

---

# Dynamic Tempering (need a better name) for Asymptotic Uniform Confidence in Bayesian ReLU Networks

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 ReLU networks, in conjunction with maximum a posteriori estimation, produce  
2 overconfidence predictions far away from the training data. Recently, it has been  
3 shown that cheap Gaussian-approximated Bayesian methods can mitigate this  
4 issue to some degree. However, asymptotically, the confidence of these Bayesian  
5 ReLU networks only converges to a constant, which can be loose depending on  
6 the network. We propose a simple heuristic, based on the tempering framework,  
7 that is applied on the distribution over logits of these networks, where we let  
8 the temperature parameter depends on the input. This method gives a stronger  
9 asymptotic confidence guarantee compared to the previous work: Far away from  
10 the training data, the confidence decays proportionally to the distance and uniform  
11 at the limit. Meanwhile, we also show this method does not change the decision  
12 boundary of the underlying network and the confidence near the training data is  
13 still close its original value.

## 14 1 New Stuff added here for easy reference

15 Consider WLOG, a three-layer binary ReLU classification network:

$$f_{\theta}(x) := w_2^{\top} r(W_1 r(W_0 x + b_0) + b_1) + b_2$$

16 where

$$\theta := \{W_0, W_1, w_2, b_0, b_1, b_2\} \in \mathbb{R}^P,$$

17 with

$$W_0 \in \mathbb{R}^{N_1 \times N_0}, W_1 \in \mathbb{R}^{N_2 \times N_1}, w_2 \in \mathbb{R}^{N_2}, \\ b_0 \in \mathbb{R}^{N_1}, b_1 \in \mathbb{R}^{N_2}, b_2 \in \mathbb{R},$$

18 and  $r$  is the component-wise ReLU function. Given a dataset  $\mathcal{D}$ , we consider the loss function  
19  $\ell : \mathbb{R}^P \rightarrow \mathbb{R}$ .

20 The idea is to add additional *inactive* ReLU units, as follows. Define  $\tilde{N}_1 := N_1 + M_1$  and  $\tilde{N}_2 :=$   
21  $N_2 + M_2$  for some  $M_1, M_2 \in \mathbb{N}$ . Then we define

$$\tilde{\theta} := \{\tilde{W}_0, \tilde{W}_1, \tilde{w}_2, \tilde{b}_0, \tilde{b}_1, b_2\} \in \mathbb{R}^{\tilde{P}},$$

22 with

$$\tilde{W}_0 \in \mathbb{R}^{\tilde{N}_1 \times N_0}, \tilde{W}_1 \in \mathbb{R}^{\tilde{N}_2 \times \tilde{N}_1}, \tilde{w}_2 \in \mathbb{R}^{\tilde{N}_2}, \\ \tilde{b}_0 \in \mathbb{R}^{\tilde{N}_1}, \tilde{b}_1 \in \mathbb{R}^{\tilde{N}_2}, b_2 \in \mathbb{R},$$

23 where for each  $i \neq 0$ , we append additional  $M_{i-1}$  columns containing zeros to  $W_i$ . In other words, we  
 24 have that  $(W_i)_{jk} = 0$  if  $k > N_{i-1}$ . Furthermore, the (regularized) loss function is now  $\tilde{\ell} : \mathbb{R}^{\tilde{P}} \rightarrow \mathbb{R}$ .

25 We can confirm that, under this construction, for any  $x$ , we have that  $f_{\tilde{\theta}}(x) = f_{\theta}(x)$ . Moreover, for  
 26 each  $i \neq 2$  the gradient of the last  $M_{i+1}$  rows of  $\tilde{W}_i$  is zero. Thus, we can choose these particular  
 27 components to ‘‘augment’’ the uncertainty of the classifier, without changing the prediction.

28 **Proposition 1** (Uncertainty Guarantee). *Suppose  $f_{\theta}, f_{\tilde{\theta}}$  are as defined above and last-layer Laplace*  
 29 *approximations are used to obtain the respective posterior  $\mathcal{N}(w_2 \mid \mu, \Sigma), \mathcal{N}(\tilde{w}_2 \mid \tilde{\mu}, \tilde{\Sigma})$ . Then, for*  
 30 *any input  $x$  and any  $M \in \mathbb{N}$  additional ReLU units, the approximate-Gaussian variance of  $f_{\tilde{\theta}}(x)$  is*  
 31 *strictly larger than that of  $f_{\theta}(x)$ .*

32 *Proof sketch.* The variance of  $f_{\theta}(x)$  is  $\phi(x)^{\top} \Sigma \phi(x)$ . Via an eigendecomposition, it can be writ-  
 33 ten as  $\sum_i^P \lambda_i(\Sigma) (Q\phi(x))_i^2 =: v$ . Using the same argument, the variance of  $f_{\tilde{\theta}}(x)$  is given by  
 34  $\sum_i^{\tilde{P}} \lambda_i(\tilde{\Sigma}) (\tilde{Q}\tilde{\phi}(x))_i^2 =: \tilde{v}$  where  $\tilde{P} > P$ . Note that  $\Sigma, \tilde{\Sigma}$  are SPD, thus the summands are all positive.  
 35 Therefore  $\tilde{v} > v$ , because  $\tilde{v}$  has more summand than  $v$ .  $\square$

### 36 1.1 Possible geometric interpretation

37 The graph of the loss  $\ell$  on Euclidean space forms an embedded submanifold of codimension one, i.e.  
 38 a hypersurface, in  $\mathbb{R}^P \times \mathbb{R} \simeq \mathbb{R}^{P+1}$ . The the graph of  $\tilde{\ell}$  is a hypersurface in the higher dimensional  
 39 space  $\mathbb{R}^{\tilde{P}+1}$ , but restricted to the subset  $\mathbb{R}^{P+1} \subset \mathbb{R}^{\tilde{P}+1}$ , it is exactly the loss landscape of  $\ell$ .

40 Suppose we use a Laplace approximation with posterior mean  $\mu$  (a local mode of the loss landscape).  
 41 For each  $x$ , the variance  $v(x)$  of the prediction  $f_{\theta}(x)$  can be approximated with the quadratic form  
 42  $g(x)^{\top} H^{-1} g(x)$ , where  $g(x)$  is the gradient of  $f_{\theta}$  w.r.t.  $\theta$  at  $\mu$  and  $H$  is the Hessian of the loss w.r.t.  
 43  $\theta$  at  $\mu$ . This means that, in the eigenbasis  $Q$  of  $H^{-1}$ , we can write the variance as

$$v(x) \approx \sum_{i=1}^P \frac{1}{\lambda_i(H)} (Qg(x))_i^2.$$

44 The key point here is that since  $\mu$  is a stationary point,  $\lambda_i(H)$  is the *principal curvatures*, since the  
 45 Hessian matrix  $H$  in this case is the same as the second fundamental form’s matrix representation.  
 46 The principal curvatures are extrinsic, thus depend on the embedding of the submanifold. We can we  
 47 think that ReLUQ is a choice of embedding of the original loss landscape onto  $\mathbb{R}^{\tilde{P}+1}$ . The goal of  
 48 this method is therefore to find such embedding, modifying the principal curvatures, s.t. the variance  
 49 has the desired properties.

## 50 2 Introduction

51 Meaningful uncertainty estimates on the predictions of deep networks are important in mission-  
 52 critical applications such as self-driving vehicles. Without them, deep networks are susceptible  
 53 to out-of-distribution data since they are known to be overconfident: except around the decision  
 54 boundary, their predictive confidence is arbitrarily close to one. Specifically for ReLU networks,  
 55 this overconfidence problem has been studied recently by Nguyen et al. [2015], Hein et al. [2019].  
 56 It can be shown that far away from the training data, the prediction of ReLU networks converges  
 57 to one. One can therefore cannot trust the confidence estimate that ReLU networks output. This is  
 58 problematic since this means that these networks cannot even detect such obvious out-of-distribution  
 59 points.

60 Recently, it has been shown that by employing Bayesian methods on ReLU networks, this issue  
 61 can be mitigated. Specifically, Kristiadi et al. [2020] showed that applying Gaussian-approximated  
 62 Bayesian methods onto pre-trained ReLU networks, e.g. via a simple Laplace approximation, could  
 63 make their confidence asymptotically constant. However, their result does not necessarily imply  
 64 that the confidence is asymptotically uniform, which is the ideal outcome. Furthermore, their bound  
 65 depends on the spectrum of the covariance matrix of the approximate posterior. Thus, depending on  
 66 the network, this bound can be loose.

67 In this paper, we propose a simple heuristic method to fix the aforementioned problems in Gaussian-  
68 approximated Bayesian ReLU networks. Our method is based on the tempering technique that is  
69 ubiquitously used in Bayesian inference [Wenzel et al., 2020, etc.]. However, instead of tempering  
70 the posterior with a constant temperature, we temper the *distribution over logits* with a temperature  
71 that depends on the input. We therefore call our method “dynamic predictive tempering” or DPT  
72 in short. We show that DPT yields (i) confidence arbitrarily close to 0.5 far away from the training  
73 data and (ii) uniform confidence in the limit. Furthermore, we show that, DPT does not change the  
74 underlying decision boundary of vanilla Bayesian prediction, thus maintain the original accuracy of  
75 the network.

76 We begin this paper with an exposition to the overconfidence problem of ReLU networks and how  
77 can Bayesian methods mitigate this in Section 3. In Section 4 we introduce DPT and analyze its  
78 properties. Related work will be discussed in Section 5 while empirical results in Section 6.

### 79 3 Bayesian ReLU networks

80 We call a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}^k$  piecewise affine if there is a decomposition of the input space  $\mathbb{R}^n$   
81 by a finite set of polytopes  $\{P_1, \dots, P_r\}$  such that the restriction  $f|_{P_i}$  is affine for each  $i = 1 \dots r$ .  
82 These polytopes are referred to as the linear regions of  $f$ . ReLU networks are neural networks that  
83 results in piecewise affine function on the input space [Arora et al., 2018]. This class of networks  
84 includes any neural network that is composed by fully-connected and convolution layers, along with  
85 ReLU or leaky-ReLU activation function and max- or average-pooling. Commonly, these networks  
86 are trained via a maximum a posteriori (MAP) estimation. Thus in the rest of this paper, unless stated  
87 otherwise, by ReLU networks we refer to MAP-trained ReLU networks.

88 ReLU networks are arguably the most widely used deep architecture due to its high accuracy. However,  
89 they are not suitable for answering any problem beyond decision problems that requires accurate  
90 uncertainty estimates, since they tend to be overconfident. Hein et al. [2019] showed that, given a  
91 ReLU network  $f_\theta : \mathbb{R}^n \rightarrow \mathbb{R}^k$ , for almost any input point  $x \in \mathbb{R}^n$ , as  $\alpha > 0$  goes to infinity, the  
92 confidence

$$\max_{i \in \{1, \dots, k\}} \text{softmax}(f_\theta(\alpha x), i) := \max_{i \in \{1, \dots, k\}} \frac{\exp f_\theta^i(\alpha x)}{\sum_{j=1}^k \exp f_\theta^j(\alpha x)}$$

93 of  $\alpha x$  goes to one. In most realistic problems, even a finite  $\alpha$  could already yield such obvious  
94 out-of-distribution data, for example in natural image classification, where the in-distribution data is  
95 fully contained in the unit box  $[0, 1]^n \subset \mathbb{R}^n$ . This problem shows that the uncertainty estimates of  
96 ReLU networks cannot be trusted and are useless for detecting out-of-distribution data far away from  
97 the training data.

98 In what follows, we focus on the binary classification case due to the lack of analytic approximation  
99 of the integral of softmax function w.r.t. a Gaussian measure. Note that, in this case, the confidence is  
100 defined by

$$\text{conf}(x) := \max_{i \in \{0, 1\}} \sigma(f_\theta(x)) = \sigma(|f_\theta(x)|). \quad (1)$$

101 The second equality follows from Proposition 123 in the appendix. We will revisit the multi-class  
102 case in the empirical analysis.

103 Recently, Kristiadi et al. [2020] argued that this overconfidence problem is due to the usage of  
104 point estimate for making predictions in a ReLU network. Approximate Bayesian methods with  
105 Gaussian posterior are therefore proposed to mitigate this problem. Let  $f_\theta : \mathbb{R}^n \rightarrow \mathbb{R}$  be a binary  
106 ReLU classifier,  $p(y = 1|x, \theta) := \sigma(f_\theta(x))$  be the likelihood<sup>1</sup>, and  $p(\theta|\mathcal{D}) \approx \mathcal{N}(\theta|\mu, \Sigma)$  be the  
107 approximate posterior, obtained via e.g. a Laplace approximation or variational Bayes. Given an  
108 input  $x$ , the distribution over the logits  $f_\theta(x)$  can further be approximated via a linearization  $f_\theta(x)$ ,  
109 seen as a function of  $\theta$  [MacKay, 1995]. That is, if we let  $g(x) := \nabla f_\theta(x)|_\mu$  to be the gradient of  
110  $f_\theta(x)$  w.r.t.  $\theta$  at  $\mu$ , we have the following first-order Taylor’s expansion:

$$f_\theta(x) \approx f_\mu(x) + g(x)^\top (\theta - \mu),$$

111 and thus we have

$$p(f_\theta(x)|x, \mathcal{D}) \approx \mathcal{N}(f_\theta(x)|f_\mu(x), g(x)^\top \Sigma g(x)). \quad (2)$$

<sup>1</sup> $\sigma(z) := 1/(1 + \exp(-z))$  is the logistic function.

112 The predictive distribution is given by the following integral

$$p(y = 1|x, \mathcal{D}) = \int \sigma(f_\theta(x))p(f_\theta(x)|x, \mathcal{D}) df_\theta(x), \quad (3)$$

113 which can be approximated via the convolution of probit function [MacKay, 1992], which gives us

$$p(y = 1|x, \mathcal{D}) \approx \sigma\left(\frac{f_\mu(x)}{(1 + \pi/8 g(x)^\top \Sigma g(x))^{\frac{1}{2}}}\right) =: \sigma(z(x)). \quad (4)$$

114 Using this framework, Kristiadi et al. [2020] analytically showed that the marginalization in (3)  
 115 gives for any input  $x$ , confidence that asymptotically converges to a constant. Furthermore, they  
 116 show empirically that for popular ReLU networks like LeNets or ResNets, the confidence already  
 117 converges in finite distance regime, i.e. when  $\alpha < \infty$ , in both binary and multi-class cases.

## 118 4 Dynamic tempering

119 The analysis of Kristiadi et al. [2020] reveals that, popular approximate Bayesian methods could  
 120 mitigate the overconfidence problem of ReLU networks to *some degree* since their proofs disregard  
 121 the distance between an input  $x$  and the training dataset  $\mathcal{D}$ . Furthermore, the confidence is only  
 122 asymptotically constant, depending on the minimum eigenvalue of the posterior covariance  $\Sigma$ . In  
 123 some networks, it might be the case that this eigenvalue is too small such that the bound becomes  
 124 loose. Based on this problem, we formulate the following two desiderata: (i) far away from the  
 125 training data, the confidence should be  $\epsilon$ -close to 0.5, and (ii) it should be uniform in the limit, not  
 126 merely a constant. Fulfilling these desiderata—thus “fixing” Bayesian methods—is our principal  
 127 focus in this section.

128 To that end, we propose a simple heuristic based on the tempering framework. However, whereas  
 129 usually the tempering is done to the posterior  $p(\theta|\mathcal{D})$  via some constant parameter  $t$ , we propose to  
 130 temper the distribution over the logits (2) given an input  $x$  with a “dynamic” parameter  $t(x)$  that  
 131 depends on  $x$  instead. That is, we follow the construction of the Bayesian ReLU networks in the  
 132 previous section, but we use the following tempered distribution instead of (2):

$$\begin{aligned} p_t(f_\theta(x)|x, \mathcal{D}) &: \propto p(f_\theta(x)|x, \mathcal{D})^{t(x)} \\ &\propto \exp\left(-\frac{1}{2g(x)^\top \Sigma g(x)} |f_\theta(x), f_\mu(x)|^2\right)^{t(x)} \\ &= \mathcal{N}(f_\theta(x)|f_\mu(x), t(x)^{-1} \text{var}(f_\theta(x))) \end{aligned} \quad (5)$$

133 where  $t$  is a function  $\mathbb{R}^n \rightarrow \mathbb{R}$  and we have used  $\text{var}(f_\theta(x)) := g(x)^\top \Sigma g(x)$  (cf. (2)). It is easy to  
 134 see that (4) now becomes

$$p_t(y = 1|x, \mathcal{D}) \approx \sigma\left(\frac{f_\mu(x)}{(1 + \pi/8 t(x)^{-1} g(x)^\top \Sigma g(x))^{\frac{1}{2}}}\right) =: \sigma(z_t(x)). \quad (6)$$

135 Note that the difference between (3) and (6) amounts to the additional scalar term  $1/t(x)$  in the  
 136 denominator. One can then substitute  $\sigma(f_\theta(x))$  in (1) with (6) to obtain the tempered Bayesian  
 137 confidence:

$$\text{conf}_t(x) := \sigma(|z_t(x)|). \quad (7)$$

138 How can we define  $t$  such that we fulfill our desiderata? Since we would like to capture the distance  
 139 between an input to the training dataset, one option is to use the information given by the data density  
 140  $p(x)$  approximated by the training data  $\mathcal{D}$ . Let  $t(x) \propto p(x)$  be proportional to  $p(x)$ . Then, it is clear  
 141 that the value  $t(x)$  will be high if  $x$  is close to  $\mathcal{D}$ , since  $\mathcal{D}$  are assumed to be sampled from  $p(x)$ .  
 142 Conversely, if  $x$  is far away from  $\mathcal{D}$ , then  $t(x)$  will be close to zero since it lies on the low-density  
 143 region of  $p(x)$ .

144 More concretely, let us assume that the density  $p(x)$  is approximated with the kernel density estimator  
 145 (KDE) with a Gaussian kernel with length-scale  $\ell$ :

$$p(x) \approx \frac{1}{m} \sum_{i=1}^m \frac{1}{\sqrt{2\pi\ell^2}} \exp\left(-\frac{1}{2\ell^2} d(x, x_i)^2\right),$$

146 where  $d(x, x_i) := \|x - x_i\|_2$  is the Euclidean metric. We can then define  $t : \mathbb{R}^n \rightarrow (0, \infty)$  as

$$t(x) := \sum_{i=1}^m \exp\left(-\frac{1}{2\ell^2} d(x, x_i)^2\right), \quad (8)$$

147 since the only requirement is that  $t \propto p$ . This construction has the following properties.

148 **Lemma 1.** *Let  $t$  be defined in (8), and  $x \in \mathbb{R}^n$  be arbitrary. Then, it holds that  $\lim_{\delta \rightarrow \infty} t(\delta x) = 0$ .*

149 *Proof.* First, we note that  $\lim_{\delta \rightarrow \infty} d(\delta x, x_i)^2 = \infty$  for each  $i = 1, \dots, m$  and that the exponent is  
150 continuous everywhere. Therefore we can write

$$\begin{aligned} \lim_{\delta \rightarrow \infty} t(\delta x) &= \lim_{\delta \rightarrow \infty} \sum_{i=1}^m \exp\left(-\frac{1}{2\ell^2} d(x, x_i)^2\right) \\ &= \sum_{i=1}^m \exp\left(-\frac{1}{2\ell^2} \lim_{\delta \rightarrow \infty} d(x, x_i)^2\right) \\ &= \sum_{i=1}^m \exp(-\infty) \\ &= \sum_{i=1}^m 0 \\ &= 0, \end{aligned}$$

151 thus ending the proof.  $\square$

152 **Lemma 2** (Hein et al., 2019). *Let  $\{P_i\}_{i=1}^r$  be the set of linear regions associated to the binary ReLU  
153 classifier  $f_\theta : \mathbb{R}^n \rightarrow \mathbb{R}$ . For any  $x \in \mathbb{R}^n$  there exists  $\beta > 0$  and  $j \in \{1, \dots, r\}$  such that  $\delta x \in P_j$   
154 for all  $\delta \geq \beta$ . Furthermore, the restriction  $f_\theta|_{P_j}$  of  $f_\theta$  to  $P_j$  can be written as the affine function:*

$$f_\theta|_{P_j}(\hat{x}) = w(\theta)^\top \hat{x} + b(\theta),$$

155 for all  $\hat{x} \in P_j$ , where  $w(\theta) \in \mathbb{R}^n$  and  $b(\theta) \in \mathbb{R}$  are some vector and scalar that depend on  $\theta$ .  $\square$

156 We shall refer to the polytope  $P_j$  given by Lemma 2 as the *outer polytope of  $f_\theta$  w.r.t.  $x$* . We are now  
157 ready to state our main theorem, which states that given a ReLU network  $f_\theta$  and any test point  $x$ , at  
158 the corresponding outer polytope, the confidence is asymptotically low. Furthermore, the confidence  
159 is uniform at the limit.

160 The following assumes that we don't have biases. Need to show similar results with biases.

161 **Theorem 1.** *Let  $f_\theta : \mathbb{R}^n \rightarrow \mathbb{R}$  be a binary ReLU classifier without bias vectors,  $\mathcal{N}(\theta|\mu, \Sigma)$  be the  
162 posterior, and let  $t : \mathbb{R}^n \rightarrow (0, \infty)$  be defined in (8). Then for any test point  $x \in \mathbb{R}^n$  and arbitrary  
163  $\epsilon > 0$ , in the outer polytope of  $f_\theta$  w.r.t.  $x$ , there exists  $\delta$  such that if*

$$\min_{i=1, \dots, m} d(\delta x, x_i) \geq \sqrt{-\ell^2 \log \frac{(-\log(\frac{1}{\epsilon+0.5} - 1))^2 \text{var}(f_\theta(x))}{m |f_\mu(x)|^2}},$$

164 then it holds that

$$\text{conf}_t(\delta x) \leq \frac{1}{2} + \epsilon.$$

165 Furthermore, in the limit with  $\delta \rightarrow \infty$ , the uniform confidence  $\text{conf}(\delta x) = 0.5$  is attained.

166 *Proof.* By Lemma 2, at the outer polytope of  $f_\theta$  w.r.t.  $x$ , say  $P$ , we can write the restriction of  $f_\theta$  to  
167  $P$  as

$$f_\theta|_P(\delta x) = w(\theta)^\top (\delta x) = \delta w(\theta)^\top x = \delta f_\theta(x),$$

168 where  $\delta \geq \beta$ . Furthermore the gradient of  $f_\theta(\delta x)$  w.r.t.  $\theta$  (evaluated at  $\mu$ ) can be written as follows:

$$g(x) = \nabla f_\theta(\delta x) = J^\top(\delta x) = \delta J^\top x,$$

169 where the  $n \times p$  matrix  $J$  is the Jacobian  $\partial w / \partial \theta$  at  $\mu$ . The term inside the logistic in (7) is therefore

$$\begin{aligned}
|z(\delta x)| &= \frac{|\delta f_\mu(x)|}{(1 + t(\delta x)^{-1} \pi/8 \delta^2 (J^\top x)^\top \Sigma (J^\top x))^{\frac{1}{2}}} \\
&\leq \frac{\delta |f_\mu(x)|}{\delta (t(\delta x)^{-1} \pi/8 \text{var}(f_\theta(x)))^{\frac{1}{2}}} \\
&= \frac{\sqrt{t(\delta x)} |f_\mu(x)|}{\sqrt{\text{var}(f_\theta(x))}}.
\end{aligned} \tag{9}$$

170 Let  $d(\delta x, \mathcal{D}) := \min_{i=1, \dots, m} d(\delta x, x_i)$ . Notice that

$$t(\delta x) = \sum_{i=1}^m \exp\left(-\frac{1}{\ell^2} d(\delta x, x_i)^2\right) \leq m \exp\left(-\frac{1}{\ell^2} d(\delta x, \mathcal{D})^2\right).$$

171 Let  $\alpha := -\log\left(\frac{1}{\epsilon+0.5} - 1\right)$ . Substituting the previous inequality into (9) and using the hypothesis on  
172  $d(\delta x, \mathcal{D})$  yields

$$\begin{aligned}
|z(\delta x)| &\leq \frac{\sqrt{m \exp\left(-\frac{1}{\ell^2} d(\delta x, \mathcal{D})^2\right)} |f_\mu(x)|}{\sqrt{\text{var}(f_\theta(x))}} \\
&\leq \frac{\sqrt{m \exp\left(-\frac{1}{\ell^2} (-\ell^2) \log\left(\frac{\alpha^2 \text{var}(f_\theta(x))}{m |f_\mu(x)|^2}\right)\right)} |f_\mu(x)|}{\sqrt{\text{var}(f_\theta(x))}} \\
&= \frac{\sqrt{\frac{\alpha^2 \text{var}(f_\theta(x))}{|f_\mu(x)|^2}} |f_\mu(x)|}{\sqrt{\text{var}(f_\theta(x))}} \\
&= \frac{\alpha \sqrt{\text{var}(f_\theta(x))} |f_\mu(x)|}{|f_\mu(x)| \sqrt{\text{var}(f_\theta(x))}} \\
&= \alpha.
\end{aligned}$$

173 Since the logistic function is monotonically increasing, it holds that

$$\begin{aligned}
\text{conf}_t(\delta x) &= \sigma(|z(\delta x)|) \leq \sigma(\alpha) \\
&= \frac{1}{1 + \exp(-\alpha)} = \frac{1}{1 + \exp\left(\log\left(\frac{1}{\epsilon+0.5} - 1\right)\right)} \\
&= \frac{1}{\frac{1}{\epsilon+0.5}} = \frac{1}{2} + \epsilon,
\end{aligned}$$

174 which is the first desired result.

175 For the second result, Lemma 1 tells us that  $\lim_{\delta \rightarrow \infty} t(\delta x) = 0$  and thus  $\lim_{\delta \rightarrow \infty} t(\delta x)^{-1} = \infty$ .  
176 Therefore, by the definition of confidence (7) and the continuity of the logistic function at zero, it  
177 holds that

$$\begin{aligned}
\lim_{\delta \rightarrow \infty} \text{conf}_t(\delta x) &= \lim_{\delta \rightarrow \infty} \sigma\left(\frac{\delta |f_\mu(x)|}{(1 + t(\delta x)^{-1} \pi/8 \delta^2 \text{var}(f_\theta(x)))^{\frac{1}{2}}}\right) \\
&= \lim_{\delta \rightarrow \infty} \sigma\left(\frac{|f_\mu(x)|}{(\delta^{-2} + t(\delta x)^{-1} \pi/8 \text{var}(f_\theta(x)))^{\frac{1}{2}}}\right) \\
&= \sigma\left(\frac{|f_\mu(x)|}{(\lim_{\delta \rightarrow \infty} \delta^{-2} + \lim_{\delta \rightarrow \infty} t(\delta x)^{-1} \pi/8 \text{var}(f_\theta(x)))^{\frac{1}{2}}}\right) \\
&= \sigma\left(\frac{|f_\mu(x)|}{(\lim_{\delta \rightarrow \infty} t(\delta x)^{-1} \pi/8 \text{var}(f_\theta(x)))^{\frac{1}{2}}}\right) \\
&= \sigma(0) \\
&= \frac{1}{2},
\end{aligned}$$

178 which concludes the proof.  $\square$

For ReLU networks with bias, I can only say things in the limit. A fully rigorous proof for the asymptotic statement is surprisingly hard.

179

180 **Theorem 2.** Let  $f_\theta : \mathbb{R}^n \rightarrow \mathbb{R}$  be a binary ReLU classifier,  $\mathcal{N}(\theta|\mu, \Sigma)$  be the posterior, and let  
 181  $t : \mathbb{R}^n \rightarrow (0, \infty)$  be defined in (8). For any test point  $x \in \mathbb{R}^n$

$$\lim_{\delta \rightarrow \infty} \text{conf}_t(\delta x) = \frac{1}{2}.$$

182 *Proof.* By Lemma 2, at the outer polytope of  $f_\theta$  w.r.t.  $x$ , say  $P$ , we can write the restriction of  $f_\theta$  to  
 183  $P$  as

$$f_\theta|_P(\delta x) = w(\theta)^\top (\delta x) + b(\theta) = \delta(w(\theta)^\top x + \delta^{-1} b(\theta)),$$

184 where  $\delta \geq \beta$ . Furthermore the gradient of  $f_\theta(\delta x)$  w.r.t.  $\theta$  (evaluated at  $\mu$ ) can be written as follows:

$$g(x) = \nabla f_\theta|_P(\delta x) = J^\top (\delta x) + v = \delta(J^\top x + \delta^{-1} v),$$

185 where the  $n \times p$  matrix  $J$  is the Jacobian  $\partial w / \partial \theta$  at  $\mu$  and  $v \in \mathbb{R}^p$  is the gradient  $\nabla b$  of  $b$  w.r.t.  $\theta$  at  $\mu$ .  
 186 Now, Lemma 1 tells us that  $\lim_{\delta \rightarrow \infty} t(\delta x) = 0$  and thus  $\lim_{\delta \rightarrow \infty} t(\delta x)^{-1} = \infty$ . Therefore, by the  
 187 definition of confidence (7) and the continuity of the logistic function at zero, it holds that

$$\begin{aligned} \lim_{\delta \rightarrow \infty} \text{conf}_t(\delta x) &= \lim_{\delta \rightarrow \infty} \sigma \left( \frac{|f_\mu|_P(\delta x)|}{(1 + t(\delta x)^{-1} \pi/8 g(\delta x)^\top \Sigma g(\delta x))^{\frac{1}{2}}} \right) \\ &= \lim_{\delta \rightarrow \infty} \sigma \left( \frac{|\delta(w(\theta)^\top x + \delta^{-1} b(\theta))|}{(1 + t(\delta x)^{-1} \pi/8 (\delta(J^\top x + \delta^{-1} v))^\top \Sigma (\delta(J^\top x + \delta^{-1} v)))^{\frac{1}{2}}} \right) \\ &= \lim_{\delta \rightarrow \infty} \sigma \left( \frac{\delta |w(\theta)^\top x + \delta^{-1} b(\theta)|}{(1 + t(\delta x)^{-1} \pi/8 \delta^2 (J^\top x + \delta^{-1} v)^\top \Sigma (J^\top x + \delta^{-1} v))^{\frac{1}{2}}} \right) \\ &= \lim_{\delta \rightarrow \infty} \sigma \left( \frac{|w(\theta)^\top x + \delta^{-1} b(\theta)|}{(\delta^{-2} + t(\delta x)^{-1} \pi/8 (J^\top x + \delta^{-1} v)^\top \Sigma (J^\top x + \delta^{-1} v))^{\frac{1}{2}}} \right) \\ &= \sigma \left( \frac{\lim_{\delta \rightarrow \infty} |w(\theta)^\top x + \delta^{-1} b(\theta)|}{(\lim_{\delta \rightarrow \infty} \delta^{-2} + \lim_{\delta \rightarrow \infty} t(\delta x)^{-1} \pi/8 (J^\top x + \delta^{-1} v)^\top \Sigma (J^\top x + \delta^{-1} v))^{\frac{1}{2}}} \right) \\ &= \sigma \left( \frac{|w(\theta)^\top x|}{(\lim_{\delta \rightarrow \infty} t(\delta x)^{-1} \pi/8 (J^\top x + \delta^{-1} v)^\top \Sigma (J^\top x + \delta^{-1} v))^{\frac{1}{2}}} \right) \\ &= \sigma(0) \\ &= \frac{1}{2}, \end{aligned}$$

188 which concludes the proof.  $\square$

189 **Proposition 2.** Let  $f_\theta : \mathbb{R}^n \rightarrow \mathbb{R}$  be a binary ReLU classifier,  $\mathcal{N}(\theta|\mu, \Sigma)$  be the posterior, and let  
 190  $t : \mathbb{R}^n \rightarrow (0, \infty)$  be defined in (8). Let the functions  $z$  and  $z_t$  be defined as in (4) and (6), respectively.  
 191 Then for any  $x \in \mathbb{R}^n$ ,

$$\sigma(z_t(x)) = 0.5 \iff \sigma(z(x)) = 0.5 \quad (10)$$

192 *Proof.* It suffices to show that  $z_t(x) = 0 \iff z(x) = 0$  since  $\sigma^{-1}(0.5) = 0$ . The keys of this  
 193 proof are the observation that the codomain of  $t$  is the set of positive real numbers, along with the  
 194 fact that  $z_t$  is just  $z$  with the additional term  $1/t(x)$  in its denominator. In particular, notice that the  
 195 numerator of  $z$  and  $z_t$  is the same.

196 ( $\implies$ ) Suppose  $z_t(x) = 0$ . The denominator of  $z_t$  is positive since  $\pi/8 > 0$ , the matrix  $\Sigma$  is positive  
 197 definite, and  $t(x) \in (0, \infty)$  by definition. Thus its numerator must be zero, which implies that  $z(x)$   
 198 must also be zero.

199 ( $\impliedby$ ) Suppose  $z(x)$  is zero. Using the same argument, the denominator of  $z$  is positive and thus its  
 200 numerator must be zero, which implies that  $z_t(x)$  must also be zero.  $\square$

201 **5 Related work**

202 **6 Experiments**

203 **7 Conclusion**

204 **Broader Impact**

205 Authors are required to include a statement of the broader impact of their work, including its ethical  
206 aspects and future societal consequences. Authors should discuss both positive and negative outcomes,  
207 if any. For instance, authors should discuss a) who may benefit from this research, b) who may be  
208 put at disadvantage from this research, c) what are the consequences of failure of the system, and d)  
209 whether the task/method leverages biases in the data. If authors believe this is not applicable to them,  
210 authors can simply state this.

211 Use unnumbered first level headings for this section, which should go at the end of the paper. **Note**  
212 **that this section does not count towards the eight pages of content that are allowed.**

213 **References**

214 Anh Nguyen, Jason Yosinski, and Jeff Clune. Deep neural networks are easily fooled: High confidence  
215 predictions for unrecognizable images. In *CVPR*, 2015.

216 Matthias Hein, Maksym Andriushchenko, and Julian Bitterwolf. Why relu networks yield high-confidence  
217 predictions far away from the training data and how to mitigate the problem. In *CVPR*, 2019.

218 Agustinus Kristiadi, Matthias Hein, and Philipp Hennig. Being Bayesian, Even Just a Bit, Fixes Overconfidence  
219 in ReLU Networks. *arXiv preprint arXiv:2002.10118*, 2020.

220 Florian Wenzel, Kevin Roth, Bastiaan S Veeling, Jakub Świątkowski, Linh Tran, Stephan Mandt, Jasper Snoek,  
221 Tim Salimans, Rodolphe Jenatton, and Sebastian Nowozin. How Good is the Bayes Posterior in Deep Neural  
222 Networks Really? *arXiv preprint arXiv:2002.02405*, 2020.

223 Raman Arora, Amitabh Basu, Poorya Mianjy, and Anirbit Mukherjee. Understanding deep neural networks with  
224 rectified linear units. In *ICLR*, 2018.

225 David JC MacKay. Probable networks and plausible predictions—a review of practical bayesian methods for  
226 supervised neural networks. *Network: computation in neural systems*, 1995.

227 David JC MacKay. The evidence framework applied to classification networks. *Neural computation*, 1992.