

---

# Locking and Quacking: Stacking Bayesian models predictions by log-pooling and superposition

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Combining predictive distributions is a central problem in Bayesian inference and  
2 machine learning. Currently, predictives are almost exclusively combined using  
3 linear density-mixtures such as Bayesian model averaging, Bayesian stacking,  
4 and mixture of experts. Nonetheless, linear mixtures impose traits that might  
5 be undesirable for some applications, such as multi-modality. While there are  
6 alternative strategies (e.g., geometric bridge or superposition), optimizing their  
7 parameters usually implies computing intractable normalizing constant repeatedly.  
8 In this extended abstract, we present two novel Bayesian model combination tools.  
9 They are generalizations of *stacking*, but combine posterior densities by log-linear  
10 pooling (*locking*) and quantum superposition (*quacking*). To optimize model  
11 weights while avoiding the burden of normalizing constants, we maximize the  
12 Hyvärinen score of the combined posterior predictions. We demonstrate locking  
13 and quacking with an illustrative example.

## 14 1 Introduction

15 A general challenge in statistics is prediction in the presence of multiple candidate models or learning  
16 algorithms: we are interested in some outcome  $y$  on a measurable space  $\mathcal{Y} \subset \mathbb{R}^d$ ; We fit different  
17 models to the data, or the same model on different parts of the dataset, and obtain a set of predictive  
18 distributions,  $\{\pi_1(y), \dots, \pi_K(y)\}$ , where each  $\pi_k(y)$  is a (conditional<sup>1</sup>) probabilistic density such  
19 that  $\int_{\mathcal{Y}} \pi_k(y) dy = 1$ . When combining models, there are three subjective decisions to make: (1)  
20 individual models, (2) the “prior” assigned to each model, and (3) the form in which individual  
21 sampling models are combined in the predictive sampling distribution. Here, we focus on the latter.

The combination operation binds individual sampling distributions into a larger encompassing sampling model. A combination operator, parametrized by some parameter  $w$ , maps a sequence of probability densities into a single probability density:

$$h(\pi_1(\cdot), \dots, \pi_K(\cdot)|w) = \pi_*(\cdot), \text{ s.t. } \pi_*(\cdot) \geq 0, \mathbb{E} \pi_*(\cdot) = 1.$$

For example, a (linear) mixture can be represented by

$$h(\pi_1(\cdot), \dots, \pi_K(\cdot)|w) = \sum_k w_k \pi_k(\cdot), \quad \sum_k w_k = 1.$$

22 In Bayesian statistics, the *mixture* is the *de facto* combination operator to combine predictive distribu-  
23 tions, and used in Bayesian model averaging [11], stacking [15], hierarchical stacking [16], and as its  
24 name has implied, mixture-of-experts [6, 7, 17], and hypothesis testing [8]. Despite its mathematical  
25 convenience, *mixture* has a few limitations:

---

<sup>1</sup>the dependence on covariate  $x$  is suppressed for brevity

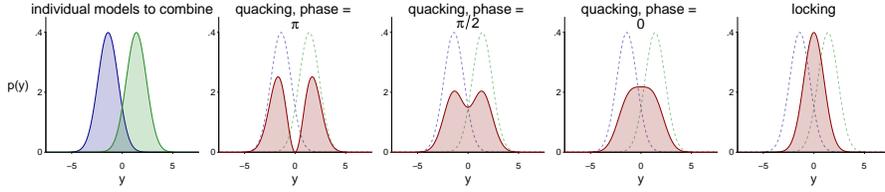


Figure 1: When combining two probabilistic predictions (panel 1), quacking combines them via superposition and locking combines them by geometric bridges.

26 Mixture is linear. If only relying on mixing to combine individual sampling models, the depth of  
 27 the combination network is restricted to one. It only examines likelihoods through their evaluations  
 28 at realized observations. Mixture of predictive densities typically results in a multimodal posterior  
 29 prediction, which comes with unnatural interpretation and poor interval coverage.

30 In this extended abstract, we investigate two operators to combine Bayesian predictives (Figure 1):

31 1. Geometric bridge (log-linear pooling),

$$h(\pi_1(\cdot), \dots, \pi_K(\cdot)|w) := \frac{\prod_k \pi_k^{w_k}(\cdot)}{\int \prod_k \pi_k^{w_k}(y) dy}, \quad w \in \mathcal{S}^k. \quad (1)$$

32 2. Superposition,

$$h(\pi_1(\cdot), \dots, \pi_K(\cdot)|w, \alpha) := \frac{|\sum_k \sqrt{w_k} \sqrt{\pi_k(\cdot)} e^{i\alpha_k}|^2}{\int_{\mathbb{R}} |\sum_k \sqrt{w_k} \sqrt{\pi_k(y')} e^{i\alpha_k}|^2 dy'}, \quad \alpha \in [0, 2\pi)^K, \quad w \in \mathcal{S}^k. \quad (2)$$

33 where  $\mathcal{S}^K$  is the  $k$ -dimensional simplex.

34 Compared with to mixture, these new operators have appealing features: When individual sampling  
 35 models are log-concave, so is their geometric bridge, hence preserving the unimodality. Moreover,  
 36 in superposition, when the phases  $\alpha$  are uniformly distributed, we get back a mixture of densities.  
 37 Even when there is only one single model, depending on the phase, the superposition and geometric  
 38 bridge can make the combined distributions spikier, or flatter— approximately a power transforma-  
 39 tion, thereby automatically calibrating the prediction confidence. Finally, unlike the mixture, the  
 40 superposition and geometric bridge can create a middle mode, leading to more flexible predictions.

41 The remaining question is then how to optimize the weights  $w$  such that the combined predictions can  
 42 better fit the data. This is challenging because of the intractable normalizing constant, and existing  
 43 log-linear pooling techniques rely on some non-testable normal approximation [1, 4, 10, 12]. In the  
 44 next section, we provide a practical solution that incorporates the Hyvärinen score [5] and Bayesian  
 45 posterior predictions.

## 46 2 Operator-oriented model averaging

### 47 2.1 Scoring rules and Hyvärinen score

48 In methods like stacking and mixture of experts, we need a scoring rule Gneiting and Raftery [2]  
 49 to evaluate the combined prediction, in which the logarithmic scoring rule is the de facto choice  
 50 for it is the only continuous proper local scoring rule. However, the log score does not apply to  
 51 log-linear pooling and superposition: unless in trivial cases, the combined predictive densities contain  
 52 an unknown normalization constant in the denominator.

To bypass the normalizing constant, we use the Hyvärinen score [5] to evaluate the unnormalized  
 combined predictive density. In general, given an unnormalized density  $q$ , how good it fits the  
 observed data  $y$  is quantified by

$$\mathcal{H}(y, p) = 2\Delta_y \log p(y) + \|\nabla_y \log p(y)\|^2.$$

53 The Hyvärinen score can be interpreted as the  $L_2$  norm of the difference between the score of the  
 54 prediction and the true data generating process.

## 55 2.2 Importance weighted estimate of the score function

56 Another distinctive feature of the full Bayesian prediction is that the posterior prediction itself is a  
 57 mixture of conditional sampling distributions. That is, in the  $k$ -th model, the posterior parameter  
 58 inference given observed data  $\mathbf{D}$  is  $p_k(\theta|\mathbf{D})$ , and the predictive density for future data  $\tilde{y}$  is  $\pi_k(\tilde{y}) =$   
 59  $\int_{\Theta} f(\tilde{y}|\theta)p_k(\theta|\mathbf{D}) d\theta$ . To compute the Hyvärinen score of this posterior prediction, we need the  
 60 pointwise score function

$$\frac{\partial}{\partial y} \log \pi_k(y) = \frac{\pi'_k(y)}{\pi_k(y)} = \frac{\int_{\Theta} \frac{\partial}{\partial y} f(y|\theta)p_k(\theta|\mathbf{D}) d\theta}{\int_{\Theta} f(y|\theta)p_k(\theta|\mathbf{D}) d\theta}. \quad (3)$$

61 We will typically use Markov chain Monte Carlo (MCMC) methods for individual model inference,  
 62 such that we have  $S$  simulation draws  $\{\theta_{sk}, 1 \leq s \leq S\}$  from the model  $k$  posterior  $p_k(\theta|\mathbf{D})$ . The  
 63 score function (3) is a ratio of integrals. We compute both the denominator and numerator by Monte  
 64 Carlo sum, and a plug-in estimate of the score function is

$$\frac{\partial}{\partial y} \log \pi_k(y) \approx g_k(y) := \frac{\sum_{s=1}^S \frac{\partial}{\partial y} f(y|\theta_{ks})}{\sum_{s=1}^S f(y|\theta_{ks})}. \quad (4)$$

65 There is no worry that the denominator and numerator are estimated using the same draws: We can  
 66 view (4) as self-normalized importance sampling with a proposal density  $p_k(\theta|\mathbf{D})$ , target density  
 67  $f(y|\theta)p_k(\theta|\mathbf{D})$ , and a function  $h(\theta) = \frac{\partial}{\partial y} f(y|\theta)/f(y|\theta)$ <sup>2</sup>. The usual convergence theory of self-  
 68 normalized importance sampling [e.g., 9] guarantees the consistency and asymptotically normality of  
 69 our score function estimate (4).

70 Similarly, for the second derivative, the Monte Carlo estimate given each model is also consistent,

$$\frac{\partial^2}{\partial y^2} \log \pi_k(y) = \frac{\pi_k(y)\pi''_k(y) - [\pi'_k(y)]^2}{[\pi_k(y)]^2} \approx h_{ik} := \frac{\sum_{s=1}^S f''(y|\theta_{ks})}{\sum_{s=1}^S f(y|\theta_{ks})} - [g_{ik}]^2.$$

## 71 2.3 Proposed method: optimizing the Hyvärinen score of the combined posterior densities

72 Our general model combination method contains five steps as follows:

73 **Step 1:** fit each model to the data and obtain  $K$  predictive densities. In practice the pos-  
 74 teriors  $p_k(\theta|\mathbf{D})$  are represented by Monte Carlo draws,  $\theta_{k1}, \dots, \theta_{kS}$ , leading to the estimate  
 75  $\pi_k(\cdot) := \frac{1}{S} \sum_{s=1}^S p_k(\cdot|\theta_{ks})$ . **Step 2:** express the unnormalized predictive density via the combination  
 76 operator. For example, in locking we have  $q(\cdot|w) := \prod_k \pi_k^{w_k}(\cdot)$ . **Step 3:** evaluate  $\nabla_y \log q(\cdot|w)$  and  
 77  $\Delta_y \log q(\cdot|w)$  at every observed  $y_i$  points. They come in *closed form functions* of  $\nabla_y \pi_k(y_i|\theta_{ks})$  and  
 78  $\Delta_y \pi_k(y_i|\theta_{ks})$ . In locking:

$$q'_i(w) := \nabla_y \log q(y_i|w) = \sum_{k=1}^K w_k \nabla_y \log (\pi_k(y_i)) \approx \sum_{k=1}^K w_k \frac{\sum_{s=1}^S \nabla_y p_k(y_i|\theta_{ks})}{\sum_{s=1}^S p_k(y_i|\theta_{ks})}, \quad (5)$$

$$q''_i(w) := \Delta_y \log q(y_i|w) \approx \sum_{k=1}^K \frac{w_k}{S} \sum_{s=1}^S (\Delta_y \log p_k(y_i|\theta_{ks})). \quad (6)$$

79 The quacking derivatives also come in closed form expression (functions of weight  $w$  and phase  $\alpha$ ),  
 80 but we omit them here.

81 **Step 4:** optimize model weight vector  $w$  by the constrained optimization

$$\hat{w}_{\text{opt}} = \min_w \left( \sum_{i=1}^n (2q''_i(w) + |q'_i(w)|^2) - \log \text{prior}(w) \right), \quad \text{s.t.} \quad \sum_{i=1}^K w_k = 1, \quad w_i \geq 0. \quad (7)$$

82 We use an non-informative prior Dirichlet (1.01) for weight regularization.

<sup>2</sup>Rewrite the score function into (3) into  $\int_{\Theta} h(\theta) \frac{f(y|\theta)p_k(\theta|\mathbf{D})}{\int_{\Theta} f(y|\theta)p_k(\theta|\mathbf{D}) d\theta} d\theta$  admits the self-normalized importance sampling estimate.

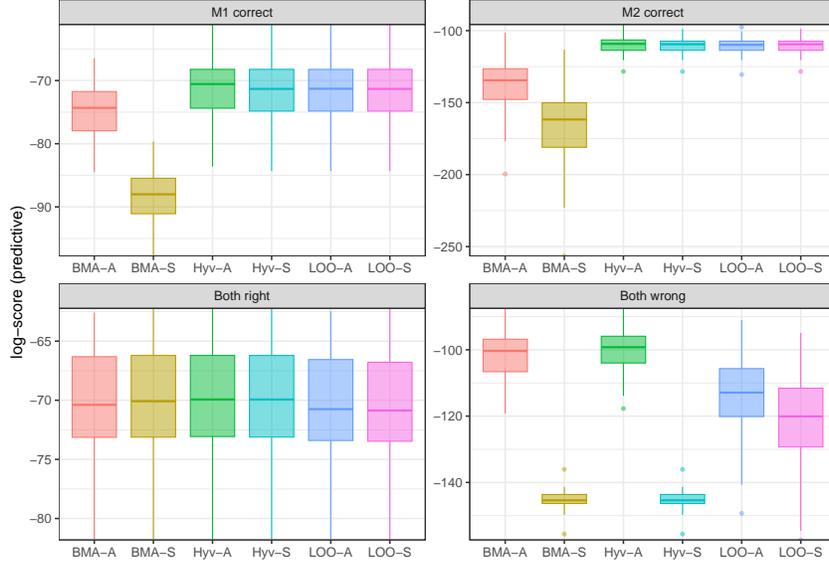


Figure 2: Log predictive scores. For each method, we show the overall predictive log score,  $\sum_{j=1}^{N_{\text{pred}}} \log \pi_*(y_j)$ . Our method (Hyvärinen model averaging, green) achieves higher log-scores in all four scenarios.

83 *Complexity.* The key blessing of applying scoring matching to Bayesian predictions is that the  
 84 Monte Carlo integral is linear, and is exchangeable with gradient operators. Hence, all we need is to  
 85 compute and store the gradient and hessian of the log likelihood (with respect to data) at the sampled  
 86 parameters once, that is  $\nabla_y \pi_k(y_i | \theta_{k,s})$  and  $\Delta_y \pi_k(y_i | \theta_{k,s})$ . In particular, the score functions have  
 87 already been computed in gradient-based MCMC sampler, such as in dynamic Hamiltonian Monte  
 88 Carlo [e.g., 3], hence nearly free. The summation in the objective function 7 contains  $nKS$  gradient  
 89 evaluations in total, and can be computed in parallel.

### 90 3 Example

91 To manifest the flexibility of our new approach, in this section we run experiments and compare the  
 92 locking method to other state of the art model averaging and selection tools. We adapt the experiment  
 93 setting from Shao et al. [13]. Consider two normal belief models.

$$\begin{aligned} \mathcal{M}_1 : Y_i &\sim \text{Normal}(\theta_1, 1), \theta_1 \sim \text{Normal}(0, v_0), \\ \mathcal{M}_2 : Y_i &\sim \text{Normal}(0, \theta_2), \theta_2 \sim \text{Inverse-}\chi^2(\nu_0, \tau_0). \end{aligned}$$

94 Following Shao et al. [13], we picked  $v_0 = 10$ ,  $\nu_0 = 0.1$  and  $\tau_0 = 1$ . We simulate  $N_{\text{train}}$  data points  
 95 from a true data generating process: a normal distribution with mean  $\mu^*$  and variance  $v^*$ . We also  
 96 generate a  $N_{\text{test}}$  independent test samples. We consider four scenarios: (1)  $\mu^* = 1$  and  $v^* = 1$   
 97 meaning that  $\mathcal{M}_1$  is correctly specified but  $\mathcal{M}_2$  is not; (2)  $\mu^* = 0$  and  $v^* = 5$  meaning that  $\mathcal{M}_2$   
 98 is correctly specified but  $\mathcal{M}_1$  is not; (3)  $\mu^* = 4$  and  $v^* = 3$ , a situation in which neither model  
 99 is correctly specified and (4)  $\mu^* = 0$  and  $v^* = 1$ , in which both are correctly specified. We ran  
 100  $M = 100$  replications of each scenario, with  $N_{\text{train}} = 200$  and  $N_{\text{test}} = 50$ .

101 We compare six methods in total, (1) model selection using marginal likelihood, (2) Bayesian  
 102 model averaging (3) model selection using leave-one-out log predictive densities [LOO-elpd, 14],  
 103 (4) Bayesian stacking [15], (5) model selection using Hyvärinen score [13], and (6) *locking* (ours).  
 104 We evaluate predictive performance of the learned combined model. To make the comparison fair,  
 105 we pick a metric that we do not directly optimize over: the log predictive density on test data. As  
 106 shown in Figure 2, our new *locking* method outperforms all its alternatives in terms of log-scores.  
 107 We show that the Hyvärinen stacking (i.e. first optimising the weights using the Hyvärinen score  
 108 on held-out data then forming a log-pool with the optimised weights) leads to higher log-predictive  
 109 densities overall.

## References

- 110 [1] Luiz M Carvalho, Daniel AM Villela, Flavio C Coelho, and Leonardo S Bastos. Bayesian  
111 inference for the weights in logarithmic pooling. *Bayesian Analysis*, 1(1):1–29, 2022.  
112
- 113 [2] Tilmann Gneiting and Adrian E Raftery. Strictly proper scoring rules, prediction, and estimation.  
114 *Journal of the American Statistical Association*, 102(477):359–378, 2007.
- 115 [3] Matthew D Hoffman, Andrew Gelman, et al. The No-U-Turn sampler: adaptively setting path  
116 lengths in hamiltonian Monte Carlo. *Journal of Machine Learning Research*, 15(1):1593–1623,  
117 2014.
- 118 [4] Zaijing Huang and Andrew Gelman. Sampling for Bayesian computation with large datasets.  
119 *Technical Report, Columbia University*, 2005.
- 120 [5] Aapo Hyvärinen. Estimation of non-normalized statistical models by score matching. *Journal*  
121 *of Machine Learning Research*, 6(4), 2005.
- 122 [6] Robert A Jacobs, Michael I Jordan, Steven J Nowlan, and Geoffrey E Hinton. Adaptive mixtures  
123 of local experts. *Neural computation*, 3(1):79–87, 1991.
- 124 [7] Michael I Jordan and Robert A Jacobs. Hierarchical mixtures of experts and the EM algorithm.  
125 *Neural computation*, 6(2):181–214, 1994.
- 126 [8] K Kamary, K Mengersen, CP Robert, and J Rousseau. Bayesian hypothesis testing as a mixture  
127 estimation model. *arXiv preprint arXiv:1412.2044*, 2018.
- 128 [9] Art B. Owen. *Monte Carlo theory, methods and examples*. 2013.
- 129 [10] David Poole and Adrian E Raftery. Inference for deterministic simulation models: the Bayesian  
130 melding approach. *Journal of the American Statistical Association*, 95(452):1244–1255, 2000.
- 131 [11] Adrian E Raftery, David Madigan, and Jennifer A Hoeting. Bayesian model averaging for linear  
132 regression models. *Journal of the American Statistical Association*, 92(437):179–191, 1997.
- 133 [12] MJ Rufo, J Martin, and CJ Pérez. Log-linear pool to combine prior distributions: A suggestion  
134 for a calibration-based approach. *Bayesian Analysis*, 7(2):411–438, 2012.
- 135 [13] Stephane Shao, Pierre E Jacob, Jie Ding, and Vahid Tarokh. Bayesian model comparison  
136 with the hyvärinen score: Computation and consistency. *Journal of the American Statistical*  
137 *Association*, 2019.
- 138 [14] Aki Vehtari, Andrew Gelman, and Jonah Gabry. Practical Bayesian model evaluation using  
139 leave-one-out cross-validation and waic. *Statistics and computing*, 27(5):1413–1432, 2017.
- 140 [15] Yuling Yao, Aki Vehtari, Daniel Simpson, and Andrew Gelman. Using stacking to average  
141 Bayesian predictive distributions. *Bayesian Analysis*, 13(3):917–1007, 2018.
- 142 [16] Yuling Yao, Gregor Pirš, Aki Vehtari, and Andrew Gelman. Bayesian hierarchical stacking:  
143 Some models are (somewhere) useful. *Bayesian Analysis*, 1(1):1–29, 2021.
- 144 [17] Seniha Esen Yuksel, Joseph N Wilson, and Paul D Gader. Twenty years of mixture of experts.  
145 *IEEE Transactions on Neural Networks and Learning Systems*, 23(8):1177–1193, 2012.