

1 We thank all reviewers for their insightful comments. First, we will address some major themes from the reviews:

2 **Non-probabilistic Methods & Constrained Optimization** Our paper is motivated by the ability to impose con-  
3 straints probabilistically, i.e. the challenge is to incorporate constraints in a model that infers output *distributions*.  
4 Bayesian techniques deviate significantly from classical optimization. As such, we did not include non-probabilistic  
5 methods like [Stewart and Ermon, 2017] as baselines, or use notation conventional to constraint optimization.

6 **Hard vs. Soft Constraints** Figure 3 shows minor violations of constraints as the priors used are soft and assign  
7 small (but  $> 0$ ) probability to infractions. The idiosyncratic nature of SVGD inference (used in Figure 3b) in learning  
8 “repulsive/diverse functions” also makes it more likely to violate a constraint; for example, Figure S1 below shows that  
9 for the same example, having only 10 SVGD particles eliminate the few violating functions.

10 More generally, in probabilistic systems, while hard constraints are theoretically  
11 possible by assigning 0 probability to violations, (i) numerical instability issues  
12 could arise, and (ii) we tend to obey Cromwell’s rule in Bayesian inference, where  
13 the support of the prior is usually the entire output space and unlikely functions are  
14 naturally weeded out. **We stress that workarounds do exist:** (i) we can specify  
15 extremely small ( $\approx 0$ ) probability to the order of numerical insignificance, so  
16 long as the prior remains differentiable, (ii) guarantees *on top of* soft constraints  
17 can be enforced, e.g. rejection sampling over OC-BNNs (which will be tractable),  
18 (iii) in the amortized setting, we can set a threshold for  $\epsilon$  directly. We note that  
19 soft constraints are often useful too, e.g. for learning (alongside the training data)  
20 where the function might be *outside of* the constrained region.

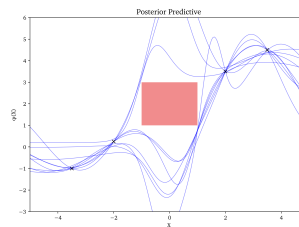


Figure S1: Same example as Figure 3b, except using 10 SVGD particles (instead of 50).

21 **Novelty (1)** While regularization techniques are common, it is not immediately  
22 clear that data-based regularization leads to “well-behaved” and useful priors, (e.g.  
23 smooth functions with suitable OOD variance), especially for the amortized variant, or for non-Gaussian likelihoods  
24 (e.g. constraints over output ranges). Tractability of sampling with input dimensionality is also not obvious (and not  
25 demonstrated by [18]), for example, we found that sampling at the border of constraints proved reasonably well at  
26 guiding the model towards good functions. **(2)** We acknowledge points about similarity to [18] made by R4. A more  
27 accurate comparison would be that our framework is more general and more versatile at incorporating a diverse range  
28 of constraint formulations, without the need to make various Gaussian approximations or sacrifice tractability. **(3)** We  
29 want to highlight the strength and novelty of our suite of experiments, which shows that OC-BNNs are useful and work  
30 well on a diverse set of real-life problems and constraints.

### 31 Additional Comments

32 **[All] Technicality:** The stochastic process setup in Section 4 is to ensure a formal and principled definition keeping  
33 with Bayesian inference, not to deceive the reader into unnecessary complexity. We acknowledge that a more intuitive  
34 explanation and/or an explicit algorithm box might suffice and technical details be left to supplementary material.

35 **[R1][R3] Def. 4.2:** Indeed, we omitted marginalizing  $y$ . Equation (1) in Definition 4.2 should read:

$$p(Y \circ C_y(\mathbf{x})|\mathbf{x}) = \int_{\mathcal{Y}} \mathbb{I}[y \circ C_y(\mathbf{x})] \underbrace{\int_{\mathcal{W}} p(Y = y|\mathbf{x}, \mathbf{w}) p(\mathbf{w}) d\mathbf{w}}_{\text{prior predictive}} dy \leq \epsilon$$

36 Also, as R1 pointed out,  $\circ$  should be swapped here: for a positive constraint,  $\circ = \neq$  (and vice versa).

37 **[R3] Section 3:** We optimize w.r.t. the parameters of a Gaussian variational representation, hence “variational”.

38 **[R1][R4] Fig. 2:** The imperfect fit is due to an idiosyncratic combination of low model capacity and sampling for this  
39 particular example. Note that the plot was shaded at a specific confidence level; it is challenging for a 10-node RBF  
40 network to fit a specific rectangle at *identical levels of confidence*. A key takeaway from Figure 2 is that we are not  
41 overly confident far away from the green rectangle, especially in the prior predictive.

42 **[R4] Section 4.1:** There is a one-to-one correspondence: each constraint  $(C_x, C_y, \circ)$  is modeled with a *single* stochastic  
43 process, whereby the points of the input region  $C_x$  is the index of the process. (The equation on Line 135, on the other  
44 hand, represents the product of multiple, independent constraints.) The confusion here may arise from the fact that the  
45 points within a single constraint are also “independent” as their correct output has been directly defined by  $C_y$ .

46 Russell Stewart and Stefano Ermon. Label-free Supervision of Neural Networks with Physics and Domain Knowledge.  
47 In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.