

1

2 **Reverse engineering neural networks to characterise their cost functions**

3

4 Takuya Isomura¹, Karl Friston²

5 1 Laboratory for Neural Computation and Adaptation, RIKEN Center for Brain Science, Wako,
6 Saitama 351-0198, Japan

7 2 Wellcome Centre for Human Neuroimaging, Institute of Neurology, University College
8 London, 12 Queen Square, London, WC1N 3AR, UK

9 Corresponding author email: takuya.isomura@riken.jp

10

11

12

13 **Abstract**

14 This work considers a class of biologically plausible cost functions for neural networks, where
15 the same cost function is minimised by both neural activity and plasticity. We show that such
16 cost functions can be cast as a variational bound on model evidence under an implicit
17 generative model. Using generative models based on Markov decision processes (MDP), we
18 show, analytically, that neural activity and plasticity perform Bayesian inference and learning,
19 respectively, by maximising model evidence. Using mathematical and numerical analyses, we
20 then confirm that biologically plausible cost functions—used in neural
21 networks—correspond to variational free energy under some prior beliefs about the
22 prevalence of latent states that generate inputs. These prior beliefs are determined by
23 particular constants (i.e., thresholds) that define the cost function. This means that the Bayes
24 optimal encoding of latent or hidden states is achieved when, and only when, the network’s
25 implicit priors match the process that generates the inputs. Our results suggest that when a
26 neural network minimises its cost function, it is implicitly minimising variational free energy
27 under optimal or sub-optimal prior beliefs. This insight is potentially important because it
28 suggests that any free parameter of a neural network’s cost function can itself be
29 optimised—by minimisation with respect to variational free energy.

30

31 **Keywords:** free-energy principle, variational Bayesian inference, learning algorithm, synaptic
32 plasticity, Markov decision process, blind source separation

33

34 **1. Introduction**

35 Cost functions are ubiquitous in scientific fields that entail optimisation—including physics,
36 chemistry, biology, engineering, and machine learning. Furthermore, any optimisation
37 problem that can be specified using a cost function can be formulated as a gradient descent.
38 In the neurosciences, this enables one to treat neuronal dynamics and plasticity as an
39 optimisation process (Marr, 1969; Albus, 1971; Schultz et al., 1997; Sutton & Barto, 1998;
40 Linsker, 1988; Brown et al., 2001). These examples highlight the importance of specifying a
41 problem in terms of cost functions, from which neural and synaptic dynamics can be derived.
42 In other words, cost functions provide a formal (i.e., normative) expression of the purpose of
43 a neural network and prescribe the dynamics of that neural network. Crucially, once the cost
44 function has been established and an initial condition has been selected, it is no longer
45 necessary to solve the dynamics. Instead, one can characterise the neural network’s
46 behaviour in terms of fixed points, basin of attraction, and structural stability—based on and
47 only on the cost function. In short, it is important to identify the cost function to understand
48 the dynamics, plasticity, and function of a neural network.

49 A ubiquitous cost function in neurobiology, theoretical biology, and machine learning is

50 model evidence, or equivalently, marginal likelihood or surprise; namely, the probability of
51 some inputs or data under a model of how those inputs were generated by unknown or
52 hidden causes. Generally, the evaluation of surprise is intractable. However, this evaluation
53 can be converted into an optimisation problem by inducing a variational bound on surprise.
54 In machine learning, this is known as an evidence lower bound (ELBO), while the same
55 quantity is known as variational free energy in statistical physics and theoretical
56 neurobiology.

57 Variational free energy minimisation is a candidate principle that governs neuronal activity
58 and synaptic plasticity (Friston et al., 2006; Friston, 2010). Here, surprise reflects the
59 improbability of sensory inputs given a model of how those inputs were caused. In turn,
60 minimising variational free energy, as a proxy for surprise, corresponds to inferring the
61 (unobservable) causes of (observable) consequences. To the extent that biological systems
62 minimise variational free energy, it is possible to say that they infer and learn the hidden
63 states and parameters that generate their sensory inputs (Helmholtz, 1925; Knill & Pouget,
64 2004; DiCarlo et al., 2012) and consequently predict those inputs (Rao & Ballard, 1999;
65 Friston, 2005). This is generally referred to as perceptual inference based upon an internal
66 generative model about the external world (Dayan et al., 1995; George & Hawkins, 2009;
67 Bastos et al., 2012).

68 Variational free energy minimisation provides a unified mathematical formulation of these
69 inference and learning processes in terms of self-organising neural networks that function as
70 Bayes optimal encoders. Moreover, organisms can use the same cost function to control their
71 surrounding environment by sampling predicted (i.e., preferred) inputs. This is known as
72 active inference (Friston et al., 2011). The ensuing free-energy principle suggests that active
73 inference and learning are mediated by changes in neural activity, synaptic strengths, and the
74 behaviour of an organism to minimise variational free energy, as a proxy for surprise.
75 Crucially, variational free energy and model evidence rest upon a generative model of
76 continuous or discrete hidden states. A number of recent studies have used Markov decision
77 process (MDP) generative models to elaborate schemes that minimise variational free energy
78 (Friston, FitzGerald et al., 2016; Friston, FitzGerald et al., 2017; Friston, Parr et al., 2017). This
79 minimisation reproduces various interesting dynamics and behaviours of real neuronal
80 networks and biological organisms. However, it remains to be established whether
81 variational free energy minimisation is an apt explanation for any given neural network, as
82 opposed to the optimisation of alternative cost functions.

83 In principle, any neural network that produces an output or a decision can be cast as
84 performing some form of inference, in terms of Bayesian decision theory. On this reading,
85 the complete class theorem suggests that any neural network can be regarded as performing
86 Bayesian inference under some prior beliefs; therefore, it can be regarded as minimising
87 variational free energy. The complete class theorem (Wald, 1947; Brown, 1981) states that
88 for any pair of decisions and cost functions, there are some prior beliefs (implicit in the
89 generative model) that render the decisions Bayes optimal. This suggests that it should be

90 theoretically possible to identify an implicit generative model within any neural network
91 architecture, which renders its cost function a variational free energy or ELBO. In what
92 follows, we show that such identification is possible for a fairly canonical form of a neural
93 network and a generic form of a generative model.

94 In brief, we adopt a reverse engineering approach to identify a plausible cost function for
95 neural networks—and show that the resulting cost function is formally equivalent to
96 variational free energy. Here, we define a cost function as a function of sensory input, neural
97 activity, and synaptic strengths and suppose that neural activity and synaptic plasticity
98 follows a gradient descent on the cost function. For simplicity, we consider single-layer
99 feed-forward neural networks comprising firing rate neuron models and focus on blind
100 source separation (BSS); namely, the problem of separating sensory inputs into multiple
101 hidden sources or causes (Belouchrani et al., 1997; Cichocki et al., 2009; Comon & Jutten,
102 2010), which provides the minimum setup for modelling causal inference. Previously, we
103 observed BSS performed by *in vitro* neural networks (Isomura et al., 2015) and reproduced
104 this self-supervised process using an MDP and variational free energy minimisation (Isomura
105 & Friston, 2018). These works suggest that variational free energy minimisation offers a
106 plausible account of the empirical behaviour of *in vitro* networks.

107 In this work, we ask whether variational free energy minimisation can account for the
108 normative behaviour of a canonical neural network that minimises its cost function, by
109 considering all possible cost functions, within a generic class. Using mathematical analysis,
110 we identify a class of cost functions—from which update rules for both neural activity and
111 synaptic plasticity can be derived—when a single-layer feed-forward neural network
112 comprises firing rate neurons whose firing intensity is determined by the sigmoid activation
113 function. The gradient descent on the ensuing cost function naturally leads to Hebbian
114 plasticity with an activity-dependent homeostatic term. We show that these cost functions
115 are formally homologous to variational free energy under an MDP. Finally, we discuss the
116 implications of this result for explaining the empirical behaviour of neuronal networks, in
117 terms of free energy minimisation under particular prior beliefs.

118

119 **2. Methods**

120 In this section, we first derive the form of a variational free energy cost function under a
121 specific generative model; namely a Markov decision process¹. We will go through the
122 derivations carefully, with a focus on the form of the ensuing Bayesian belief updating. The

¹ Strictly speaking, the generative model used in this paper is a hidden Markov model (HMM) because we do not consider probabilistic transitions between hidden states that depend upon control variables. However, for consistency with the literature on variational treatments of discrete state space models, we retain the MDP formalism; noting that we are using a special case (with unstructured state transitions).

123 form of this update will re-emerge later, when reverse engineering the cost functions implicit
 124 in neural networks. This section starts with a description of Markov decision processes—as a
 125 general kind of generative model—and then considers the minimisation of variational free
 126 energy under these models.

127 **2.1 Generative models.** Under an MDP model (Fig. 1A), a minimal BSS setup (in a
 128 discrete-space) reduces to the likelihood mapping from N_s hidden sources or states $s_t \equiv$
 129 $(s_t^{(1)}, \dots, s_t^{(N_s)})^T$ to N_o observations $o_t \equiv (o_t^{(1)}, \dots, o_t^{(N_o)})^T$. Each source and observation
 130 takes a value of one (ON state) or zero (OFF state) at each time step; i.e., $s_t^{(j)}, o_t^{(i)} \in \{1,0\}$.
 131 Throughout this paper, j denotes the j -th hidden state, while i denotes the i -th observation.
 132 The probability of $s_t^{(j)}$ follows a categorical distribution $P(s_t^{(j)}) = \text{Cat}(D^{(j)})$, where
 133 $D^{(j)} \equiv (D_1^{(j)}, D_0^{(j)}) \in \mathbb{R}^2$ with $D_1^{(j)} + D_0^{(j)} = 1$.

134 The probability of an outcome is determined by the likelihood mapping from all hidden
 135 states to each kind of observation in terms of a categorical distribution, $P(o_t^{(i)} | s_t, A^{(i)}) =$
 136 $\text{Cat}(A^{(i)})$. Here, each element of the tensor $A^{(i)} \in \mathbb{R}^{2 \times 2^{N_s}}$ parameterises the probability
 137 that $P(o_t^{(i)} = k | s_t = \vec{l})$, where $k \in \{1,0\}$ are possible observations and $\vec{l} \in \{1,0\}^{N_s}$
 138 encodes a particular combination of hidden states. The prior distribution of each column of
 139 $A^{(i)}$, denoted by $A_{\vec{l}}^{(i)}$, has a Dirichlet distribution $P(A_{\vec{l}}^{(i)}) = \text{Dir}(a_{\vec{l}}^{(i)})$ with concentration
 140 parameter $a_{\vec{l}}^{(i)} \in \mathbb{R}^2$. We use Dirichlet distributions, as they are tractable and widely used
 141 for random variables that take a continuous value between 0 and 1. Furthermore, learning
 142 the likelihood mapping leads to biologically plausible update rules, which have the form of
 143 associative or Hebbian plasticity: please see below and (Friston et al., 2016) for details.

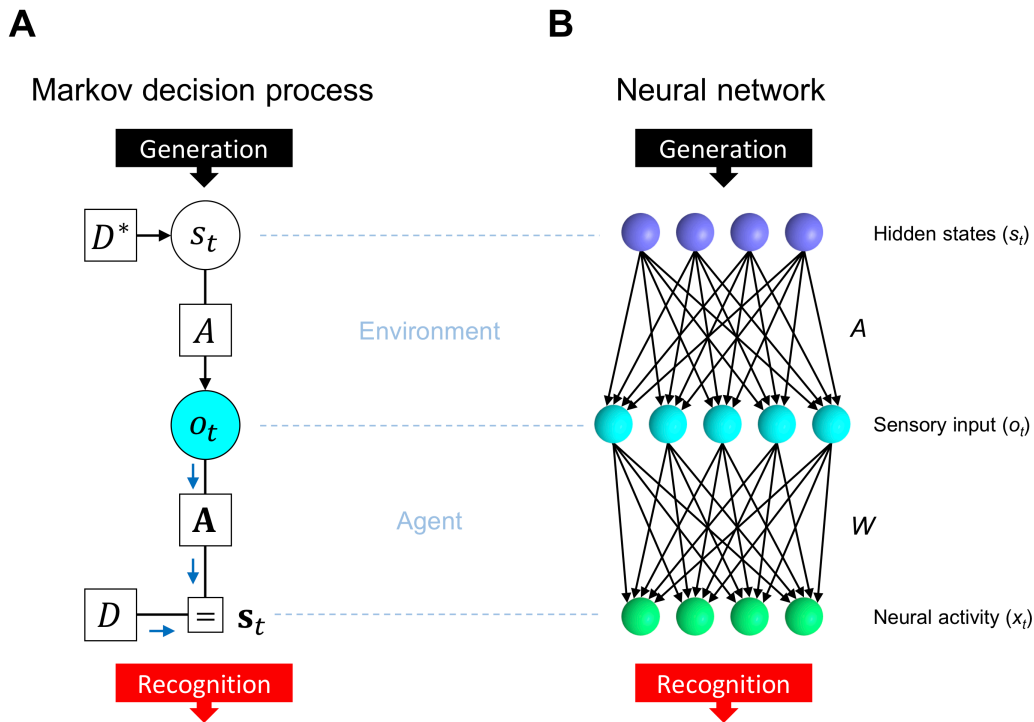
144 We use $\tilde{o} \equiv (o_1, \dots, o_t)$ and $\tilde{s} \equiv (s_1, \dots, s_t)$ to denote sequences of observations and
 145 hidden states, respectively. With this notation in place, the generative model (i.e., the joint
 146 distribution over outcomes, hidden states, and the parameters of their likelihood mapping)
 147 can be expressed as

148
$$P(\tilde{o}, \tilde{s}, A) = P(A) \prod_{\tau=1}^t P(o_\tau | s_\tau, A) P(s_\tau)$$

$$149 \quad = \prod_{i=1}^{N_o} P(A^{(i)}) \cdot \prod_{\tau=1}^t \left\{ \prod_{i=1}^{N_o} P(o_{\tau}^{(i)} | s_{\tau}, A^{(i)}) \prod_{j=1}^{N_s} P(s_{\tau}^{(j)}) \right\}. \quad (1)$$

150 Throughout this paper, t denotes the current time, while τ denotes an arbitrary time from
 151 the past to the present, $1 \leq \tau \leq t$.

152



153

154 **Figure 1.** Comparison between an MDP scheme and a neural network. (A) MDP scheme
 155 expressed as a Forney factor graph (Forney, 2001; Dauwels, 2007) based upon the
 156 formulation in (Friston, Parr et al., 2017). In this BSS setup, the prior D determines hidden
 157 states s_t , while s_t determines observation o_t through the likelihood mapping A . Inference
 158 corresponds to the inversion of this generative process. Here, D^* indicates the true prior
 159 while D indicates the prior under which the network operates. If $D = D^*$, the inference is
 160 optimal; otherwise, it is biased. (B) Neural network comprising a single layer feed-forward
 161 network with a sigmoid activation function. The network receives sensory inputs $o_t =$
 162 $(o_t^{(1)}, \dots, o_t^{(N_o)})^T$ that are generated from hidden states $s_t = (s_t^{(1)}, \dots, s_t^{(N_s)})^T$ and outputs
 163 neural activities $x_t = (x_{t1}, \dots, x_{tN_x})^T$. Here, x_{tj} should encode the posterior expectation
 164 about a binary state $s_t^{(j)}$.

165

166 **2.2 Minimisation of variational free energy.** In this MDP scheme, the aim is to minimise
 167 surprise by minimising variational free energy as a proxy; i.e., performing approximate or
 168 variational Bayesian inference. From the generative model, we can motivate a mean-field
 169 approximation to the posterior (i.e., recognition) density as follows:

$$170 \quad Q(\tilde{s}, A) = Q(A)Q(\tilde{s}) = \prod_{i=1}^{N_o} Q(A^{(i)}) \cdot \prod_{\tau=1}^t \prod_{j=1}^{N_s} Q(s_{\tau}^{(j)}), \quad (2)$$

171 where $A^{(i)}$ is the likelihood mapping (i.e., tensor), and the marginal posterior distributions
 172 of $s_{\tau}^{(j)}$ and $A^{(i)}$ have a categorical $Q(s_{\tau}^{(j)}) = \text{Cat}(\mathbf{s}_{\tau}^{(j)})$ and Dirichlet distribution
 173 $Q(A^{(i)}) = \text{Dir}(\mathbf{a}^{(i)})$, respectively. For simplicity, we assume that $A^{(i)}$ factorises into the
 174 product of the likelihood mappings from the j -th hidden state to the i -th observation: $A_k^{(i)} \approx$
 175 $A_k^{(i,1)} \otimes \dots \otimes A_k^{(i,N_s)}$ (where \otimes denotes the outer product and $A^{(i,j)} \in \mathbb{R}^{2 \times 2}$). This (mean
 176 field) approximation simplifies the computation of the state posteriors.

177 In what follows, a bold case variable indicates the posterior expectation of the
 178 corresponding variable in italics. For example, $s_{\tau}^{(j)}$ takes the value 0 or 1, while the
 179 posterior expectation $\mathbf{s}_{\tau}^{(j)} \in \mathbb{R}^2$ is the expected value of $s_{\tau}^{(j)}$ that lies between 0 and 1.
 180 Moreover, $\mathbf{a}^{(i,j)} \in \mathbb{R}^{2 \times 2}$ denotes positive concentration parameters. Below, we use the
 181 posterior expectation of $\ln A^{(i,j)}$ to encode posterior beliefs about the likelihood, which are
 182 given by

$$183 \quad \ln \mathbf{A}^{(i,j)} \equiv \mathbb{E}_{Q(A^{(i,j)})}[\ln A^{(i,j)}] = \psi(\mathbf{a}_l^{(i,j)}) - \psi(\mathbf{a}_{1l}^{(i,j)} + \mathbf{a}_{0l}^{(i,j)})$$

$$184 \quad = \ln \mathbf{a}_l^{(i,j)} - \ln(\mathbf{a}_{1l}^{(i,j)} + \mathbf{a}_{0l}^{(i,j)}) + \mathcal{O}\left(\left(\mathbf{a}_l^{(i,j)}\right)^{-1}\right), \quad (3)$$

185 where $l \in \{1,0\}$. Here, $\psi(\cdot) \equiv \Gamma'(\cdot)/\Gamma(\cdot)$ denotes the digamma function, which arises
 186 naturally from the definition of the Dirichlet distribution. Please see (Friston et al., 2016) for
 187 details. $\mathbb{E}_{Q(A^{(i,j)})}[\cdot]$ denotes the expectation over the posterior of $A^{(i,j)}$.

188 The ensuing variational free energy of this generative model is then given by

$$189 \quad F(\tilde{o}, Q(\tilde{s}), Q(A))$$

$$190 \quad \equiv \sum_{\tau=1}^t \left\{ \mathbb{E}_{Q(s_{\tau})Q(A)}[-\ln P(o_{\tau}|s_{\tau}, A)] + \mathcal{D}_{\text{KL}}[Q(s_{\tau})||P(s_{\tau})] \right\} + \mathcal{D}_{\text{KL}}[Q(A)||P(A)]$$

$$\begin{aligned}
 191 \quad &= \underbrace{\sum_{j=1}^{N_s} \sum_{\tau=1}^t \mathbf{s}_{\tau}^{(j)} \cdot \left(- \sum_{i=1}^{N_o} \ln \mathbf{A}^{(i,j)} \cdot o_{\tau}^{(i)} + \ln \mathbf{s}_{\tau}^{(j)} - \ln D^{(j)} \right)}_{\text{accuracy+state complexity}} \\
 192 \quad &+ \underbrace{\sum_{i=1}^{N_o} \sum_{j=1}^{N_s} \{ (\mathbf{a}^{(i,j)} - a^{(i,j)}) \cdot \ln \mathbf{A}^{(i,j)} - \ln \mathcal{B}(\mathbf{a}^{(i,j)}) \}}_{\text{parameter complexity}}, \quad (4)
 \end{aligned}$$

193 where $\ln \mathbf{A}^{(i,j)} \cdot o_{\tau}^{(i)}$ denotes the inner product of $\ln \mathbf{A}^{(i,j)}$ and a one-hot encoded vector
 194 of $o_{\tau}^{(i)}$, $\mathcal{D}_{\text{KL}}[\cdot || \cdot]$ is the complexity as scored by the Kullback–Leibler divergence (Kullback
 195 & Leibler, 1951), and $\mathcal{B}(\mathbf{a}^{(i,j)}) \equiv \mathcal{B}(\mathbf{a}_1^{(i,j)}) \mathcal{B}(\mathbf{a}_0^{(i,j)})$ with $\mathcal{B}(\mathbf{a}_l^{(i,j)}) \equiv \Gamma(\mathbf{a}_{1l}^{(i,j)}) \Gamma(\mathbf{a}_{0l}^{(i,j)}) /$
 196 $\Gamma(\mathbf{a}_{1l}^{(i,j)} + \mathbf{a}_{0l}^{(i,j)})$ is the beta function. The first term in the final equality comprises the
 197 accuracy ($-\mathbf{s}_{\tau}^{(j)} \cdot \sum_{i=1}^{N_o} \ln \mathbf{A}^{(i,j)} \cdot o_{\tau}^{(i)}$) and (state) complexity ($\mathbf{s}_{\tau}^{(j)} \cdot (\ln \mathbf{s}_{\tau}^{(j)} - \ln D^{(j)})$). The
 198 accuracy term is simply the expected log likelihood of an observation, while complexity
 199 scores the divergence between prior and posterior beliefs. In other words, complexity
 200 reflects the degree of belief updating or degrees of freedom required to provide an accurate
 201 account of observations. Both belief updates to states and parameters incur a complexity
 202 cost: the state complexity increases with time t , while parameter complexity increases in the
 203 order of $\ln t$ —and is thus negligible when t is large (see Supplementary Methods S1 for
 204 details). This means that we can ignore parameter complexity, when the scheme has
 205 experienced a sufficient number of outcomes. We will drop the parameter complexity in
 206 subsequent sections. In the remainder of this section, we show how the minimisation of
 207 variational free energy transforms (i.e., updates) priors into posteriors, when the parameter
 208 complexity is evaluated explicitly.

209 Inference optimises posterior expectations about the hidden states by minimising
 210 variational free energy. The optimal posterior expectations are obtained by solving the
 211 variation of F to give

$$212 \quad \mathbf{s}_t^{(j)} = \sigma \left(\sum_{i=1}^{N_o} \ln \mathbf{A}^{(i,j)} \cdot o_t^{(i)} + \ln D^{(j)} \right) = \sigma(\ln \mathbf{A}^{(\cdot,j)} \cdot o_t + \ln D^{(j)}), \quad (5)$$

213 where $\sigma(\cdot)$ is the softmax function. As $s_t^{(j)}$ is a binary value in this work, the posterior
 214 expectation of $s_t^{(j)}$ taking a value of one (ON state) can be expressed as

$$\begin{aligned}
 215 \quad \mathbf{s}_{t1}^{(j)} &= \frac{\exp(\ln \mathbf{A}_1^{(i,j)} \cdot o_t + \ln D_1^{(j)})}{\exp(\ln \mathbf{A}_1^{(i,j)} \cdot o_t + \ln D_1^{(j)}) + \exp(\ln \mathbf{A}_0^{(i,j)} \cdot o_t + \ln D_0^{(j)})} \\
 216 \quad &= \text{sig}(\ln \mathbf{A}_1^{(i,j)} \cdot o_t - \ln \mathbf{A}_0^{(i,j)} \cdot o_t + \ln D_1^{(j)} - \ln D_0^{(j)}) \quad (6)
 \end{aligned}$$

217 using the sigmoid function $\text{sig}(z) \equiv 1/(1 + \exp(-z))$. Thus, the posterior expectation of
 218 $s_t^{(j)}$ taking a value 0 (OFF state) is $\mathbf{s}_{t0}^{(j)} = 1 - \mathbf{s}_{t1}^{(j)}$. Here, $D_1^{(j)}$ and $D_0^{(j)}$ are constants
 219 denoting the prior beliefs about hidden states. Bayes optimal encoding is obtained when,
 220 and only when, the prior beliefs match the genuine prior distribution; i.e., $D_1^{(j)} = D_0^{(j)} = 0.5$
 221 in this BSS setup. This concludes our treatment of inference about hidden states under this
 222 minimal scheme. Note that the updates in Equation (5) have a biological plausibility in the
 223 sense that the posterior expectations can be associated with nonnegative sigmoid-shape
 224 firing rates (also known as neurometric functions (Tolhurst et al., 1983; Newsome et al.,
 225 1989)), while the arguments of the sigmoid (softmax) function can be associated with
 226 neuronal depolarisation; rendering the softmax function a voltage-firing rate activation
 227 function. Please see (Friston, FitzGerald et al., 2017) for a more comprehensive
 228 discussion—and simulations using this kind of variational message passing to reproduce
 229 empirical phenomena; such as place fields, mismatch negativity responses, phase-precession,
 230 pre-play activity, etc in systems neuroscience.

231 In terms of learning, by solving the variation of F with respect to $\mathbf{a}^{(i,j)}$, the optimal
 232 posterior expectations about the parameters are given by

$$233 \quad \mathbf{a}^{(i,j)} = a^{(i,j)} + \sum_{\tau=1}^t o_{\tau}^{(i)} \otimes \mathbf{s}_{\tau}^{(j)} = a^{(i,j)} + \overline{t o_t^{(i)} \otimes \mathbf{s}_t^{(j)}}, \quad (7)$$

234 where $a^{(i,j)}$ is the prior, $o_{\tau}^{(i)} \otimes \mathbf{s}_{\tau}^{(j)}$ expresses the outer product of a one-hot encoded
 235 vector of $o_{\tau}^{(i)}$ and $\mathbf{s}_{\tau}^{(j)}$, and $\overline{t o_t^{(i)} \otimes \mathbf{s}_t^{(j)}} \equiv \frac{1}{t} \sum_{\tau=1}^t o_{\tau}^{(i)} \otimes \mathbf{s}_{\tau}^{(j)}$. Thus, the optimal posterior
 236 expectation of matrix \mathbf{A} is

$$237 \quad \begin{cases} \mathbf{A}_{11}^{(i,j)} = \frac{\mathbf{a}_{11}^{(i,j)}}{\mathbf{a}_{11}^{(i,j)} + \mathbf{a}_{01}^{(i,j)}} = \frac{\overline{t o_t^{(i)} \mathbf{s}_{t1}^{(j)}} + \mathbf{a}_{11}^{(i,j)}}{t \overline{\mathbf{s}_{t1}^{(j)}} + \mathbf{a}_{11}^{(i,j)} + \mathbf{a}_{01}^{(i,j)}} = \frac{\overline{t o_t^{(i)} \mathbf{s}_{t1}^{(j)}}}{\overline{\mathbf{s}_{t1}^{(j)}}} + \mathcal{O}\left(\frac{1}{t}\right) \\ \mathbf{A}_{10}^{(i,j)} = \frac{\mathbf{a}_{10}^{(i,j)}}{\mathbf{a}_{10}^{(i,j)} + \mathbf{a}_{00}^{(i,j)}} = \frac{\overline{t o_t^{(i)} \mathbf{s}_{t0}^{(j)}} + \mathbf{a}_{10}^{(i,j)}}{t \overline{\mathbf{s}_{t0}^{(j)}} + \mathbf{a}_{10}^{(i,j)} + \mathbf{a}_{00}^{(i,j)}} = \frac{\overline{t o_t^{(i)} \mathbf{s}_{t0}^{(j)}}}{\overline{\mathbf{s}_{t0}^{(j)}}} + \mathcal{O}\left(\frac{1}{t}\right) \end{cases}, \quad (8)$$

238 where $\overline{o_t^{(i)} \mathbf{s}_{t1}^{(j)}} = \frac{1}{t} \sum_{\tau=1}^t o_\tau^{(i)} \mathbf{s}_{\tau1}^{(j)}$, $\overline{\mathbf{s}_{t1}^{(j)}} = \frac{1}{t} \sum_{\tau=1}^t \mathbf{s}_{\tau1}^{(j)}$, $\overline{o_t^{(i)} \mathbf{s}_{t0}^{(j)}} = \frac{1}{t} \sum_{\tau=1}^t o_\tau^{(i)} \mathbf{s}_{\tau0}^{(j)}$, and $\overline{\mathbf{s}_{t0}^{(j)}} =$
 239 $\frac{1}{t} \sum_{\tau=1}^t \mathbf{s}_{\tau0}^{(j)}$. Further, $\mathbf{A}_{01}^{(i,j)} = 1 - \mathbf{A}_{11}^{(i,j)}$ and $\mathbf{A}_{00}^{(i,j)} = 1 - \mathbf{A}_{10}^{(i,j)}$. The prior of parameters
 240 $\alpha^{(i,j)}$ is in the order of 1 and is thus negligible when t is large. The matrix $\mathbf{A}^{(i,j)}$ express the
 241 optimal posterior expectations of $o_t^{(i)}$ taking the ON state when $s_t^{(j)}$ is ON ($\mathbf{A}_{11}^{(i,j)}$) or OFF
 242 ($\mathbf{A}_{10}^{(i,j)}$), or $o_t^{(i)}$ taking the OFF state when $s_t^{(j)}$ is ON ($\mathbf{A}_{01}^{(i,j)}$) or OFF ($\mathbf{A}_{00}^{(i,j)}$). Although this
 243 expression may seem complicated, it is fairly straightforward. The posterior expectations of
 244 the likelihood simply accumulate posterior expectations about the co-occurrence of states
 245 and their outcomes. These accumulated (Dirichlet) parameters are then normalised to give a
 246 likelihood or probability. Crucially, one can observe the associative or Hebbian aspect of this
 247 belief update, expressed here in terms of the outer products between outcomes and
 248 posteriors about states in Equation (7). We now turn to the equivalent update for neural
 249 activities and synaptic weights of a neural network.

250

251 **2.3 Neural activity and Hebbian plasticity models.** Next, we consider the neural activity and
 252 synaptic plasticity in the neural network (Fig. 1B). We assume that the j -th neuron's activity
 253 x_{tj} is given by

$$254 \quad \dot{x}_{tj} \propto \underbrace{-f'(x_{tj})}_{\text{leakage}} + \underbrace{W_{j1}o_t - W_{j0}o_t}_{\text{synaptic input}} + \underbrace{h_{j1} - h_{j0}}_{\text{threshold}}. \quad (9)$$

255 We suppose that $W_{j1} \in \mathbb{R}^{N_o}$ and $W_{j0} \in \mathbb{R}^{N_o}$ comprise row vectors of synapses, and $h_{j1} \in$
 256 \mathbb{R} and $h_{j0} \in \mathbb{R}$ are adaptive thresholds that depend on the values of W_{j1} and W_{j0} ,
 257 respectively. One may regard W_{j1} and W_{j0} as excitatory and inhibitory synapses,
 258 respectively. We further assume that the nonlinear leakage $f'(\cdot)$ (i.e., the leak current) is
 259 the inverse of the sigmoid function (i.e., the logit function), such that the fixed point of x_{tj} is
 260 given by

$$261 \quad x_{tj} = \text{sig}(W_{j1}o_t - W_{j0}o_t + h_{j1} - h_{j0})$$

$$262 \quad = \frac{\exp(W_{j1}o_t + h_{j1})}{\exp(W_{j1}o_t + h_{j1}) + \exp(W_{j0}o_t + h_{j0})}. \quad (10)$$

263 We further assume that synaptic strengths are updated following Hebbian plasticity with an
 264 activity-dependent homeostatic term as follows:

$$265 \quad \begin{cases} \Delta W_{j1}(t) \equiv W_{j1}(t+1) - W_{j1}(t) \propto \text{Hebb}_1(x_{tj}, o_t, W_{j1}) + \text{Home}_1(x_{tj}, W_{j1}) \\ \Delta W_{j0}(t) \equiv W_{j0}(t+1) - W_{j0}(t) \propto \text{Hebb}_0(x_{tj}, o_t, W_{j0}) + \text{Home}_0(x_{tj}, W_{j0}) \end{cases} \quad (11)$$

266 where $Hebb_1$ and $Hebb_0$ denote Hebbian plasticity as determined by the product of
 267 sensory inputs and neural outputs, and $Home_1$ and $Home_0$ denote homeostatic plasticity
 268 determined by output neural activity.

269 In the MDP scheme, posterior expectations about hidden states and parameters are
 270 usually associated with neural activity and synaptic strengths. Here, we can observe a formal
 271 similarity between the solutions for the state posterior (Equation (6)) and the activity in the
 272 neural network (Equation (10)). By this analogy, x_{tj} can be regarded as encoding the
 273 posterior expectation of the ON state $\mathbf{s}_{t1}^{(j)}$. Moreover, W_{j1} and W_{j0} correspond to
 274 $\ln \mathbf{A}_{11}^{(c,j)} - \ln(\vec{\mathbf{1}} - \mathbf{A}_{11}^{(c,j)}) = \text{sig}^{-1}(\mathbf{A}_{11}^{(c,j)})$ and $\ln \mathbf{A}_{10}^{(c,j)} - \ln(\vec{\mathbf{1}} - \mathbf{A}_{10}^{(c,j)}) = \text{sig}^{-1}(\mathbf{A}_{10}^{(c,j)})$,
 275 respectively, in the sense that they express the amplitude of o_t influencing x_{tj} or $\mathbf{s}_{t1}^{(j)}$.

276 Here, $\vec{\mathbf{1}} = (1, \dots, 1) \in \mathbb{R}^{N_o}$ is a vector of ones. In particular, the optimal posterior of a
 277 hidden state taking a value of one (Equation (6)) is given by the ratio of the beliefs about ON
 278 and OFF states, expressed as a sigmoid function. Thus, to be a Bayes optimal encoder, the
 279 fixed point of neural activity needs to be a sigmoid function. This requirement is
 280 straightforwardly ensured when $f'(x_{tj})$ is the inverse of the sigmoid function (see Equation
 281 (13) below). Under this condition, the fixed point or solution for x_{tk} (Equation (10))
 282 compares inputs from ON and OFF pathways, and thus x_{tj} straightforwardly encodes the
 283 posterior of the j -th hidden state being ON (i.e., $x_{tj} \rightarrow \mathbf{s}_{t1}^{(j)}$). In short, the above neural
 284 network is effectively inferring the hidden state.

285 If the activity of the neural network is performing inference, does the Hebbian plasticity
 286 correspond to Bayes optimal learning? In other words, does the synaptic update rule in
 287 Equation (11) ensure that the neural activity and synaptic strengths asymptotically encode
 288 Bayes optimal posterior beliefs about hidden states ($x_{tj} \rightarrow \mathbf{s}_{t1}^{(j)}$) and parameters ($W_{j1} \rightarrow$
 289 $\text{sig}^{-1}(\mathbf{A}_{11}^{(c,j)})$), respectively? To this end, below we will identify a class of cost functions from
 290 which the neural activity and synaptic plasticity can be derived, and consider the conditions
 291 under which the cost function becomes consistent with variational free energy.

292

293 **2.4 Neural network cost functions.** Here, we consider a class of functions that constitute a
 294 cost function for both neural activity and synaptic plasticity. We start by assuming that the
 295 update of the j -th neuron's activity (Equation (9)) is determined by the gradient of cost
 296 function L_j ; i.e., $\dot{x}_{tj} \propto -\partial L_j / \partial x_{tj}$. By integrating the right-hand side of Equation (9), we
 297 obtain a class of cost functions as

$$\begin{aligned}
 298 \quad L_j &= \sum_{\tau=1}^t (f(x_{\tau j}) - x_{\tau j}W_{j1}o_{\tau} - (1 - x_{\tau j})W_{j0}o_{\tau} - x_{\tau j}h_{j1} - (1 - x_{\tau j})h_{j0}) + \mathcal{O}(1) \\
 299 \quad &= \sum_{\tau=1}^t \left(f(x_{\tau j}) - \begin{pmatrix} x_{\tau j} \\ 1 - x_{\tau j} \end{pmatrix}^T \begin{pmatrix} W_{j1} \\ W_{j0} \end{pmatrix} o_{\tau} + \begin{pmatrix} h_{j1} \\ h_{j0} \end{pmatrix} \right) + \mathcal{O}(1), \quad (12)
 \end{aligned}$$

300 where the $\mathcal{O}(1)$ term, which depends on W_{j1} and W_{j0} , is of a lower order than the other
 301 terms (as they are $\mathcal{O}(t)$) and is thus negligible when t is large. Please see Supplementary
 302 Methods S3 for the case where we explicitly evaluate the $\mathcal{O}(1)$ term, to demonstrate the
 303 formal correspondence between the initial values of synaptic strengths and the parameter
 304 prior $p(A)$. The cost function of the entire network is defined by $L \equiv \sum_{j=1}^{N_x} L_j$.
 305 When $f'(x_{\tau j})$ is the inverse of the sigmoid function, we have

$$306 \quad f(x_{\tau j}) = x_{\tau j} \ln x_{\tau j} + (1 - x_{\tau j}) \ln(1 - x_{\tau j}) \quad (13)$$

307 up to a constant term. We further assume that the synaptic weight update rule is derived
 308 from the same cost function L_j . Thus, the synaptic plasticity is given by

$$309 \quad \begin{cases} \dot{W}_{j1} \propto -\frac{1}{t} \frac{\partial L_j}{\partial W_{j1}} = \overline{x_{tj} o_t} + \overline{x_{tj} h'_{j1}} \\ \dot{W}_{j0} \propto -\frac{1}{t} \frac{\partial L_j}{\partial W_{j0}} = \overline{(1 - x_{tj}) o_t} + \overline{1 - x_{tj} h'_{j0}} \end{cases}, \quad (14)$$

310 where $\overline{x_{tj} o_t} \equiv \frac{1}{t} \sum_{\tau=1}^t x_{\tau j} o_{\tau}$, $\overline{x_{tj}} \equiv \frac{1}{t} \sum_{\tau=1}^t x_{\tau j}$, $\overline{(1 - x_{tj}) o_t} \equiv \frac{1}{t} \sum_{\tau=1}^t (1 - x_{\tau j}) o_{\tau}$, $\overline{1 - x_{tj}} \equiv$
 311 $\frac{1}{t} \sum_{\tau=1}^t (1 - x_{\tau j})$, $h'_{j1} \equiv \partial h_{j1} / \partial W_{j1}$, and $h'_{j0} \equiv \partial h_{j0} / \partial W_{j0}$. Note that the update of W_{j1} is
 312 not directly influenced by W_{j0} , and *vice versa*, because they encode parameters in physically
 313 distinct pathways (i.e., the updates are local learning rules (Lee et al., 2000)). The update rule
 314 for W_{j1} can be viewed as Hebbian plasticity mediated by an additional activity-dependent
 315 term expressing homeostatic plasticity. Moreover, the update of W_{j0} can be viewed as
 316 anti-Hebbian plasticity with a homeostatic term, in the sense that W_{j0} is reduced when
 317 input (o_t) and output (x_{tj}) fire together. The fixed points of W_{j1} and W_{j0} are given by

$$318 \quad \begin{cases} W_{j1} = h_1'^{-1} \left(-\frac{\overline{x_{tj} o_t}}{\overline{x_{tj}}} \right) \\ W_{j0} = h_0'^{-1} \left(-\frac{\overline{(1 - x_{tj}) o_t}}{1 - \overline{x_{tj}}} \right) \end{cases}. \quad (15)$$

319 Crucially, these synaptic strength updates are a subclass of the general synaptic plasticity rule
 320 in Equation (11); see also Supplementary Methods S2 for the mathematical explanation.

321 Therefore, if the synaptic update rule is derived from the cost function underlying neural
 322 activity, the synaptic update rule has a biologically plausible form comprising Hebbian
 323 plasticity and activity-dependent homeostatic plasticity.

324

325 **2.5 Comparison with variational free energy.** Here, we establish a formal relationship
 326 between the cost function L and variational free energy. We define $\widehat{W}_{j1} \equiv \text{sig}(W_{j1})$ and
 327 $\widehat{W}_{j0} \equiv \text{sig}(W_{j0})$ as the sigmoid functions of synaptic strengths. We consider the case in
 328 which neural activity is expressed as a sigmoid function and thus Equation (13) holds. As

329 $W_{j1} = \ln \widehat{W}_{j1} - \ln(\vec{1} - \widehat{W}_{j1})$, Equation (12) becomes

$$\begin{aligned}
 330 \quad L = \sum_{j=1}^{N_x} \sum_{\tau=1}^t \begin{pmatrix} x_{\tau j} \\ 1 - x_{\tau j} \end{pmatrix}^T & \left\{ \begin{pmatrix} \ln x_{\tau j} \\ \ln(1 - x_{\tau j}) \end{pmatrix} - \begin{pmatrix} \ln \widehat{W}_{j1} & \ln(\vec{1} - \widehat{W}_{j1}) \\ \ln \widehat{W}_{j0} & \ln(\vec{1} - \widehat{W}_{j0}) \end{pmatrix} \begin{pmatrix} o_{\tau} \\ \vec{1} - o_{\tau} \end{pmatrix} - \begin{pmatrix} h_{j1} \\ h_{j0} \end{pmatrix} \right. \\
 331 & \left. + \begin{pmatrix} \ln(\vec{1} - \widehat{W}_{j1}) \\ \ln(\vec{1} - \widehat{W}_{j0}) \end{pmatrix} \vec{1} \right\} + \mathcal{O}(1), \quad (16)
 \end{aligned}$$

332 where $\vec{1} = (1, \dots, 1) \in \mathbb{R}^{N_o}$. One can immediately see a formal correspondence between
 333 this cost function and variational free energy (Equation (4)). That is, when we assume that
 334 $x_{tj} = \mathbf{s}_{t1}^{(j)}$, $\widehat{W}_{j1} = \mathbf{A}_{11}^{(j)}$, and $\widehat{W}_{j0} = \mathbf{A}_{10}^{(j)}$, Equation (16) has exactly the same form as the
 335 sum of the accuracy and state complexity, which is the leading order term of variational free
 336 energy (see the first term in the last equality of Equation (4)).

337 Specifically, when the thresholds satisfy $h_{j1} = \ln(\vec{1} - \widehat{W}_{j1}) \cdot \vec{1} + \ln D_1^{(j)}$ and $h_{j0} =$
 338 $\ln(\vec{1} - \widehat{W}_{j0}) \cdot \vec{1} + \ln D_0^{(j)}$, Equation (16) becomes equivalent to Equation (4) up to the $\ln t$
 339 order term (that disappears when t is large). Therefore, in this case, the fixed points of neural
 340 activity and synaptic strengths become the posteriors; thus, x_{tj} asymptotically becomes the
 341 Bayes optimal encoder for a large t limit (provided with D that matches the genuine prior
 342 D^*).

343 In other words, we can define perturbation terms $\phi_{j1} \equiv h_{j1} - \ln(\vec{1} - \widehat{W}_{j1}) \cdot \vec{1}$ and
 344 $\phi_{j0} \equiv h_{j0} - \ln(\vec{1} - \widehat{W}_{j0}) \cdot \vec{1}$ as functions of W_{j1} and W_{j0} , respectively, and can express the
 345 cost function as

$$346 \quad L = \sum_{j=1}^{N_x} \sum_{\tau=1}^t \left(\frac{x_{\tau j}}{1-x_{\tau j}} \right)^T \left\{ \begin{pmatrix} \ln x_{\tau j} \\ \ln(1-x_{\tau j}) \end{pmatrix} - \begin{pmatrix} \ln \widehat{W}_{j1} & \ln(\vec{1} - \widehat{W}_{j1}) \\ \ln \widehat{W}_{j0} & \ln(\vec{1} - \widehat{W}_{j0}) \end{pmatrix} \begin{pmatrix} o_{\tau} \\ \vec{1} - o_{\tau} \end{pmatrix} - \begin{pmatrix} \phi_{j1} \\ \phi_{j0} \end{pmatrix} \right\} + \mathcal{O}(1). \quad (17)$$

347 Here, without loss of generality, we can suppose that the constant terms in ϕ_{j1} and ϕ_{j0}
 348 are chosen to ensure that $\exp(\phi_{j1}) + \exp(\phi_{j0}) = 1$. Under this condition,
 349 $(\exp(\phi_{j1}), \exp(\phi_{j0}))$ can be viewed as the prior belief about hidden states

$$350 \quad \begin{cases} \phi_{j1} = \ln D_1^{(j)} \\ \phi_{j0} = \ln D_0^{(j)} \end{cases} \quad (18)$$

351 and thus Equation (17) is formally equivalent to the accuracy and state complexity terms of
 352 variational free energy.

353 This means that when the prior belief about states ($D^{(j)}$) is a function of the parameter
 354 posteriors ($\mathbf{A}^{(j)}$), the generic cost function under consideration can be expressed in the
 355 form of variational free energy, up to the $\mathcal{O}(\ln t)$ term. A generic cost function L is
 356 sub-optimal from the perspective of Bayesian inference unless ϕ_{j1} and ϕ_{j0} are tuned
 357 appropriately to express the unbiased (i.e., optimal) prior belief. In this BSS setup, $\phi_{j1} =$
 358 $\phi_{j0} = \text{const}$ is optimal; thus, a generic L would asymptotically give an upper bound of
 359 variational free energy with the optimal prior belief about states when t is large.

360

361 **2.6 Analysis on synaptic update rules.** To explicitly solve the fixed points of W_{j1} and W_{j0}
 362 that provide the global minimum of L , we suppose ϕ_{j1} and ϕ_{j0} as linear functions of W_{j1}
 363 and W_{j0} , respectively, given by

$$364 \quad \begin{cases} \phi_{j1} = \alpha_{j1} + W_{j1}\beta_{j1} \\ \phi_{j0} = \alpha_{j0} + W_{j0}\beta_{j0} \end{cases}, \quad (19)$$

365 where $\alpha_{j1}, \alpha_{j0} \in \mathbb{R}$ and $\beta_{j1}, \beta_{j0} \in \mathbb{R}^{N_o}$ are constants. By solving the variation of L with
 366 respect to W_{j1} and W_{j0} , we find the fixed point of synaptic strengths as

$$367 \quad \begin{cases} W_{j1} = \text{sig}^{-1} \left(\frac{\overline{x_{tj} o_t}}{\overline{x_{tj}}} + \beta_{j1} \right) \\ W_{j0} = \text{sig}^{-1} \left(\frac{\overline{(1-x_{tj}) o_t}}{1-\overline{x_{tj}}} + \beta_{j0} \right) \end{cases}. \quad (20)$$

368 Since the update from t to $t+1$ is expressed as $\text{sig}(W_{j1} + \Delta W_{j1}) - \text{sig}(W_{j1}) = \widehat{W}_{j1} \odot$
 369 $(\vec{1} - \widehat{W}_{j1}) \odot \Delta W_{j1} + \mathcal{O}(|\Delta W_{j1}|^2)$ and $\text{sig}(W_{j1} + \Delta W_{j1}) - \text{sig}(W_{j1}) \approx x_{(t+1)j} o_{t+1} / \overline{x_{tj}} -$

370 $x_{(t+1)j} \overline{x_{tj} o_t} / \overline{x_{tj}}^2 = x_{(t+1)j} o_{t+1} / \overline{x_{tj}} - (\widehat{W}_{j1} - \beta_{j1}) x_{(t+1)j} / \overline{x_{tj}}$, we recover the following

371 synaptic plasticity:

$$\begin{cases}
 \Delta W_{j1} = \underbrace{\frac{\widehat{W}_{j1}^{\odot -1} \odot (1 - \widehat{W}_{j1})^{\odot -1}}{x_{tj}}}_{\text{adaptive learning rate}} \odot \left\{ \underbrace{x_{(t+1)j} o_{t+1}}_{\text{Hebbian plasticity}} - \underbrace{(\widehat{W}_{j1} - \beta_{j1}) x_{(t+1)j}}_{\text{homeostatic plasticity}} \right\} \\
 \Delta W_{j0} = \underbrace{\frac{\widehat{W}_{j0}^{\odot -1} \odot (1 - \widehat{W}_{j0})^{\odot -1}}{1 - x_{tj}}}_{\text{adaptive learning rate}} \odot \left\{ \underbrace{(1 - x_{(t+1)j}) o_{t+1}}_{\text{anti-Hebbian plasticity}} - \underbrace{(\widehat{W}_{j0} - \beta_{j0})(1 - x_{(t+1)j})}_{\text{homeostatic plasticity}} \right\}
 \end{cases}, \quad (21)$$

373 where \odot denotes the element-wise (Hadamard) product and $\widehat{W}_{j1}^{\odot -1}$ denotes the
 374 element-wise inverse of \widehat{W}_{j1} . This synaptic plasticity rule is a subclass of the generic synaptic
 375 plasticity rule in Equation (11).

376 In summary, we demonstrated that under a few minimal assumptions and ignoring small
 377 contributions to weight updates, the neural network under consideration can be regarded as
 378 minimising an approximation to model evidence, because the cost function can be
 379 formulated in terms of variational free energy. In what follows, we will rehearse our analytic
 380 results and then use numerical analyses to illustrate Bayes optimal inference (and learning)
 381 in a neural network when, and only when, it has the right priors.

382

383 3. Results

384 **3.1 Analytical form of neural network cost functions.** The analysis in the preceding section
 385 rests on the following assumptions:

386 (1) Updates of neural activity and synaptic weights are determined by a gradient descent on
 387 a cost function L .

388 (2) Neural activity is updated by the weighted sum of sensory inputs, and its fixed point is
 389 expressed as the sigmoid function.

390 Under these assumptions, we can express the cost function for a neural network as follows
 391 (see Equation (17)):

$$392 \quad L = \sum_{j=1}^{N_x} \sum_{\tau=1}^t \begin{pmatrix} x_{\tau j} \\ 1 - x_{\tau j} \end{pmatrix}^T \left\{ \begin{pmatrix} \ln x_{\tau j} \\ \ln(1 - x_{\tau j}) \end{pmatrix} - \begin{pmatrix} \ln \widehat{W}_{j1} & \ln(\vec{1} - \widehat{W}_{j1}) \\ \ln \widehat{W}_{j0} & \ln(\vec{1} - \widehat{W}_{j0}) \end{pmatrix} \begin{pmatrix} o_{\tau} \\ \vec{1} - o_{\tau} \end{pmatrix} - \begin{pmatrix} \phi_{j1} \\ \phi_{j0} \end{pmatrix} \right\} + \mathcal{O}(1),$$

393 where $\widehat{W}_{j1} = \text{sig}(W_{j1})$ and $\widehat{W}_{j0} = \text{sig}(W_{j0})$ hold, and ϕ_{j1} and ϕ_{j0} are functions of W_{j1}
 394 and W_{j0} , respectively. The log likelihood function (accuracy term) and divergence of hidden
 395 states (complexity term) of variational free energy emerge naturally under the assumption of
 396 a sigmoid activation function. The cost function above has additional terms denoted by ϕ_{j1}
 397 and ϕ_{j0} . In other words, we can say that the cost function L is variational free energy under

398 a sub-optimal prior belief about hidden states, depending on W_{j1} and W_{j0} : $\ln P(s_t^{(j)}) =$
 399 $\ln D^{(j)} = \phi_j$, where $\phi_j \equiv (\phi_{j1}, \phi_{j0})$. This prior alters the landscape of the cost function in a
 400 sub-optimal manner and thus provides a biased solution for neural activities and synaptic
 401 strengths, which differ from the Bayes optimal encoders.

402 For analytical tractability, we further assume the following:

403 (3) *The perturbation terms (ϕ_{j1} and ϕ_{j0}) that constitute the difference between the cost*
 404 *function and variational free energy with optimal prior beliefs can be expressed as linear*
 405 *equations of W_{j1} and W_{j0} .*

406 From assumption 3, Equation (17) becomes

$$407 \quad L = \sum_{j=1}^{N_x} \left[\sum_{\tau=1}^t \begin{pmatrix} x_{\tau j} \\ 1 - x_{\tau j} \end{pmatrix}^T \left\{ \begin{pmatrix} \ln x_{\tau j} \\ \ln(1 - x_{\tau j}) \end{pmatrix} - \begin{pmatrix} \ln \widehat{W}_{j1} & \ln(\vec{1} - \widehat{W}_{j1}) \\ \ln \widehat{W}_{j0} & \ln(\vec{1} - \widehat{W}_{j0}) \end{pmatrix} \begin{pmatrix} o_{\tau} \\ \vec{1} - o_{\tau} \end{pmatrix} \right. \right. \\
 408 \quad \left. \left. - \begin{pmatrix} \alpha_{j1} + W_{j1}\beta_{j1} \\ \alpha_{j0} + W_{j0}\beta_{j0} \end{pmatrix} \right\} \right] + \mathcal{O}(1), \quad (22)$$

409 where $\{\alpha_{j1}, \alpha_{j0}, \beta_{j1}, \beta_{j0}\}$ are constants. The cost function has degrees of freedom with
 410 respect to the choice of constants $\{\alpha_{j1}, \alpha_{j0}, \beta_{j1}, \beta_{j0}\}$, which correspond to the prior belief
 411 about states $D^{(j)}$. The neural activity and synaptic strengths that give the minimum of a
 412 generic physiological cost function L are biased by these constants, which may be analogous
 413 to physiological constraints (see Discussion for details).

414 The cost function of the neural networks considered is characterised only by ϕ_j . Thus,
 415 after fixing ϕ_j by fixing constrains $(\alpha_{j1}, \alpha_{j0})$ and (β_{j1}, β_{j0}) , the remaining degrees of
 416 freedom are the initial synaptic weights. These correspond to the prior distribution of
 417 parameters $P(A)$ in the variational Bayesian formulation (please see Supplementary
 418 Methods 3).

419 The fixed point of synaptic strengths that give the minimum of L is given analytically as
 420 Equation (20), expressing that (β_{j1}, β_{j0}) deviates the centre of the nonlinear
 421 mapping—from Hebbian products to synaptic strengths—from the optimal position (shown
 422 in Equation (8)). As shown in Equation (14), the derivative of L with respect to W_{j1} and W_{j0}
 423 recovers the synaptic update rules that comprise Hebbian and activity-dependent
 424 homeostatic terms. Although Equation (14) expresses the dynamics of synaptic strengths
 425 that converge to the fixed point, it is consistent with a plasticity rule that gives the synaptic
 426 change from t to $t+1$ (Equation (21)).

427 Hence, based on assumptions 1 and 2, we find that the cost function approximates
 428 variational free energy; see also Supplementary Table S1 for their correspondence. Under
 429 this condition, neural activity encodes the posterior expectation about hidden states, $x_{\tau j} =$

430 $\mathbf{s}_{\tau 1}^{(j)} = Q\left(s_{\tau}^{(j)} = 1\right)$, and synaptic strengths encode the posterior expectation of the
 431 parameters, $\widehat{W}_{j1} = \text{sig}(W_{j1}) = \mathbf{A}_{11}^{(\cdot j)}$ and $\widehat{W}_{j0} = \text{sig}(W_{j0}) = \mathbf{A}_{10}^{(\cdot j)}$. In addition, based on
 432 assumption 3, the accuracy of approximation depends on the deviation of constants
 433 $\{\alpha_{j1}, \alpha_{j0}, \beta_{j1}, \beta_{j0}\}$ from their optimal values. From a Bayesian perspective, these constants
 434 can be viewed as prior beliefs, $\ln P\left(s_t^{(j)}\right) = \ln D^{(j)} = (\alpha_{j1} + W_{j1}\beta_{j1}, \alpha_{j0} + W_{j0}\beta_{j0})$, when
 435 we assume that $(x_{tj}, 1 - x_{tj})$ represents the state posterior $\mathbf{s}_t^{(j)}$. When and only when
 436 $(\alpha_{j1}, \alpha_{j0}) = (-\ln 2, -\ln 2)$ and $(\beta_{j1}, \beta_{j0}) = (\vec{0}, \vec{0})$, the cost function becomes variational
 437 free energy with optimal prior beliefs (for BSS), whose global minimum ensures Bayes
 438 optimal encoding.

439 In short, we identify a class of biologically plausible cost functions from which the update
 440 rules for both neural activity and synaptic plasticity can be derived. When the activation
 441 function for neural activity is a sigmoid function, a cost function in this class is expressed
 442 straightforwardly as variational free energy. With respect to the choice of constants
 443 expressing physiological constraints in the neural network, the cost function has degrees of
 444 freedom that may be viewed as (potentially sub-optimal) prior beliefs from the Bayesian
 445 perspective. Now, we illustrate the implicit inference and learning in neural networks
 446 through simulations of BSS.

447

448 **3.2 Numerical simulations.** Here, we simulated the dynamics of neural activity and synaptic
 449 strengths when they followed a gradient descent on the cost function in Equation (22). We
 450 considered a BSS comprising two hidden sources (or states) and 32 observations (or sensory
 451 inputs), formulated as an MDP. The two hidden sources comprised four patterns: $s_t =$
 452 $s_t^{(1)} \otimes s_t^{(2)} = (0,0), (1,0), (0,1), (1,1)$. An observation $o_t^{(i)}$ was generated through the
 453 likelihood mapping $A^{(i)}$, defined as

$$454 \begin{cases} P(o_t^{(i)} = 1 | s_t, A^{(i)}) = A_{1\cdot}^{(i)} = \left(0, \frac{3}{4}, \frac{1}{4}, 1\right) & \text{for } 1 \leq i \leq 16 \\ P(o_t^{(i)} = 1 | s_t, A^{(i)}) = A_{1\cdot}^{(i)} = \left(0, \frac{1}{4}, \frac{3}{4}, 1\right) & \text{for } 17 \leq i \leq 32 \end{cases} \quad (23)$$

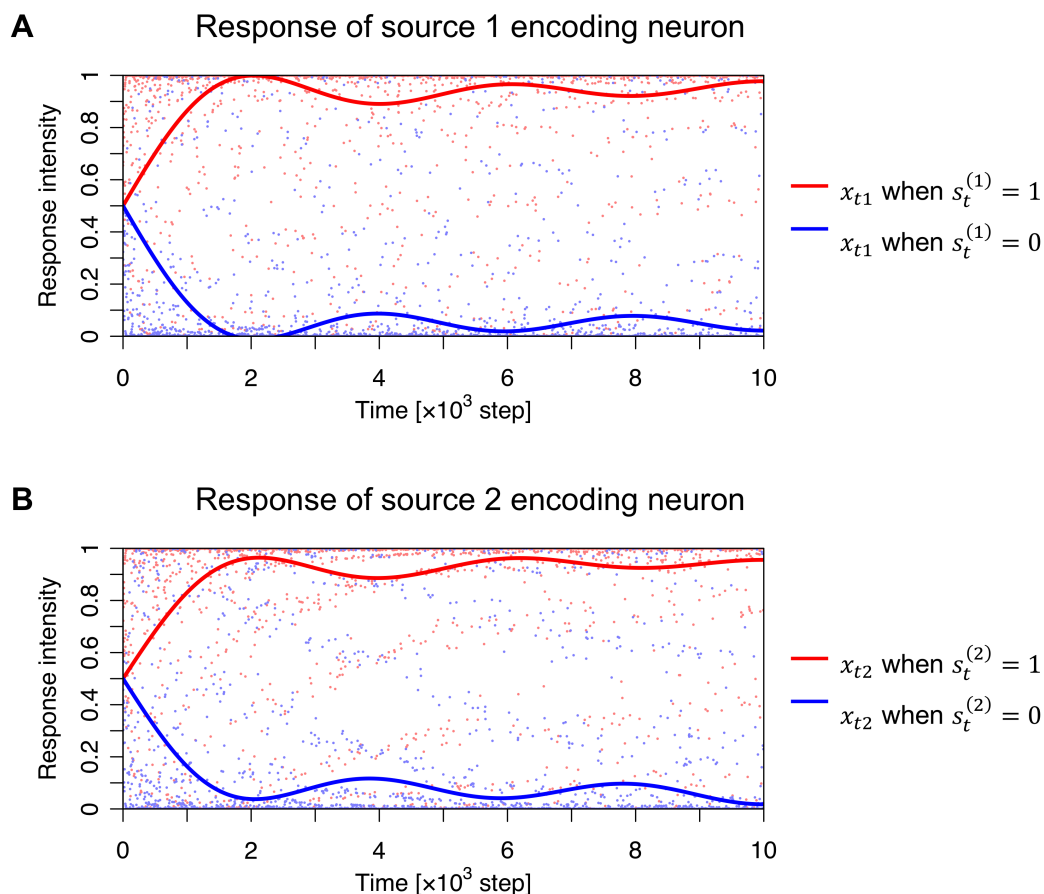
455 Here, for example, $A_{10}^{(i)} = 3/4$ for $1 \leq i \leq 16$ is the probability of $o_t^{(i)}$ taking one when
 456 $s_t = (1,0)$. The simulations continued over $T = 10^4$ time steps. Notably, this simulation
 457 setup is exactly the same experimental setup as that we used for *in vitro* neural networks
 458 (Isomura et al., 2015; Isomura, Friston, 2018). We leverage this setup to clarify the

459 relationship among our empirical work, a feed-forward neural network model, and
460 variational Bayesian formulations.

461 First, as in (Isomura & Friston, 2018), we demonstrated that a network with a cost function
462 with optimised constants ($(\alpha_{j_1}, \alpha_{j_0}) = (-\ln 2, -\ln 2)$ and $(\beta_{j_1}, \beta_{j_0}) = (\vec{0}, \vec{0})$) can perform
463 BSS successfully (Fig. 2). The responses of neuron 1 came to recognise source 1 after training,
464 indicating that neuron 1 learnt to encode source 1 (Fig. 2A). Meanwhile, neuron 2 learnt to
465 infer source 2 (Fig. 2B). This demonstrates that minimisation of the cost function, with
466 optimal constants, is equivalent to variational free energy minimisation, and hence is
467 sufficient to emulate BSS. Next, we quantified the dependency of BSS performance on the
468 form of the cost function, by varying the above-mentioned constants (Fig. 3).

469 We varied $(\alpha_{j_1}, \alpha_{j_0})$ in a range of $0.05 \leq \exp(\alpha_{j_1}) \leq 0.95$, while maintaining
470 $\exp(\alpha_{j_1}) + \exp(\alpha_{j_0}) = 1$, and found that changing $(\alpha_{j_1}, \alpha_{j_0})$ from $(-\ln 2, -\ln 2)$ led to
471 a failure of BSS. Because neuron 1 encodes source 1 with optimal α , the correlation
472 between source 1 and the response of neuron 1 is close to one, while the correlation
473 between source 2 and the response of neuron 1 is nearly zero. In the case of sub-optimal α ,
474 these correlations fall to around 0.5, indicating that the response of neuron 1 encodes a
475 mixture of source 1 and source 2 (Fig. 3A). Moreover, a failure of BSS can be induced when
476 the elements of β take values far from zero (Fig. 3B). When the elements of β are
477 generated from a zero-mean Gaussian distribution, the accuracy of BSS—measured using the
478 correlation between sources and responses—decreases as the standard deviation increases.

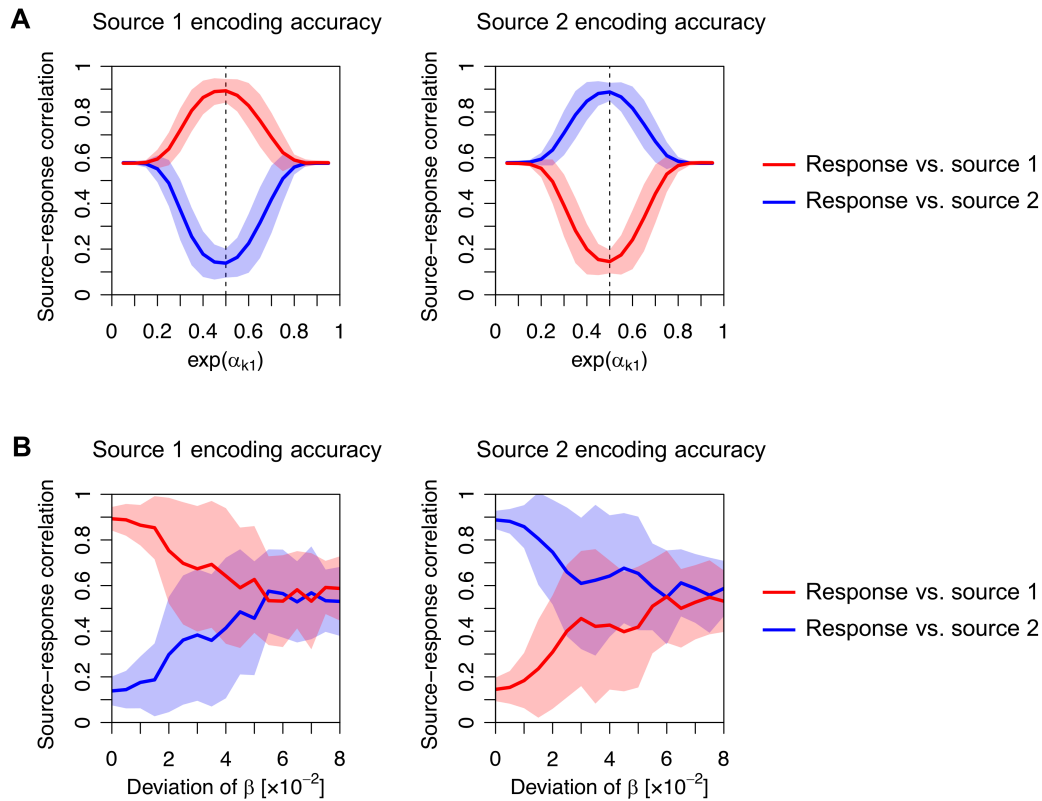
479



480

481 **Figure 2.** Emergence of response selectivity for a source. **(A)** Evolution of neuron 1's
482 responses that learn to encode source 1, in the sense that the response is high when source
483 1 takes a value of one (red dots), and it is low when source 1 takes a value of zero (blue dots).
484 Lines correspond to smoothed trajectories obtained using a discrete cosine transform. **(B)**
485 Emergence of neuron 2's response that learns to encode source 2. These results indicate that
486 the neural network succeeded in separating two independent sources. The code is provided
487 as Supplementary Source Code.

488



489

490 **Figure 3.** Dependence of source encoding accuracy on constants. Left panels show the
 491 magnitudes of the correlations between sources and responses of a neuron expected to
 492 encode source 1: $|\text{corr}(s_t^{(1)}, x_{t1})|$ and $|\text{corr}(s_t^{(2)}, x_{t1})|$. The right panels show the
 493 magnitudes of the correlations between sources and responses of a neuron expected to
 494 encode source 2: $|\text{corr}(s_t^{(1)}, x_{t2})|$ and $|\text{corr}(s_t^{(2)}, x_{t2})|$. **(A)** Dependence on the constant
 495 α that controls the excitability of a neuron, when β is fixed to zero. The dashed line (0.5)
 496 indicates the optimal value of $\exp(\alpha_{j1})$. **(B)** Dependence on constant β , when α is fixed as
 497 $(\alpha_{j1}, \alpha_{j0}) = (-\ln 2, -\ln 2)$. Elements of β were randomly generated from a Gaussian
 498 distribution with zero mean. The standard deviation of β was varied (horizontal axis), where
 499 zero deviation was optimal. Lines and shaded areas indicate the mean and standard
 500 deviation of the source-response correlation, evaluated with 50 different sequences. The
 501 code is provided as Supplementary Source Code.

502

503 Our numerical analysis, under assumptions 1–3 mentioned above, shows that a network
 504 needs to employ a cost function that entails optimal prior beliefs to perform BSS, or
 505 equivalently, causal inference. Such a cost function is obtained when its constants, which do
 506 not appear in the variational free energy with the optimal generative model for BSS, become
 507 negligible. The important message here is that, in this setup, a cost function equivalent to
 508 variational free energy is necessary for Bayes optimal inference (Friston et al., 2006; Friston,

509 2010).

510

511 **3.3 Phenotyping networks.** We have shown that variational free energy (under the MDP
512 scheme) is within the class of biologically plausible cost functions found in neural networks.
513 The neural network's parameters $\phi_j = \ln D^{(j)}$ determine how the synaptic strengths change
514 depending on the history of sensory inputs and neural outputs; thus, the choice of ϕ_j
515 provides degrees of freedom in the shape of the generic cost functions under consideration
516 that determine the purpose or function of the neural network. Among various ϕ_j , only
517 $\phi_j = (-\ln 2, -\ln 2)$ can make the cost function variational free energy with optimal prior
518 beliefs for BSS. Hence, one could regard generic neural networks (of the sort considered in
519 this paper) as performing approximate Bayesian inference under priors that may or may not
520 be optimal. This result is as predicted by the complete class theorem as it implies that any
521 response of a neural network is Bayes optimal under some prior beliefs (and cost function).
522 Therefore, under the theorem, in principle, any neural network of this kind is optimal, when
523 its prior beliefs are consistent with the process that generates outcomes. This perspective
524 indicates the possibility of characterising a neural network model—and indeed a real
525 neuronal network—in terms of its implicit prior beliefs.

526 These considerations raise the possibility of using empirically observed neuronal
527 responses to infer the prior beliefs implicit in a neuronal network. For example, the synaptic
528 matrix (W_{j1}, W_{j0}) can be estimated statistically from response data. By plotting its trajectory
529 over the training period as a function of the history of a Hebbian product, one can estimate
530 the cost function constants. If these constants express a near-optimal ϕ_j , it can be
531 concluded that the network has, effectively, the right sort of priors for BSS. As we have
532 shown analytically and numerically, a cost function with $(\alpha_{j1}, \alpha_{j0})$ far from $(-\ln 2, -\ln 2)$
533 or a large deviation of (β_{j1}, β_{j0}) does not provide the Bayes optimal encoder for
534 performing BSS. Since actual neuronal networks can perform BSS (Isomura et al., 2015;
535 Isomura & Friston, 2018), it can be envisaged that the implicit cost function will exhibit a
536 near-optimal ϕ_j .

537 One can pursue this analysis further and model the responses or decisions of a neural
538 network using the above-mentioned Bayes optimal MDP scheme under different priors. Thus,
539 the priors in the MDP scheme can be adjusted to maximise the likelihood of empirical
540 responses. This sort of approach has been used in system neuroscience to characterise the
541 choice behaviour in terms of subject specific priors. Please refer to (Schwartenbeck & Friston,
542 2016) for further details.

543 Finally, from a practical perspective for optimising neural networks, understanding the
544 formal relationship between cost functions and variational free energy enables us to specify
545 the optimum value of any free parameter to realize some functions. In the present setting,
546 we can effectively optimise the constants by updating the priors themselves, such that they
547 minimise the variational free energy for BSS. Under the Dirichlet form for the priors, the

548 implicit threshold constants of the objective function can then be optimised using the
549 following updates:

$$550 \quad \phi_j = \ln D^{(j)} = \psi(\mathbf{d}^{(j)}) - \psi(\mathbf{d}_1^{(j)} + \mathbf{d}_0^{(j)}),$$

$$551 \quad \mathbf{d}^{(j)} = d^{(j)} + \sum_{\tau=1}^t \mathbf{s}_\tau^{(j)}. \quad (24)$$

552 Please refer to (Schwartenbeck & Friston, 2016) for further details. In effect, this update will
553 simply add the Dirichlet concentration parameters, $\mathbf{d}^{(j)} = (\mathbf{d}_1^{(j)}, \mathbf{d}_0^{(j)})$, to the priors in
554 proportion to the temporal summation of the posterior expectations about the hidden states.
555 Therefore, by committing to cost functions that underlie variational inference and learning,
556 any free parameter can be updated in a Bayes optimal fashion when a suitable generative
557 model is available.

558

559 **4. Discussion**

560 In this work, we investigated a class of biologically plausible cost functions for neural
561 networks. A single-layer feed-forward neural network with a sigmoid activation function that
562 receives sensory inputs generated by hidden states (i.e., BSS setup) was considered. We
563 identified a class of cost functions by assuming that neural activity and synaptic plasticity
564 minimise a common function L . The derivative of L with respect to synaptic strengths
565 furnishes a synaptic update rule following Hebbian plasticity, equipped with
566 activity-dependent homeostatic term. We have shown that the dynamics of a single-layer
567 feed-forward neural network—that minimises its cost function—is asymptotically equivalent
568 to that of variational Bayesian inference under a particular but generic (latent variable)
569 generative model. Hence, the cost function of the neural network can be viewed as
570 variational free energy, and biological constraints that characterise the neural network—in
571 the form of thresholds and neuronal excitability—become prior beliefs about hidden states.
572 This relationship holds regardless of the true generative process of the external world. We
573 have focused on discrete latent variable models that can be regarded as special (reduced)
574 cases of partially observable Markov decision processes (POMDP). However, because our
575 treatment is predicated on the complete class theorem (Brown, 1981; Wald, 1947), the same
576 conclusions should, in principle, be reached when using continuous state space models.
577 Within the class of discrete state space models, it is fairly straightforward to generate
578 continuous outcomes from discrete latent states; as exemplified by discrete variational
579 autoencoders (Rolfe, 2016) or mixed models, as described in (Friston, Parr et al., 2017).

580 One can understand the nature of the constants $\{\alpha_{j1}, \alpha_{j0}, \beta_{j1}, \beta_{j0}\}$ from the biological
581 and Bayesian perspectives as follows: $(\alpha_{j1}, \alpha_{j0})$ determines the firing threshold and thus

582 controls the mean firing rates. In other words, these parameters control the amplitude of
583 excitatory and inhibitory inputs, which may be analogous to the roles of GABAergic inputs
584 and neuromodulators in biological neuronal networks (Pawlak et al., 2010; Frémaux &
585 Gerstner, 2016; Kuśmierz et al., 2017). At the same time, $(\alpha_{j1}, \alpha_{j0})$ encodes prior beliefs
586 about states, which exert a large influence on the state posterior. The state posterior is
587 biased if $(\alpha_{j1}, \alpha_{j0})$ is selected in a sub-optimal manner—in relation to the process that
588 generates inputs. Meanwhile, (β_{j1}, β_{j0}) determines the accuracy of synaptic strengths that
589 represent the likelihood mapping of an observation $o_t^{(i)}$ taking 1 (ON state) depending on
590 hidden states (please compare Equation (8) and Equation (20)). Under a usual MDP setup
591 where the state prior does not depend on the parameter posterior, the encoder becomes
592 Bayes optimal when and only when $(\beta_{j1}, \beta_{j0}) = (\vec{0}, \vec{0})$. These constants can represent
593 biological constraints on synaptic strengths, such as the range of spine growth, spinal
594 fluctuations, or the effect of synaptic plasticity induced by spontaneous activity independent
595 of external inputs. Although the fidelity of each synapse is limited due to such internal
596 fluctuations, the accumulation of information over a large number of synapses should allow
597 accurate encoding of hidden states in the current formulation.

598 In previous reports, we have shown that *in vitro* neural networks—comprising a cortical
599 cell culture—perform BSS when receiving electrical stimulations generated from two hidden
600 sources (Isomura et al., 2015). Furthermore, we showed that minimising variational free
601 energy under an MDP is sufficient to reproduce the learning observed in an *in vitro* network
602 (Isomura & Friston, 2018). Our framework for identifying biologically plausible cost functions
603 could be relevant for identifying the principles that underlie learning or adaptation processes
604 in biological neuronal networks, using empirical response data. Here, we illustrated this
605 potential in terms of the choice of function ϕ_j in the cost functions L . In particular, if ϕ_j is
606 close to a constant $(-\ln 2, -\ln 2)$, the cost function is expressed straightforwardly as a
607 variational free energy with small state prior biases. In the future work, we plan to apply this
608 scheme to empirical data and examine the biological plausibility of variational free energy
609 minimisation.

610 The correspondence highlighted in this work enables one to identify a generative model
611 (comprising likelihood and priors) that a neural network is using. The formal correspondence
612 between neural network and variational Bayesian formations rests on the asymptotic
613 equivalence between the neural network's cost functions and variational free energy (under
614 some priors). Although variational free energy can take an arbitrary form, the
615 correspondence provides biologically plausible constraints for neural networks that implicitly
616 encode prior distributions. Hence, this formulation is potentially useful for identifying the
617 implicit generative models that underlie the dynamics of real neuronal circuits. In other
618 words, one can quantify the dynamics and plasticity of a neuronal circuit in terms of
619 variational Bayesian inference and learning under an implicit generative model.

620 The dependence between the likelihood function and the state prior vanishes when the
621 network uses the optimal threshold to perform inference with a generative process that does
622 not involve dependence between the likelihood and the state prior. In other words, the
623 dependence arises from the sub-optimality of the choice of the state prior. This means that
624 the dependence is due to the degrees of freedom in the choice of the threshold that a neural
625 network and its cost function possess. Nevertheless, minimisation of the cost function can
626 render the network Bayes optimal in the variational Bayesian sense, including the choice of
627 the state prior, as described in the previous section. This is because only variational free
628 energy with the optimal priors provides the minimum among a class of neural network cost
629 functions under consideration.

630 Although we have described the generative process in terms of an MDP, we have ignored
631 state transitions. This means the generative model in this paper reduces to a simple latent
632 variable model, with categorical states and outcomes. As noted above, we refer to MDP
633 models because they predominate in descriptions of variational (Bayesian) belief updating;
634 e.g., (Friston, FitzGerald et al., 2017). Clearly, many generative processes entail state
635 transitions, leading to hidden Markov models (HMM). When state transitions depend upon
636 control variables, we have a POMDP. To deal with such cases, extensions of the current
637 framework are required, which we hope to consider in future work.

638 In summary, we first identified a class of biologically plausible cost functions for neural
639 networks that underlie changes in both neural activity and synaptic plasticity. We then
640 identified an asymptotic equivalence between these cost functions and the cost functions
641 used in variational Bayesian formations. Given this equivalence, changes in the activity and
642 synaptic strengths of a neuronal network can be viewed as Bayesian belief updating; namely,
643 a process of transforming priors over hidden states and parameters into posteriors,
644 respectively. Hence, a cost function in this class becomes Bayes optimal when activity
645 thresholds correspond to appropriate priors in an implicit generative model. In short, the
646 neural and synaptic dynamics of neural networks can be cast as inference and learning,
647 under a variational Bayesian formation. This is potentially important for two reasons. First, it
648 means that there are some threshold parameters for any neural network (in the class
649 considered) that can be optimised for applications to data, when there are precise prior
650 beliefs about the process generating those data. Second, in virtue of the complete class
651 theorem, one can reverse engineer the priors that any neural network is adopting. This may
652 be interesting when real neuronal networks can be modelled using neural networks of the
653 class that we have considered. In other words, if one can fit neuronal responses—using a
654 neural network model parameterised in terms of threshold constants—it becomes possible
655 to evaluate the implicit priors using the above equivalence. This may find a useful application
656 when applied to *in vitro* (or *in vivo*) neuronal networks (Isomura, Friston, 2018; Levin, 2013)
657 or, indeed, dynamic causal modelling of distributed neuronal responses from non-invasive
658 data (Daunizeau et al., 2011). In this context, the neural network can, in principle, be used as
659 a dynamic causal model to estimate threshold constants and implicit priors. This ‘reverse
660 engineering’ speaks to estimating the priors used by real neuronal systems, under ideal

661 Bayesian assumptions; sometimes referred to as meta Bayesian inference (Daunizeau et al.,
662 2010).

663

664 **Acknowledgements**

665 T.I. is funded by RIKEN Center for Brain Science. K.J.F. is funded by a Wellcome Principal
666 Research Fellowship (Ref: 088130/Z/09/Z). The funders had no role in study design, data
667 collection and analysis, decision to publish, or preparation of the manuscript.

668

669 **References**

- 670 Albus, J. S. (1971). A theory of cerebellar function. *Math. Biosci.* **10**, 25-61.
- 671 Bastos, A. M., Usrey, W. M., Adams, R. A., Mangun, G. R., Fries, P. & Friston, K. J. (2012).
672 Canonical microcircuits for predictive coding. *Neuron* **76**, 695-711.
- 673 Belouchrani, A., Abed-Meraim, K., Cardoso, J.F. & Moulines, E. (1997). A blind source
674 separation technique using second-order statistics. *IEEE Trans. Signal Process.* **45**,
675 434-444.
- 676 Brown, G. D., Yamada, S. & Sejnowski, T. J. (2001). Independent component analysis at the
677 neural cocktail party. *Trends Neurosci.* **24**, 54-63.
- 678 Brown, L. D. (1981). A complete class theorem for statistical problems with finite-sample
679 spaces. *Ann. Stat.* **9**, 1289-1300.
- 680 Cichocki, A., Zdunek, R., Phan, A. H. & Amari, S. I. (2009). *Nonnegative Matrix and Tensor*
681 *Factorizations: Applications to Exploratory Multi-way Data Analysis and Blind Source*
682 *Separation*. John Wiley & Sons.
- 683 Comon, P. & Jutten, C. (2010). *Handbook of Blind Source Separation: Independent*
684 *Component Analysis and Applications*. Academic Press.
- 685 Daunizeau, J., David, O., & Stephan, K. E. (2011). Dynamic causal modelling: a critical review
686 of the biophysical and statistical foundations. *Neuroimage* **58**, 312-322.
- 687 Daunizeau, J., Den Ouden, H. E., Pessiglione, M., Kiebel, S. J., Stephan, K. E., & Friston, K. J.
688 (2010). Observing the observer (I): meta-bayesian models of learning and
689 decision-making. *PLoS One* **5**, e15554.
- 690 Dauwels, J. (2007). On variational message passing on factor graphs. *Info. Theory, 2007. ISIT*
691 *2007. IEEE Int. Sympo., IEEE*.
- 692 Dayan, P., Hinton, G. E., Neal, R. M. & Zemel, R. S. (1995). The Helmholtz machine. *Neural*
693 *Comput.* **7**, 889-904.
- 694 DiCarlo, J. J., Zoccolan, D. & Rust, N. C. (2012). How does the brain solve visual object
695 recognition? *Neuron* **73**, 415-434.
- 696 Forney, G. D. (2001). Codes on graphs: Normal realizations. *IEEE Trans. Info. Theory* **47**,
697 520-548.
- 698 Frémaux, N. & Gerstner, W. (2016). Neuromodulated spike-timing-dependent plasticity, and
699 theory of three-factor learning rules. *Front. Neural Circuits* **9**.

- 700 Friston, K. (2005). A theory of cortical responses. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **360**,
701 815-836.
- 702 Friston, K., Kilner, J. & Harrison, L. (2006). A free energy principle for the brain. *J. Physiol.*
703 *Paris* **100**, 70-87.
- 704 Friston, K. (2010). The free-energy principle: a unified brain theory?. *Nat. Rev. Neurosci.* **11**,
705 127-138.
- 706 Friston, K., Mattout, J., & Kilner, J. (2011). Action understanding and active inference. *Biol.*
707 *Cybern.* **104**, 137-160.
- 708 Friston, K., FitzGerald, T., Rigoli, F., Schwartenbeck, P. & Pezzulo, G. (2016). Active inference
709 and learning. *Neurosci. Biobehav. Rev.* **68**, 862-879.
- 710 Friston, K., FitzGerald, T., Rigoli, F., Schwartenbeck, P. & Pezzulo, G. (2017). Active inference:
711 A process theory. *Neural Comput.* **29**, 1-49.
- 712 Friston, K. J., Parr, T. & de Vries, B. D. (2017). The graphical brain: belief propagation and
713 active inference. *Netw. Neurosci.* **1**, 381-414.
- 714 George, D. & Hawkins, J. (2009). Towards a mathematical theory of cortical micro-circuits.
715 *PLoS Comput. Biol.* **5**, e1000532.
- 716 von Helmholtz, H. (1925). *Treatise on physiological optics (Vol. 3)*. The Optical Society of
717 America.
- 718 Isomura, T., Kotani, K. & Jimbo, Y. (2015). Cultured cortical neurons can perform blind source
719 separation according to the free-energy principle. *PLoS Comput. Biol.* **11**, e1004643.
- 720 Isomura, T. & Friston, K. (2018). In vitro neural networks minimise variational free energy. *Sci.*
721 *Rep.* **8**, 16926.
- 722 Knill, D. C. & Pouget, A. (2004). The Bayesian brain: the role of uncertainty in neural coding
723 and computation. *Trends Neurosci.* **27**, 712-719.
- 724 Kullback, S. & Leibler, R. A. (1951). On information and sufficiency. *Ann. Math. Stat.* **22**,
725 79-86.
- 726 Kuśmierz, Ł., Isomura, T. & Toyozumi, T. (2017). Learning with three factors: modulating
727 Hebbian plasticity with errors. *Curr. Opin. Neurobiol.* **46**, 170-177.
- 728 Lee, T. W., Girolami, M., Bell, A. J. & Sejnowski, T. J. (2000). A unifying information-theoretic
729 framework for independent component analysis. *Comput. Math. Appl.* **39**, 1-21.
- 730 Levin, M. (2013). Reprogramming cells and tissue patterning via bioelectrical pathways:
731 molecular mechanisms and biomedical opportunities. *Wiley Interdiscip. Rev. Syst. Biol.*
732 *Med.* **5**, 657-676.
- 733 Linsker, R. (1988). Self-organization in a perceptual network. *Computer* **21**, 105-117.
- 734 Mary, D. (1969). A theory of cerebellar cortex. *J. Physiol.* **202**, 437-470.
- 735 Newsome, W. T., Britten, K. H., & Movshon, J. A. (1989). Neuronal correlates of a perceptual
736 decision. *Nature* **341**, 52-54.
- 737 Pawlak, V., Wickens, J. R., Kirkwood, A. & Kerr, J. N. (2010). Timing is not everything:
738 neuromodulation opens the STDP gate. *Front. Syn. Neurosci.* **2**, 146.
- 739 Rao, R. P. & Ballard, D. H. (1999). Predictive coding in the visual cortex: a functional
740 interpretation of some extra-classical receptive-field effects. *Nat. Neurosci.* **2**, 79-87.
- 741 Schultz, W., Dayan, P., & Montague, P. R. (1997). A neural substrate of prediction and reward.

742 *Science* **275**, 1593-1599.

743 Schwartenbeck, P., & K. Friston. (2016). Computational phenotyping in psychiatry: a worked
744 example. *eNeuro* **3**, e0049-16.2016.

745 Sutton, R. S. & Barto, A. G. (1998). *Reinforcement Learning*. MIT Press, Cambridge, MA, USA.

746 Tolhurst, D. J., Movshon, J. A., & Dean, A. F. (1983). The statistical reliability of signals in
747 single neurons in cat and monkey visual cortex. *Vision Res.* **23**, 775-785.

748 Wald, A. (1947). An essentially complete class of admissible decision functions. *Ann. Math.*
749 *Stat.* **18**, 549-555.

750

751

752

753 **Supplementary Information**

754 **Reverse engineering neural networks to characterise their cost functions**

755 Takuya Isomura, Karl Friston

756

757 **Supplementary Tables**

758

759

Table S1. Correspondence of variables and functions.

Neural network formation		Variational Bayes formation	
Neural activity	x_{tj}	\Leftrightarrow	$\mathbf{s}_{t1}^{(j)}$ State posterior
Sensory input	o_t	\Leftrightarrow	o_t Observation
Synaptic strength	W_{j1}	\Leftrightarrow	$\text{sig}^{-1}(\mathbf{A}_{11}^{(\cdot,j)})$
	\widehat{W}_{j1}	\Leftrightarrow	$\mathbf{A}_{11}^{(\cdot,j)}$ Parameter posterior
Perturbation term	ϕ_{j1}	\Leftrightarrow	$\ln D_1^{(j)}$ State prior
Threshold	h_{j1}	\Leftrightarrow	$\ln(\vec{\mathbf{1}} - \mathbf{A}_{11}^{(\cdot,j)}) \cdot \vec{\mathbf{1}} + \ln D_1^{(j)}$
Initial synaptic strengths	$\lambda_{j1} \odot \widehat{W}_{j1}^{init}$	\Leftrightarrow	$a_{11}^{(\cdot,j)}$ Parameter prior

760

761 **Supplementary Methods**

762 **S1. Order of the parameter complexity**

763 The order of the parameter complexity term

$$764 \quad \mathcal{D}_A \equiv \sum_{i=1}^{N_o} \sum_{j=1}^{N_s} \sum_{l \in \{1,0\}} \left\{ (\mathbf{a}_l^{(i,j)} - a_l^{(i,j)}) \cdot \ln \mathbf{A}_l^{(i,j)} - \ln \mathcal{B}(\mathbf{a}_l^{(i,j)}) \right\} \quad (25)$$

765 is computed. To avoid the divergence of $\ln \mathbf{A}_l^{(i,j)}$, all the elements of $\mathbf{A}_l^{(i,j)}$ are assumed to

766 be larger than a positive constant ε . This means that all the elements of $\mathbf{a}_l^{(i,j)}$ are in the

767 order of t . The first term of Equation (25) becomes $\left(\mathbf{a}_l^{(i,j)} - a_l^{(i,j)}\right) \cdot \ln \mathbf{A}_l^{(i,j)} = \mathbf{a}_l^{(i,j)} \cdot$

768 $\ln \mathbf{A}_l^{(i,j)} + \mathcal{O}(1)$ since $a_l^{(i,j)} \cdot \ln \mathbf{A}_l^{(i,j)}$ is in the order of 1. Moreover, from Equation (3),

769 $\mathbf{a}_l^{(i,j)} \cdot \ln \mathbf{A}_l^{(i,j)} = \mathbf{a}_l^{(i,j)} \cdot \left(\ln \mathbf{a}_l^{(i,j)} - \ln \left(\mathbf{a}_{1l}^{(i,j)} + \mathbf{a}_{0l}^{(i,j)} \right) + \mathcal{O} \left(\left(\mathbf{a}_l^{(i,j)} \right)^{-1} \right) \right) = \mathbf{a}_l^{(i,j)} \cdot$

770 $\ln \left(\mathbf{A}_l^{(i,j)} \right) + \mathcal{O}(1)$. Meanwhile, the second term of Equation (25) comprises the logarithms of

771 gamma functions as $\ln \mathcal{B} \left(\mathbf{a}_l^{(i,j)} \right) = \ln \Gamma \left(\mathbf{a}_{1l}^{(i,j)} \right) + \ln \Gamma \left(\mathbf{a}_{0l}^{(i,j)} \right) - \ln \Gamma \left(\mathbf{a}_{1l}^{(i,j)} + \mathbf{a}_{0l}^{(i,j)} \right)$. From

772 Stirling's formula,

773
$$\Gamma \left(\mathbf{a}_{1l}^{(i,j)} \right) = \sqrt{2\pi} \left(\mathbf{a}_{1l}^{(i,j)} \right)^{-\frac{1}{2}} \left(\frac{\mathbf{a}_{1l}^{(i,j)}}{e} \right)^{\mathbf{a}_{1l}^{(i,j)}} \left(1 + \mathcal{O} \left(\left(\mathbf{a}_l^{(i,j)} \right)^{-1} \right) \right) \quad (26)$$

774 holds. The logarithm of $\Gamma \left(\mathbf{a}_{1l}^{(i,j)} \right)$ is evaluated as

775
$$\begin{aligned} \ln \Gamma \left(\mathbf{a}_{1l}^{(i,j)} \right) &= \frac{1}{2} \ln 2\pi - \frac{1}{2} \ln \mathbf{a}_{1l}^{(i,j)} + \mathbf{a}_{1l}^{(i,j)} \left(\ln \mathbf{a}_{1l}^{(i,j)} - 1 \right) + \ln \left(1 + \mathcal{O} \left(\left(\mathbf{a}_l^{(i,j)} \right)^{-1} \right) \right) \\ 776 &= \mathbf{a}_{1l}^{(i,j)} \ln \mathbf{a}_{1l}^{(i,j)} - \mathbf{a}_{1l}^{(i,j)} + \mathcal{O}(\ln t). \end{aligned} \quad (27)$$

777 Similarly, $\ln \Gamma \left(\mathbf{a}_{0l}^{(i,j)} \right) = \mathbf{a}_{0l}^{(i,j)} \ln \mathbf{a}_{0l}^{(i,j)} - \mathbf{a}_{0l}^{(i,j)} + \mathcal{O}(\ln t)$ and $\ln \Gamma \left(\mathbf{a}_{1l}^{(i,j)} + \mathbf{a}_{0l}^{(i,j)} \right) =$

778 $\left(\mathbf{a}_{1l}^{(i,j)} + \mathbf{a}_{0l}^{(i,j)} \right) \ln \left(\mathbf{a}_{1l}^{(i,j)} + \mathbf{a}_{0l}^{(i,j)} \right) - \left(\mathbf{a}_{1l}^{(i,j)} + \mathbf{a}_{0l}^{(i,j)} \right) + \mathcal{O}(\ln t)$ hold. Thus, we obtain

779
$$\begin{aligned} \ln \mathcal{B} \left(\mathbf{a}_l^{(i,j)} \right) &= \mathbf{a}_{1l}^{(i,j)} \ln \mathbf{a}_{1l}^{(i,j)} + \mathbf{a}_{0l}^{(i,j)} \ln \mathbf{a}_{0l}^{(i,j)} - \left(\mathbf{a}_{1l}^{(i,j)} + \mathbf{a}_{0l}^{(i,j)} \right) \ln \left(\mathbf{a}_{1l}^{(i,j)} + \mathbf{a}_{0l}^{(i,j)} \right) + \mathcal{O}(\ln t) \\ 780 &= \mathbf{a}_l^{(i,j)} \cdot \ln \left(\mathbf{A}_l^{(i,j)} \right) + \mathcal{O}(\ln t). \end{aligned} \quad (28)$$

781 Hence, Equation (25) becomes

782
$$\mathcal{D}_A = \sum_{i=1}^{N_o} \sum_{j=1}^{N_s} \sum_{l \in \{1,0\}} \left\{ \mathbf{a}_l^{(i,j)} \cdot \ln \left(\mathbf{A}_l^{(i,j)} \right) + \mathcal{O}(1) - \left(\mathbf{a}_l^{(i,j)} \cdot \ln \left(\mathbf{A}_l^{(i,j)} \right) + \mathcal{O}(\ln t) \right) \right\} = \mathcal{O}(\ln t). \quad (29)$$

783 Therefore, we obtain

784
$$F(\tilde{o}, Q(\tilde{s}), Q(A)) = \sum_{j=1}^{N_s} \sum_{\tau=1}^t \mathbf{s}_\tau^{(j)} \cdot \left\{ \ln \mathbf{s}_\tau^{(j)} - \sum_{i=1}^{N_o} \ln \mathbf{A}^{(i,j)} \cdot o_\tau^{(i)} - \ln D^{(j)} \right\} + \mathcal{O}(\ln t). \quad (30)$$

785 Under the current generative model comprising binary hidden states and binary
 786 observations, the optimal posterior expectation of \mathbf{A} can be obtained up to the order of
 787 $\ln t / t$ even when the $\mathcal{O}(\ln t)$ term in Equation (30) is ignored. Solving the variation of F
 788 with respect to $\mathbf{A}_{1l}^{(i,j)}$ yields the optimal posterior expectation. From $\mathbf{A}_{0l}^{(i,j)} = 1 - \mathbf{A}_{1l}^{(i,j)}$, we
 789 find

$$\begin{aligned}
 790 \quad \delta F &= \sum_{i=1}^{N_o} \sum_{j=1}^{N_s} \sum_{\tau=1}^t \mathbf{s}_{\tau}^{(j)} \cdot \left\{ -\delta \ln \mathbf{A}_{1\cdot}^{(i,j)} o_{\tau}^{(i)} - \delta \ln \left(\vec{1} - \mathbf{A}_{1\cdot}^{(i,j)} \right) \left(1 - o_{\tau}^{(i)} \right) \right\} \\
 791 &= t \sum_{i=1}^{N_o} \sum_{j=1}^{N_s} \left\{ -\left(\delta \mathbf{A}_{1\cdot}^{(i,j)} \odot \left(\mathbf{A}_{1\cdot}^{(i,j)} \right)^{\odot -1} \right) \cdot \overline{o_t^{(i)} \otimes \mathbf{s}_t^{(j)}} + \left(\delta \mathbf{A}_{1\cdot}^{(i,j)} \odot \left(\vec{1} - \mathbf{A}_{1\cdot}^{(i,j)} \right)^{\odot -1} \right) \cdot \overline{\left(1 - o_t^{(i)} \right) \mathbf{s}_t^{(j)}} \right\} \\
 792 &= t \sum_{i=1}^{N_o} \sum_{j=1}^{N_s} \left(\delta \mathbf{A}_{1\cdot}^{(i,j)} \odot \left(\mathbf{A}_{1\cdot}^{(i,j)} \right)^{\odot -1} \odot \left(\vec{1} - \mathbf{A}_{1\cdot}^{(i,j)} \right)^{\odot -1} \right) \cdot \left(\mathbf{A}_{1\cdot}^{(i,j)} \odot \overline{\mathbf{s}_t^{(j)}} - \overline{o_t^{(i)} \mathbf{s}_t^{(j)}} \right) \quad (31)
 \end{aligned}$$

793 up to the order of $\ln t$. Here, $\left(\mathbf{A}_{1\cdot}^{(i,j)} \right)^{\odot -1}$ denotes the element-wise inverse of $\mathbf{A}_{1\cdot}^{(i,j)}$. From
 794 $\delta F = 0$, we find

$$795 \quad \mathbf{A}_{1\cdot}^{(i,j)} = \overline{o_t^{(i)} \mathbf{s}_t^{(j)}} \odot \left(\overline{\mathbf{s}_t^{(j)}} \right)^{\odot -1} + \mathcal{O} \left(\frac{\ln t}{t} \right). \quad (32)$$

796 Therefore, we obtain the same result as Equation (8) up to the order of $\ln t / t$.

797

798 S2. Derivation of synaptic plasticity rule

799 We consider synaptic strengths at time t , $W_{j1} = W_{j1}(t)$, and define the change as
 800 $\Delta W_{j1} \equiv W_{j1}(t+1) - W_{j1}(t)$. From Equation (15), $h_1'(W_{j1})$ satisfies both

$$801 \quad h_1'(W_{j1} + \Delta W_{j1}) - h_1'(W_{j1}) = h_1''(W_{j1}) \odot \Delta W_{j1} + \mathcal{O} \left(|\Delta W_{j1}|^2 \right) \quad (33)$$

802 and

$$\begin{aligned}
 803 \quad h_1'(W_{j1} + \Delta W_{j1}) - h_1'(W_{j1}) &= -\frac{x_{(t+1)j} o_{t+1} + t \overline{x_{tj} o_t}}{x_{(t+1)j} + t \overline{x_{tj}}} + \frac{\overline{x_{tj} o_t}}{\overline{x_{tj}}} \\
 804 &\approx -\frac{x_{(t+1)j} o_{t+1}}{t \overline{x_{tj}}} + \frac{\overline{x_{tj} o_t}}{t \overline{x_{tj}^2}} x_{(t+1)j} = -\frac{1}{t \overline{x_{tj}}} \left(x_{(t+1)j} o_{t+1} - h_1'(W_{j1}) x_{(t+1)j} \right). \quad (34)
 \end{aligned}$$

805 Thus, we find

$$806 \quad \Delta W_{j1} = \underbrace{-\frac{h_1''(W_{j1})^{\odot-1}}{tx_{tj}}}_{\text{adaptive learning rate}} \odot \left(\underbrace{x_{(t+1)j}o_{t+1}}_{\text{Hebbian term}} - \underbrace{h_1'(W_{j1})x_{(t+1)j}}_{\text{homeostatic term}} \right). \quad (35)$$

807 Similarly,

$$808 \quad \Delta W_{j0} = \underbrace{-\frac{h_0''(W_{j0})^{\odot-1}}{t(1-x_{tj})}}_{\text{adaptive learning rate}} \odot \left(\underbrace{(1-x_{(t+1)j})o_{t+1}}_{\text{anti-Hebbian term}} - \underbrace{h_0'(W_{j0})(1-x_{(t+1)j})}_{\text{homeostatic term}} \right). \quad (36)$$

809 These plasticity rules express (anti-) Hebbian plasticity with a homeostatic term.

810

811 **S3. Correspondence between parameter prior distribution and initial synaptic strengths**

812 In general, optimising a model of observable quantities—including a neural network—can
 813 be cast inference, if there exists a learning mechanism that updates the hidden states and
 814 parameters of that model based on observations. (Exact and variational) Bayesian inference
 815 treats the hidden states and parameters as random variables, and thus transforms prior
 816 distributions $P(s_t), P(A)$ into posteriors $Q(s_t), Q(A)$. In other words, Bayesian inference is
 817 a process of transforming the prior to the posterior based on observations o_1, \dots, o_t under a
 818 generative model. From this perspective, the incorporation of prior knowledge about the
 819 hidden states and parameters is an important aspect of Bayesian inference.

820 The minimisation of a cost function by a neural network updates its activity and synaptic
 821 strengths based on observations under the given network properties (e.g., activation
 822 function and thresholds). According to the complete class theorem, this process can always
 823 be viewed as Bayesian inference. In the main text, we demonstrated that a class of cost
 824 functions—for a single-layer feed-forward network with a sigmoid activation function—has a
 825 form equivalent to variational free energy under a particular latent variable model. Here,
 826 neural activity x_t and synaptic strengths W come to encode the posterior distributions
 827 over hidden states $Q'(s_t)$ and parameters $Q'(A)$, respectively, where $Q'(s_t)$ and $Q'(A)$
 828 follow categorical and Dirichlet distributions, respectively. Moreover, we identified that the
 829 perturbation factors ϕ_j —that characterise the threshold function—correspond to the
 830 logarithm of the state prior $P(s_t)$ expressed as a categorical distribution.

831 However, one might ask whether the posteriors obtained using the network $Q'(s_t), Q'(A)$
 832 are formally different from those obtained using variational Bayesian inference $Q(s_t), Q(A)$,
 833 since only the latter explicitly considers the prior distribution of parameters $P(A)$. Thus, one
 834 may wonder if the network merely influences update rules that are similar to variational
 835 Bayes but does not transform the priors $P(s_t), P(A)$ into the posteriors $Q(s_t), Q(A)$,
 836 despite the asymptotic equivalence of the cost functions.

837 Below, we show that the initial values of synaptic strengths $W_{j_1}^{init}, W_{j_0}^{init}$ correspond to
 838 the parameter prior $P(A)$ expressed as a Dirichlet distribution, to show that a neural
 839 network indeed transforms the priors into the posteriors. For this purpose, we specify the
 840 order 1 term in Equation (12) to make the dependence on the initial synaptic strengths
 841 explicit. Specifically, we modify Equation (12) as

$$\begin{aligned}
 842 \quad L_j &= \sum_{\tau=1}^t \left(f(x_{\tau j}) - \begin{pmatrix} x_{\tau j} \\ 1 - x_{\tau j} \end{pmatrix}^T \left(\begin{pmatrix} W_{j_1} \\ W_{j_0} \end{pmatrix} o_{\tau} + \begin{pmatrix} h_{j_1} \\ h_{j_0} \end{pmatrix} \right) \right) \\
 843 &\quad + (W_{j_1}, W_{j_0}) (\lambda_{j_1} \odot \widehat{W}_{j_1}^{init}, \lambda_{j_0} \odot \widehat{W}_{j_0}^{init})^T \\
 844 &\quad + (\ln(\vec{1} - \widehat{W}_{j_1}), \ln(\vec{1} - \widehat{W}_{j_0})) (\lambda_{j_1}, \lambda_{j_0})^T, \tag{37}
 \end{aligned}$$

845 where $\widehat{W}_{j_1}^{init} \equiv \text{sig}(W_{j_1}^{init})$ and $\widehat{W}_{j_0}^{init} \equiv \text{sig}(W_{j_0}^{init})$ are the sigmoid functions of the initial
 846 synaptic strengths, and $\lambda_{j_1}, \lambda_{j_0} \in \mathbb{R}^{N_o}$ are row vectors of the inverse learning rate factors
 847 that express the insensitivity of the synaptic strengths to the activity-dependent synaptic
 848 plasticity. The third term of Equation (37) expresses the integral of \widehat{W}_{j_1} and \widehat{W}_{j_0} (with
 849 respect to W_{j_1} and W_{j_0} , respectively). This ensures that when $t = 0$ (i.e., when the first term
 850 on the right-hand side of Equation (37) is zero), the derivative of L_j is given by $\partial L_j / \partial W_{j_1} =$
 851 $\lambda_{j_1} \odot \widehat{W}_{j_1}^{init} - \lambda_{j_1} \odot \widehat{W}_{j_1}$, and thus $(W_{j_1}, W_{j_0}) = (W_{j_1}^{init}, W_{j_0}^{init})$ provides the fixed point of
 852 L_j .

853 Similar to the transformation from Equation (12) to Equation (17), we compute Equation
 854 (37) as

$$\begin{aligned}
 855 \quad L &= \sum_{j=1}^{N_x} \sum_{\tau=1}^t \begin{pmatrix} x_{\tau j} \\ 1 - x_{\tau j} \end{pmatrix}^T \left\{ \begin{pmatrix} \ln x_{\tau j} \\ \ln(1 - x_{\tau j}) \end{pmatrix} - \begin{pmatrix} \ln \widehat{W}_{j_1} & \ln(\vec{1} - \widehat{W}_{j_1}) \\ \ln \widehat{W}_{j_0} & \ln(\vec{1} - \widehat{W}_{j_0}) \end{pmatrix} \begin{pmatrix} o_{\tau} \\ \vec{1} - o_{\tau} \end{pmatrix} - \begin{pmatrix} \phi_{j_1} \\ \phi_{j_0} \end{pmatrix} \right\} \\
 856 &\quad + \sum_{j=1}^{N_x} \left\{ (\ln \widehat{W}_{j_1}, \ln(\vec{1} - \widehat{W}_{j_1})) (\lambda_{j_1} \odot \widehat{W}_{j_1}^{init}, \lambda_{j_1} \odot (\vec{1} - \widehat{W}_{j_1}^{init}))^T \right. \\
 857 &\quad \left. + (\ln \widehat{W}_{j_0}, \ln(\vec{1} - \widehat{W}_{j_0})) (\lambda_{j_0} \odot \widehat{W}_{j_0}^{init}, \lambda_{j_0} \odot (\vec{1} - \widehat{W}_{j_0}^{init}))^T \right\}. \tag{38}
 \end{aligned}$$

858 Note that we used $W_{j_1} = \ln \widehat{W}_{j_1} - \ln(\vec{1} - \widehat{W}_{j_1})$. Crucially, analogous to the correspondence
 859 between \widehat{W}_{j_1} and the Dirichlet parameters of the parameter posterior $\mathbf{a}_{11}^{(\cdot, j)}$, $\lambda_{j_1} \odot \widehat{W}_{j_1}^{init}$

860 can be formally associated with the Dirichlet parameters of the parameter prior $a_{11}^{(\cdot,j)}$. Hence,
 861 one can see the formal correspondence between the second and third terms on the
 862 right-hand side of Equation (38) and the expectation of the log parameter prior in Equation
 863 (4):

$$\begin{aligned}
 864 \quad E_{Q(A)}[\ln P(A)] &= \sum_{i=1}^{N_o} \sum_{j=1}^{N_s} \ln \mathbf{A}^{(i,j)} \cdot a^{(i,j)} \\
 865 \quad &= \sum_{i=1}^{N_o} \sum_{j=1}^{N_s} \left\{ \ln \mathbf{A}_{\cdot 1}^{(i,j)} \cdot a_{\cdot 1}^{(i,j)} + \ln \mathbf{A}_{\cdot 0}^{(i,j)} \cdot a_{\cdot 0}^{(i,j)} \right\}. \quad (39)
 \end{aligned}$$

866 Furthermore, the synaptic update rules are derived from Equation (38) as

$$867 \quad \begin{cases} \dot{W}_{j1} \propto -\frac{1}{t} \frac{\partial L}{\partial W_{j1}} = \overline{x_{tj} o_t} - \overline{x_{tj} \widehat{W}_{j1}} + \overline{x_{tj} \phi'_{j1}} + \frac{1}{t} (\lambda_{j1} \odot \widehat{W}_{j1}^{init} - \lambda_{j1} \odot \widehat{W}_{j1}) \\ \dot{W}_{j0} \propto -\frac{1}{t} \frac{\partial L}{\partial W_{j0}} = \overline{(1-x_{tj}) o_t} - \overline{(1-x_{tj}) \widehat{W}_{j0}} + \overline{(1-x_{tj}) \phi'_{j0}} + \frac{1}{t} (\lambda_{j0} \odot \widehat{W}_{j0}^{init} - \lambda_{j0} \odot \widehat{W}_{j0}) \end{cases} \quad (40)$$

868 The fixed point of Equation (40) is provided as

$$869 \quad \begin{cases} W_{j1} = \text{sig}^{-1} \left((\overline{t x_{tj} \mathbf{1}} + \lambda_{j1})^{\odot -1} \odot (\overline{t x_{tj} o_t} + \overline{t x_{tj} \phi'_{j1}} + \lambda_{j1} \odot \widehat{W}_{j1}^{init}) \right) \\ W_{j0} = \text{sig}^{-1} \left((\overline{t(1-x_{tj}) \mathbf{1}} + \lambda_{j0})^{\odot -1} \odot (\overline{t(1-x_{tj}) o_t} + \overline{t(1-x_{tj}) \phi'_{j0}} + \lambda_{j0} \odot \widehat{W}_{j0}^{init}) \right) \end{cases} \quad (41)$$

870 Note that the synaptic strengths at $t = 0$ are computed as $W_{j1} = \text{sig}^{-1} \left((\lambda_{j1})^{\odot -1} \odot$

871 $(\lambda_{j1} \odot \widehat{W}_{j1}^{init}) \right) = W_{j1}^{init}$. Again, one can see the formal correspondence between the final

872 values of the synaptic strengths given by Equation (41) in the neural network formation and

873 the parameter posterior given by Equation (8) in the variational Bayesian formation. As the

874 Dirichlet parameter of the posterior $\mathbf{a}_{11}^{(\cdot,j)}$ is decomposed into the outer product $\overline{o_t} \otimes \mathbf{s}_{t1}^{(j)}$

875 and the prior $a_{11}^{(\cdot,j)}$, they are associated with $\overline{x_{tj} o_t}$ and $\lambda_{j1} \odot \widehat{W}_{j1}^{init}$, respectively. Thus,

876 Equation (8) corresponds to Equation (41). Hence, for a given constant set

877 $\{W_{j1}^{init}, W_{j0}^{init}, \lambda_{j1}, \lambda_{j0}\}$, we identify the corresponding parameter prior $P(A^{(\cdot,j)}) =$

878 $\text{Dir}(a^{(\cdot,j)})$, given by

879
$$\mathbf{a}^{(\cdot,j)} \equiv \begin{pmatrix} a_{11}^{(\cdot,j)} & a_{10}^{(\cdot,j)} \\ a_{01}^{(\cdot,j)} & a_{00}^{(\cdot,j)} \end{pmatrix} = \begin{pmatrix} \lambda_{j1} \odot \widehat{W}_{j1}^{init} & \lambda_{j0} \odot \widehat{W}_{j0}^{init} \\ \lambda_{j1} \odot (\vec{1} - \widehat{W}_{j1}^{init}) & \lambda_{j0} \odot (\vec{1} - \widehat{W}_{j0}^{init}) \end{pmatrix}. \quad (42)$$

880 In summary, one can establish the formal correspondence between neural network and
881 variational Bayesian formations, in terms of the cost functions (Equation (4) vs. Equation
882 (38)), priors (Equation (18) and Equation (42)), and posteriors (Equation (8) vs. Equation (41)).
883 This means that a neural network successively transforms priors $P(s_t), P(A)$ into posteriors
884 $Q(s_t), Q(A)$, as parameterised with neural activity, and initial and final synaptic strengths
885 (and thresholds). Crucially, when increasing number of observations, this process is
886 asymptotically equivalent to that of variational Bayesian inference, under a specific likelihood
887 function.

888

889