

1 We thank the reviewers for their time and effort reviewing our paper. We are pleased that you found our work to “solve
2 an important problem” (R1), our method to be “elegant” (R2), and our paper to be “well written” (R3). Below we
3 respond to the key comments and issues in order of appearance.

4 **R1: “Diverse mini-batch Active Learning” and open-source code** We will discuss the paper and consider it as an
5 additional baseline. At the latest, source code will be made publicly available upon publication.

6 **R2: Batch active learning by imputing labels** When we write “greedy,” we are referring to the sequential setting
7 in which we alternate between querying a single data point and updating the model. When we say “naive batch” we
8 mean ranking and taking the top b points. The proposed imputation approach could be considered as an alternative
9 batch construction method not currently discussed in the paper. The key issue with this method is indeed the cost of
10 performing a model update after each label is imputed. While a linear increase in runtime due to repeated model updates
11 may be acceptable for constructing smaller-sized batches, in the large-scale settings we consider it simply becomes
12 infeasible. Interestingly, there is a close relationship between the imputation approach and our method: in fact, our
13 method also uses the expected (i.e., imputed) labels under the model to construct the batch in a principled way; however,
14 importantly, it does not require making repeated predictions over the entire pool set or updating the model after each
15 point is added to the batch. We will add a thorough discussion on this topic to a future version of the paper.

16 **R2: Logistic vs. probit regression** We agree with R2 that the discussion on logistic regression is clearer by considering
17 probit regression to begin with. We will update the section accordingly.

18 **R2: Extending entropy-based methods to the batch setting is not necessarily difficult** We agree the wording is
19 misleading. The main issue is how to do this efficiently for complex, non-linear models. We will clarify this point.

20 **R2: Runtime discussion** In general, constructing the batches has negligible cost (cf. Section 5 for a discussion on
21 computational complexity) compared to updating the model, so increasing the batch size allows to decrease overall
22 computational cost. We will add a more detailed comparison with other methods to the experimental section.

23 **R3: “Nonmyopic active learning of Gaussian processes: an exploration-exploitation approach”** While GPs suffer
24 from scalability issues and cannot be applied to many of the domains we considered, we agree that batch active learning
25 (AL) with GPs is under-discussed in the paper. Krause & Guestrin (2007) consider mutual information (MI; also
26 known as BALD) as an acquisition criterion, which makes this work related to ours. However, they immediately reduce
27 the batch formulation to the sequential greedy case. We highlight similarities and differences to sequential greedy
28 MI/BALD in Sections 4 and 6, and will add a discussion on extending sequential greedy methods to the batch setting
29 (cf. *R2: Extending entropy-based methods to the batch setting is not necessarily difficult*). Less relevant for our work
30 are GP-specific details (as we focus on BNNs) and proofs relying on submodularity (as the diminishing returns property
31 does not necessarily hold when performing approximate inference and stochastic optimization). We will expand the
32 related work section accordingly.

33 **R3: Clarifying first equality in eq. (4)** The first equality is achieved by applying Bayes’ rule to the posterior (eq. (1)
34 in the main paper), taking the logarithm, and applying linearity of expectations:

$$\begin{aligned} \mathbb{E}_{\mathcal{Y}_p|\mathcal{X}_p, \mathcal{D}_0} [\log p(\boldsymbol{\theta}|\mathcal{D}_0 \cup (\mathcal{X}_p, \mathcal{Y}_p))] &= \mathbb{E}_{\mathcal{Y}_p|\mathcal{X}_p, \mathcal{D}_0} [\log p(\boldsymbol{\theta}|\mathcal{D}_0) + \log p(\mathcal{Y}_p|\mathcal{X}_p, \boldsymbol{\theta}) - \log p(\mathcal{Y}_p|\mathcal{X}_p, \mathcal{D}_0)] \\ &= \log p(\boldsymbol{\theta}|\mathcal{D}_0) + \mathbb{E}_{\mathcal{Y}_p|\mathcal{X}_p, \mathcal{D}_0} [\log p(\mathcal{Y}_p|\mathcal{X}_p, \boldsymbol{\theta})] + \mathbb{H}[\mathcal{Y}_p|\mathcal{X}_p, \mathcal{D}_0], \end{aligned} \quad (1)$$

35 where we used $\mathbb{E}_{\mathcal{Y}_p} [-\log p(\mathcal{Y}_p|\mathcal{X}_p, \mathcal{D}_0)] = \mathbb{H}[\mathcal{Y}_p|\mathcal{X}_p, \mathcal{D}_0]$. We will make this derivation clearer in the next version.

36 **R3: Motivation for the batch AL setting can be improved** This scenario is practical in a number of real-world
37 applications, particularly when the cost of acquiring labels is high but can be parallelized. Examples include crowd-
38 sourcing a complex labeling task, leveraging parallel simulations on a compute cluster, or performing experiments that
39 require resources with time-limited availability (e.g. a wet-lab in natural sciences). In all these cases, being able to
40 generate a quality batch of query points without having to wait for previous labels to be acquired can be extremely
41 advantageous. The importance of the scenario is further evidenced by the volume of existing literature and ongoing
42 research on this topic. See for example Hoi et al. (2006), Guo & Schuurmans (2008), Wei et al. (2015), Sener &
43 Savarese (2018), Kirsch et al. (2019), and many more references therein, all of which are concerned with the batch AL
44 setting. We thank R3 for making this point, and will add content motivating the batch AL setting to the paper.

45 **R3: Evaluated datasets** As a methodology paper, our goal is to demonstrate the usefulness of our proposed method in
46 a broad range of scenarios. We performed experiments on several small- and large-scale regression and classification
47 datasets. While we agree with R3 that the datasets evaluated in the experimental section do not necessarily reflect
48 real-world scenarios, the experimental protocols we used resemble benchmarks from multiple important (batch) AL
49 papers (e.g., Hernandez-Lobato & Adams, 2015; Gal et al., 2017; Sener & Savarese, 2018). In fact, our work goes
50 beyond what is typically tractable with Bayesian approaches (e.g., *cifar10* and *year*), demonstrating the usefulness and
51 scalability of our method. Since the performance of the method is independent of labelling cost, we argue it is sufficient
52 to use real-world settings as motivating examples (see discussion above) and leave specific applications to future work.