53.58	88.05 $\pm$ 1.13	114.52 $\pm$ 0.31
SGFS	167.02 $\pm$ 7.82	247.15 $\pm$ 23.51	139.04 $\pm$ 0.44	338.55 $\pm$ 10.78	94.14 $\pm$ 0.87	115.89 $\pm$ 0.38
RMS	263.30 $\pm$ 17.65	216.91 $\pm$ 19.22	107.74 $\pm$ 1.50	231.32 $\pm$ 7.86	98.70 $\pm$ 0.86	106.59 $\pm$ 0.48
linPost	268.64 $\pm$ 31.29	666.00 $\pm$ 1.54	131.02 $\pm$ 0.07	101.74 $\pm$ 0.88	<b>78.88 <math>\pm</math> 0.65</b>	96.28 $\pm$ 0.04
linfullDiagPost	3815.50 $\pm$ 5.02	2467.19 $\pm$ 2.15	281.12 $\pm$ 0.11	98.82 $\pm$ 1.08	86.10 $\pm$ 0.66	116.92 $\pm$ 0.03
DropoutRMS2	74.95 $\pm$ 2.73	640.87 $\pm$ 62.50	126.26 $\pm$ 2.97	233.18 $\pm$ 10.56	89.60 $\pm$ 1.01	98.47 $\pm$ 0.78
epsGreedyRMS2	100.96 $\pm$ 4.52	191.58 $\pm$ 12.05	119.31 $\pm$ 0.89	258.57 $\pm$ 13.62	96.72 $\pm$ 1.06	98.19 $\pm$ 0.77
BBB2	182.82 $\pm$ 45.68	2200.13 $\pm$ 0.10	230.13 $\pm$ 1.20	724.98 $\pm$ 50.96	89.68 $\pm$ 1.17	114.42 $\pm$ 0.28
RMS2	<b>81.16 <math>\pm</math> 4.01</b>	289.97 $\pm$ 47.62	140.76 $\pm$ 5.93	289.76 $\pm$ 16.58	97.29 $\pm$ 0.92	100.11 $\pm$ 1.03
RMS3	85.29 $\pm$ 5.92	230.73 $\pm$ 25.62	137.28 $\pm$ 6.00	248.51 $\pm$ 13.22	96.21 $\pm$ 1.01	100.81 $\pm$ 1.11
BootstrappedBNNs1	175.48 $\pm$ 8.84	137.82 $\pm$ 4.71	122.08 $\pm$ 3.40	213.83 $\pm$ 9.20	101.76 $\pm$ 0.73	99.78 $\pm$ 0.66
linDiagPrecPost	424.79 $\pm$ 57.69	655.44 $\pm$ 2.09	131.12 $\pm$ 0.07	65.05 $\pm$ 0.71	78.74 $\pm$ 0.75	97.16 $\pm$ 0.05
constSGD	324.85 $\pm$ 74.27	<b>74.80 <math>\pm</math> 11.64</b>	<b>83.47 <math>\pm</math> 0.69</b>	780.33 $\pm$ 47.50	99.69 $\pm$ 0.98	107.91 $\pm$ 0.41
optimal	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00
DropoutRMS	83.19 $\pm$ 3.65	739.90 $\pm$ 68.49	127.05 $\pm$ 3.41	243.65 $\pm$ 11.65	89.37 $\pm$ 0.75	97.92 $\pm$ 0.84
unif	4399.38 $\pm$ 7.76	8721.46 $\pm$ 2.24	381.11 $\pm$ 0.05	1448.95 $\pm$ 23.20	135.41 $\pm$ 1.68	121.57 $\pm$ 0.02
BootstrappedBNNs2	98.82 $\pm$ 4.06	150.57 $\pm$ 6.38	122.30 $\pm$ 3.48	256.00 $\pm$ 16.35	101.06 $\pm$ 0.85	100.42 $\pm$ 1.01

Table 2: Simple regret incurred by the algorithms in Section 3 on the bandits described in Section A.

	Mushroom	Statlog	Covertime	Financial	Jester	Adult	Song	Census
BootstrappedBNNs2	10.91 $\pm$ 1.22	7.00 $\pm$ 0.11	37.81 $\pm$ 0.26	2.47 $\pm$ 0.18	59.30 $\pm$ 0.45	75.53 $\pm$ 0.14	88.78 $\pm$ 0.02	
linfullDiagPrecPost	4.28 $\pm$ 0.38	<b>7.31 <math>\pm</math> 0.02</b>	34.00 $\pm$ 0.03	3.94 $\pm$ 0.04	61.80 $\pm$ 0.52	77.04 $\pm$ 0.03	81.16 $\pm$ 0.02	
epsGreedyRMS1	3.87 $\pm$ 0.29	8.94 $\pm$ 0.05	37.47 $\pm$ 0.09	3.43 $\pm$ 0.12	60.50 $\pm$ 0.48	79.92 $\pm$ 0.04	99.98 $\pm$ 0.02	54.87 $\pm$ 0.26
epsGreedyRMS3	8.35 $\pm$ 0.88	7.40 $\pm$ 0.12	36.64 $\pm$ 0.22	3.17 $\pm$ 0.10	59.22 $\pm$ 0.58	<b>75.64 <math>\pm</math> 0.12</b>	88.73 $\pm$ 0.02	37.75 $\pm$ 0.25
epsGreedyRMS2	8.47 $\pm$ 0.65	7.90 $\pm$ 0.23	36.74 $\pm$ 0.21	3.20 $\pm$ 0.14	59.03 $\pm$ 0.55	75.61 $\pm$ 0.13	88.77 $\pm$ 0.02	37.81 $\pm$ 0.30
pNoise2	9.28 $\pm$ 1.03	7.91 $\pm$ 0.21		2.62 $\pm$ 0.15	60.75 $\pm$ 0.58	76.79 $\pm$ 0.34		
linfullPost	<b>2.88 <math>\pm</math> 0.47</b>	7.34 $\pm$ 0.02	<b>33.99 <math>\pm</math> 0.02</b>	5.61 $\pm$ 0.06	61.58 $\pm$ 0.48	76.74 $\pm$ 0.03	81.30 $\pm$ 0.02	
pNoise	9.11 $\pm$ 1.06	7.33 $\pm$ 0.15	37.15 $\pm$ 0.23	2.74 $\pm$ 0.21	59.15 $\pm$ 0.63	75.70 $\pm$ 0.13	88.76 $\pm$ 0.02	38.52 $\pm$ 0.26
BBB2	19.80 $\pm$ 3.27	8.19 $\pm$ 0.03	35.67 $\pm$ 0.02	7.53 $\pm$ 0.05	59.17 $\pm$ 0.45	83.36 $\pm$ 0.05	96.25 $\pm$ 0.03	
BBB3	25.79 $\pm$ 3.49	8.93 $\pm$ 0.03	36.79 $\pm$ 0.02	10.49 $\pm$ 0.09	58.89 $\pm$ 0.55	86.43 $\pm$ 0.04	96.36 $\pm$ 0.02	
SGFS	20.01 $\pm$ 1.73	16.07 $\pm$ 0.52	50.97 $\pm$ 0.41	7.96 $\pm$ 0.53	74.50 $\pm$ 0.57	96.46 $\pm$ 0.31	94.50 $\pm$ 0.21	56.52 $\pm$ 0.48
BBB	9.60 $\pm$ 1.96	7.52 $\pm$ 0.02	34.70 $\pm$ 0.02	4.96 $\pm$ 0.07	58.98 $\pm$ 0.47	79.75 $\pm$ 0.05	96.25 $\pm$ 0.03	
linGreedy	12.86 $\pm$ 1.61	11.93 $\pm$ 0.64	35.15 $\pm$ 0.11	<b>2.14 <math>\pm</math> 0.11</b>	<b>58.54 <math>\pm</math> 0.55</b>	82.57 $\pm$ 0.69	<b>80.18 <math>\pm</math> 0.04</b>	<b>30.72 <math>\pm</math> 0.02</b>
RMS2	11.24 $\pm$ 1.41	8.35 $\pm$ 0.52	37.35 $\pm$ 0.26	2.84 $\pm$ 0.20	59.68 $\pm$ 0.52	75.99 $\pm$ 0.17	88.79 $\pm$ 0.02	38.63 $\pm$ 0.35
RMS3	12.28 $\pm$ 1.77	8.13 $\pm$ 0.47	37.22 $\pm$ 0.22	2.70 $\pm$ 0.21	59.79 $\pm$ 0.71	75.94 $\pm$ 0.14	88.77 $\pm$ 0.02	38.73 $\pm$ 0.26
linDiagPrecPost	7.92 $\pm$ 0.99	7.51 $\pm$ 0.02	34.42 $\pm$ 0.03	4.65 $\pm$ 0.04	58.71 $\pm$ 0.65	79.88 $\pm$ 0.03	80.68 $\pm$ 0.03	
RMS	8.94 $\pm$ 1.43	8.64 $\pm$ 0.08	37.87 $\pm$ 0.13	2.69 $\pm$ 0.15	60.95 $\pm$ 0.53	79.87 $\pm$ 0.05	99.95 $\pm$ 0.02	54.96 $\pm$ 0.25
linPost	5.58 $\pm$ 0.74	7.61 $\pm$ 0.02	34.36 $\pm$ 0.02	7.30 $\pm$ 0.06	60.09 $\pm$ 0.48	79.15 $\pm$ 0.03	80.76 $\pm$ 0.02	
RMS2b	11.90 $\pm$ 1.32	10.14 $\pm$ 0.69	38.83 $\pm$ 0.27	2.42 $\pm$ 0.15	58.58 $\pm$ 0.53	76.45 $\pm$ 0.19	82.89 $\pm$ 0.33	
DropoutRMS2	13.07 $\pm$ 1.60	8.11 $\pm$ 0.47	37.45 $\pm$ 0.27	2.45 $\pm$ 0.16	58.91 $\pm$ 0.47	75.94 $\pm$ 0.13	88.75 $\pm$ 0.03	38.25 $\pm$ 0.32
unif	100.00 $\pm$ 0.14	100.00 $\pm$ 0.02	100.00 $\pm$ 0.02	100.00 $\pm$ 1.55	100.00 $\pm$ 1.09	100.00 $\pm$ 0.02	100.00 $\pm$ 0.02	100.00 $\pm$ 0.01
constSGD	8.37 $\pm$ 1.28	10.01 $\pm$ 0.58	36.87 $\pm$ 0.26	3.53 $\pm$ 0.18	58.72 $\pm$ 0.55	75.89 $\pm$ 0.14	103.30 $\pm$ 0.00	31.50 $\pm$ 0.05
optimal	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00
DropoutRMS	12.16 $\pm$ 1.35	7.67 $\pm$ 0.37	37.25 $\pm$ 0.26	2.51 $\pm$ 0.12	58.74 $\pm$ 0.59	75.99 $\pm$ 0.16	88.86 $\pm$ 0.04	38.43 $\pm$ 0.34

Table 3: Cumulative regret for the Wheel Bandit with increasing values of  $\delta$ .

	$\delta = 0.5$	$\delta = 0.7$	$\delta = 0.9$	$\delta = 0.95$	$\delta = 0.99$
DropoutRMS	50.91 $\pm$ 4.98	73.49 $\pm$ 4.09	87.99 $\pm$ 4.51	107.75 $\pm$ 2.88	102.53 $\pm$ 0.67
linfullPost	1.00 $\pm$ 0.01	1.28 $\pm$ 0.03	<b>3.03 <math>\pm</math> 0.08</b>	<b>5.35 <math>\pm</math> 0.18</b>	42.26 $\pm$ 2.18
NeuralLinear	3.53 $\pm$ 0.09	3.20 $\pm$ 0.31	6.81 $\pm$ 2.44	18.36 $\pm$ 4.47	77.80 $\pm$ 4.71
RMS	5.27 $\pm$ 1.64	14.84 $\pm$ 3.83	63.75 $\pm$ 6.19	103.18 $\pm$ 2.92	103.25 $\pm$ 0.55
linPost	<b>0.70 <math>\pm</math> 0.02</b>	<b>1.15 <math>\pm</math> 0.04</b>	3.20 $\pm$ 0.11	5.53 $\pm$ 0.15	<b>26.51 <math>\pm</math> 1.46</b>
RMS2b	43.98 $\pm$ 4.97	61.20 $\pm$ 4.69	100.58 $\pm$ 3.32	109.03 $\pm$ 1.90	102.96 $\pm$ 0.61
pNoise2	41.14 $\pm$ 5.11	69.27 $\pm$ 4.92	96.17 $\pm$ 3.69	108.82 $\pm$ 2.16	101.99 $\pm$ 0.78
BBB2	8.85 $\pm$ 4.27	13.49 $\pm$ 4.89	78.22 $\pm$ 7.71	90.77 $\pm$ 6.08	102.29 $\pm$ 0.85
BBB3	5.74 $\pm$ 2.73	17.14 $\pm$ 5.38	59.62 $\pm$ 8.10	104.87 $\pm$ 4.95	101.81 $\pm$ 1.05
GP-TS	2.33 $\pm$ 0.06	2.84 $\pm$ 0.24	15.07 $\pm$ 2.24	56.03 $\pm$ 2.55	95.27 $\pm$ 1.18
linfullDiagPost	1.53 $\pm$ 0.02	1.80 $\pm$ 0.03	3.53 $\pm$ 0.07	6.18 $\pm$ 0.14	35.96 $\pm$ 2.46
epsGreedyRMS3	14.14 $\pm$ 1.56	53.38 $\pm$ 2.93	70.58 $\pm$ 1.41	84.06 $\pm$ 1.64	98.22 $\pm$ 1.39
epsGreedyRMS2	14.64 $\pm$ 1.91	47.17 $\pm$ 3.28	72.71 $\pm$ 1.67	85.18 $\pm$ 1.94	95.79 $\pm$ 1.47
epsGreedyRMS1	2.72 $\pm$ 0.14	5.03 $\pm$ 0.39	18.07 $\pm$ 2.30	58.76 $\pm$ 3.08	96.51 $\pm$ 1.11
DropoutRMS2	60.32 $\pm$ 3.93	64.44 $\pm$ 4.33	98.55 $\pm$ 3.75	110.42 $\pm$ 2.63	102.59 $\pm$ 0.63
pNoise	49.38 $\pm$ 5.39	63.81 $\pm$ 3.27	90.70 $\pm$ 3.71	106.05 $\pm$ 2.56	102.86 $\pm$ 0.93
unif	100.00 $\pm$ 0.06	100.00 $\pm$ 0.08	100.00 $\pm$ 0.11	100.00 $\pm$ 0.13	100.00 $\pm$ 0.27
linDiagPrecPost	0.68 $\pm$ 0.02	1.21 $\pm$ 0.05	3.44 $\pm$ 0.11	5.60 $\pm$ 0.21	23.64 $\pm$ 1.38
linGreedy	57.14 $\pm$ 4.88	73.25 $\pm$ 3.96	104.57 $\pm$ 3.67	112.27 $\pm$ 1.87	100.47 $\pm$ 1.22
constSGD	23.49 $\pm$ 3.80	30.13 $\pm$ 5.61	64.76 $\pm$ 6.56	90.09 $\pm$ 5.90	103.13 $\pm$ 0.38
linfullDiagPrecPost	0.95 $\pm$ 0.02	1.33 $\pm$ 0.04	3.35 $\pm$ 0.08	5.71 $\pm$ 0.16	33.38 $\pm$ 2.15
BootstrappedBNNs2	45.65 $\pm$ 4.20	71.45 $\pm$ 3.23	89.37 $\pm$ 3.56	103.60 $\pm$ 2.56	100.86 $\pm$ 1.05
BootstrappedBNNs1	40.41 $\pm$ 4.08	66.97 $\pm$ 3.94	94.75 $\pm$ 3.10	105.75 $\pm$ 2.54	103.53 $\pm$ 0.69
linDiagPost	0.70 $\pm$ 0.02	1.18 $\pm$ 0.04	3.57 $\pm$ 0.12	5.63 $\pm$ 0.16	20.07 $\pm$ 0.78
RMS2	52.83 $\pm$ 3.90	79.47 $\pm$ 4.06	89.93 $\pm$ 3.44	109.57 $\pm$ 2.61	103.84 $\pm$ 0.46
RMS3	50.95 $\pm$ 3.93	67.98 $\pm$ 3.46	97.76 $\pm$ 3.38	105.98 $\pm$ 3.03	101.71 $\pm$ 0.97
optimal	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00
SGFS	13.33 $\pm$ 1.93	19.15 $\pm$ 2.75	33.23 $\pm$ 2.57	37.86 $\pm$ 1.57	49.45 $\pm$ 1.84
BBB	11.65 $\pm$ 4.44	8.96 $\pm$ 2.96	69.91 $\pm$ 7.61	84.01 $\pm$ 6.56	103.27 $\pm$ 0.63



## A REAL-WORLD DATASETS

**Mushroom.** The Mushroom Dataset (Schlimmer, 1981) contains 22 attributes per mushroom, and two classes: poisonous and safe. As in Blundell et al. (2015), we create a bandit problem where the agent must decide whether to eat or not a given mushroom. Eating a safe mushroom provides reward +5. Eating a poisonous mushroom delivers reward +5 with probability 1/2 and reward -35 otherwise. If the agent does not eat a mushroom, then the reward is 0. We set  $n = 50000$ .

**Statlog.** The Shuttle Statlog Dataset (Asuncion & Newman, 2007) provides the value of  $d = 9$  indicators during a space shuttle flight, and the goal is to predict the state of the radiator subsystem of the shuttle. There are  $k = 7$  possible states, and if the agent selects the right state, then reward 1 is generated. Otherwise, the agent obtains no reward ( $r = 0$ ). The most interesting aspect of the dataset is that one action is the optimal one in 80% of the cases, and some algorithms may commit to this action instead of further exploring. In this case,  $n = 43500$ .

**Coverttype.** The Coverttype Dataset (Asuncion & Newman, 2007) classifies the cover type of northern Colorado forest areas in  $k = 7$  classes, based on  $d = 54$  features, including elevation, slope, aspect, and soil type. Again, the agent obtains reward 1 if the correct class is selected, and 0 otherwise. We run the bandit for  $n = 150000$ .

**Financial.** We created the Financial Dataset by pulling the stock prices of  $d = 21$  publicly traded companies in NYSE and Nasdaq, for the last 14 years ( $n = 3713$ ). For each day, the context was the price difference between the beginning and end of the session for each stock. We synthetically created the arms, to be a linear combination of the contexts, representing  $k = 8$  different potential portfolios. By far, this was the smallest dataset, and many algorithms over-explored at the beginning with no time to amortize their investment (Thompson Sampling does not account for the horizon).

**Jester.** We create a recommendation system bandit problem as follows. The Jester Dataset (Goldberg et al., 2001) provides continuous ratings in  $[-10, 10]$  for 100 jokes from 73421 users. We find a complete subset of  $n = 19181$  users rating all 40 jokes. Following Riquelme et al. (2017), we take  $d = 32$  of the ratings as the context of the user, and  $k = 8$  as the arms. The agent recommends one joke, and obtains the reward corresponding to the rating of the user for the selected joke.

**Adult.** The Adult Dataset (Kohavi; Asuncion & Newman, 2007) comprises personal information from the US Census Bureau database, and the standard prediction task is to determine if a person makes over \$50K a year or not. However, we consider the  $k = 14$  different occupations as feasible actions, based on  $d = 94$  covariates (many of them binarized). As in previous datasets, the agent obtains reward 1 for making the right prediction, and 0 otherwise. We set  $n = 45222$ .

**Census.** The US Census (1990) Dataset (Asuncion & Newman, 2007) contains a number of personal features (age, native language, education...) which we summarize in  $d = 389$  covariates, including binary dummy variables for categorical features. Our goal again is to predict the occupation of the individual among  $k = 9$  classes. The agent obtains reward 1 for making the right prediction, and 0 otherwise, for each of the  $n = 250000$  randomly selected data points.

**Song.** The YearPredictionMSD Dataset is a subset of the Million Song Dataset (Bertin-Mahieux et al., 2011). The goal is to predict the year a given song was released (1922-2011) based on  $d = 90$  technical audio features. We divided the years in  $k = 10$  contiguous year buckets containing the same number of songs, and provided decreasing Gaussian rewards as a function of the distance between the interval chosen by the agent and the one containing the year the song was actually released. We initially selected  $n = 250000$  songs at random from the training set.

The Statlog, Coverttype, Adult, and Census datasets were tested in Elmachet et al. (2017).