

Overview: Backpropagation and Predictive Coding

Samuel Lippl

April 23, 2018

Contents

1	Backpropagation	1
2	Predictive Coding	2
	Open Problems: Overview	
1	Layered structure	1
2	Connection weights as learnable parameters	4
3	Justification of (42)	6
4	Hebbian Σ	6

1 Backpropagation

1.1 We consider an additive feed-forward neural network in layered structure $\theta \in \mathbb{R}_{c,v}^{\mathcal{V} \times \mathcal{V}}$.

1 OPEN PROBLEM (LAYERED STRUCTURE). Construction of layers and the corresponding order on \mathcal{V} .

The step function is therefore

$$T_{\theta}^s = Wf \triangleright s \quad (1)$$

and the canonical network function is

$$\mathbb{R}^{\mathcal{V}^{(0)}} \ni i \mapsto N_{\theta}^i = Wf \triangleright i \downarrow \in \mathbb{R}^{\mathcal{V}} \quad (2)$$

where we denote the corresponding output function by

$$O_{\theta}^i \in \mathbb{R}^{\mathcal{V}^{(L)}} \quad (3)$$

We will often leave θ and i implicit where they are clear from context.

1.2 (BACKPROPAGATION OBJECTIVE FUNCTION) The backpropagation algorithm is concerned with optimizing the prediction that is generated by a neural network with respect to some loss function. We will first consider the quadratic loss function between a prediction O generated from some input and the real output $o \in \mathbb{R}^{\mathcal{V}^{(L)}}$ ¹

$$E_{\theta}^{i;o} := E(O_{\theta}^i, o) := \frac{1}{2}(O_{\theta}^i - o)^T(O_{\theta}^i - o) \quad (4)$$

¹ The given framework is easily extendable to a training set with more than one observation.

θ is therefore modified to reduce E which we do by gradient descent:

$$\theta \mapsto \theta + \alpha \frac{\partial E}{\partial \theta} \quad (5)$$

As the loss function is generally one-dimensional, we can consider $\frac{\partial E}{\partial \theta}$ to be analogous to θ .

1.3 THEOREM (BACKPROPAGATION IDENTITY). Defining

$$\delta_{\text{diff}} := O_{\theta}^i - o \quad (6)$$

and

$$\delta := f'(N) \circ (\theta \delta_{\text{diff}} \mathbb{1}) \quad (7)$$

we obtain the following identity:

$$\frac{\partial E}{\partial \theta} = \delta f(N)^T \quad (8)$$

Proof. (8) follows from a repeated application of the chain rule. We first prove that $\delta = \frac{\partial E}{\partial N}$. We start by observing that

$$\frac{\partial E}{\partial N(\mathcal{V}^{(L)})} = \frac{\partial E}{\partial O} = \delta_{\text{diff}} \quad (9)$$

For $v \notin \mathcal{V}^{(L)}$, we can write recursively (as we have a feed-forward network)

$$\frac{\partial E}{\partial N(v)} = \frac{\partial E}{\partial N(v \rightarrow)} \frac{\partial N(v \rightarrow)}{\partial N(v)} = \frac{\partial E}{\partial N(v \rightarrow)} \theta_{v \rightarrow (v \rightarrow)} \cdot f'(N(v \rightarrow)) \quad (10)$$

As all other entries in $\theta_{v \rightarrow}$ are 0, this corresponds to the recursive definition in (7)².

As

$$\frac{\partial N(w)}{\partial \theta_{v \rightarrow w}} = f(w) \quad (11)$$

and therefore

$$\frac{\partial E}{\partial \theta_{v \rightarrow w}} = \delta(v) f(w) \quad (12)$$

we obtain (8) which proves the theorem. \square

1.4 REMARK. Backpropagation is only defined on feedforward networks, where the proof of theorem 1.3 demonstrates the point at which a recurrent network would fail the algorithm. While we can also formulate the backpropagation algorithm for non-additive neural networks, we do not win much by that. It is the particular additive structure that allows the concise recursive notation.

2 Predictive Coding

I will present the predictive coding model according to [4] in a generalized framework.

2.1 (GENERATIVE MODEL) We first consider a so-called *generative model* $\xi \in \mathbb{R}^{\mathcal{V}_{\text{gen}} \times \mathcal{V}_{\text{gen}}}$ as a neural network which is a stochastic additive feed-forward network that is developed depending on input states $s(\mathcal{V}_{\text{gen}}^{(0)})$ by

$$s(v) = \mu_v(s(\rightarrow v)) + \epsilon_v \quad \epsilon \sim \mathcal{N}(0, \Sigma) \quad (13)$$

where $\Sigma \in \mathbb{R}^{\mathcal{V}_{\text{gen}}}$ is a positive definite covariance matrix and the expected value is

$$\mu_v(s(\rightarrow v)) = \xi_{\rightarrow v}^T s(\rightarrow v) \quad (14)$$

²The multiplication by $f'(N)$ does not have to be part of the recursive definition.

2.2 REMARK. Σ is often a diagonal matrix or only contains correlations across the same layers.

2.3 (OBJECTIVE FUNCTION OF PREDICTIVE CODING) While backpropagation is focused on optimizing the prediction, predictive coding attempts to find the most likely state of the entire generative model. As is often case, we consider the log-likelihood and our objective function is

$$F_{\xi, \Sigma}(s) := \ln P \left(s \left(\mathcal{V} \setminus \mathcal{V}_{\text{gen}}^{(0)} \right) | s \left(\mathcal{V}_{\text{gen}}^{(0)} \right) \right) \quad (15)$$

which we obtain in proposition 2.4.

2.4 PROPOSITION. Optimizing F corresponds to optimizing

$$-\frac{1}{2} (\ln (\det \Sigma) + \epsilon^T \Sigma^{-1} \epsilon) \quad (16)$$

where we have only omitted a constant term. We therefore denote (16) by F .

Proof. As $\epsilon \in \mathbb{R}^{\mathcal{V} \setminus \mathcal{V}^{(0)}}$,

$$F_{\xi, \Sigma}(s) = \ln p(\epsilon | 0, \Sigma) \quad (17)$$

where p is the multivariate normal distribution with expectation 0 and covariance Σ :

$$p(\epsilon | 0, \Sigma) = \frac{1}{\sqrt{(2\pi)^{|\mathcal{V}|} \det \Sigma}} \exp \left(-\frac{1}{2} \epsilon^T \Sigma^{-1} \epsilon \right) \quad (18)$$

Note that F in the form of (17) still depends on s as

$$\epsilon = s - \xi^T f(s) \quad (19)$$

where we omit all $v \in \mathcal{V}^{(L)}$ as their error does not behave stochastically but is trivially zero.

Logarithmizing (18) yields

$$\ln p = -\ln \left(\sqrt{(2\pi)^{|\mathcal{V}|}} \right) - \ln \left(\sqrt{\det \Sigma} \right) - \frac{1}{2} \epsilon^T \Sigma^{-1} \epsilon \quad (20)$$

where we omit the constant term which, together with some basic transformations yields (16). \square

2.5 (INFERENCE) Clearly, if we are given $s(\mathcal{V}^{(0)})$, the F -optimal prediction of all other states would be

$$\theta^T s(\mathcal{V}^{(0)}) \downarrow \quad (21)$$

However, we generally cannot observe $\mathcal{V}^{(0)}$ which corresponds to a cause that may be an abstract concept such as a face and therefore even constructed by our brain. In the general framework, there is a subset $\mathcal{I} \subseteq \mathcal{V}$ of which we know the value and have to infer the state of the remaining network ξ by optimizing F over the PUs $\mathcal{V} \setminus \mathcal{I} =: \mathcal{H}$.

2.6 THEOREM (PREDICTIVE CODING MODEL I: PREDICTION). The neural network with the step function on $s \in \mathbb{R}^{\mathcal{V} \times \mathcal{E}}$

$$T_{\xi}^{(s, \epsilon)}(\mathcal{V}) = s - \epsilon + \text{diag}(f'(s)) \xi \epsilon \quad (22)$$

$$T_{\xi}^{(s, \epsilon)}(\mathcal{E}) = \Sigma^{-1} (s - \xi^T f(s)) \quad (23)$$

performs gradient descent on F . By leaving \mathcal{I} constant, we may perform inference according to 2.5.

Proof. We observe that

$$\frac{\partial F}{\partial N} = \frac{\partial F}{\partial \varepsilon} \frac{\partial \varepsilon}{\partial N} \quad (24)$$

and, as

$$\varepsilon = \Sigma^{-1} \epsilon$$

we can write

$$F = -\frac{1}{2} \varepsilon^T \Sigma \varepsilon \quad (25)$$

Therefore

$$\frac{\partial F}{\partial \varepsilon} = -\varepsilon^T \Sigma \quad (26)$$

and

$$\frac{\partial \varepsilon}{\partial N} = \Sigma^{-1} (I - \xi^T \text{diag}(f'(N))) \quad (27)$$

where $I \in \mathbb{R}^{\mathcal{V} \times \mathcal{V}}$ is the identity matrix and $J = \text{diag}(f'(N))$ is the diagonal matrix where $J_{v \rightarrow v} = f'(N(v))$. We can now put together (24,26,27) to obtain (22) \square

2.7 (LEARNABLE Σ) Theorem 2.6 provides an implicit definition of a neural network which is neither additive nor feed-forward.

2 OPEN PROBLEM (CONNECTION WEIGHTS AS LEARNABLE PARAMETERS). I am still missing the section where I give a reason for the following.

Additionally, we would like to treat Σ as learnable parameters and therefore want to identify them with connection weights. The following definition provides a neural network with a weight matrix that contains all learnable parameters. It is important to note that while the network is not additive, it would be relatively easy to render the network additive by defining additional nodes for f and f' . I do not see the point of such a definition as the network in theorem 2.10 is biologically plausible according to [4].

2.8 DEFINITION (PREDICTIVE CODING MODEL II: PREDICTIVE NETWORK). If we define error nodes for all PUs but the first layer of the generative model

$$\varepsilon := \left\{ \varepsilon_v \mid v \in \mathcal{V}_{\text{gen}} \setminus \mathcal{V}_{\text{gen}}^{(0)} \right\} \quad (28)$$

the matrices

$$I_{\mathcal{V}} = \begin{pmatrix} 1_{\zeta} & & \\ & \ddots & \\ & & 1_{\zeta} \end{pmatrix} \in \mathbb{R}_{\zeta, \mathcal{V}}^{\mathcal{V} \times \mathcal{V}}, I_{\varepsilon} = \begin{pmatrix} 1_{\zeta} & & \\ & \ddots & \\ & & 1_{\zeta} \end{pmatrix} \in \mathbb{R}_{\zeta, \varepsilon}^{\varepsilon \times \varepsilon} \quad (29)$$

the Processing Units

$$\mathcal{V} = \mathcal{V}_{\text{gen}} \cup \varepsilon \quad (30)$$

and $\xi^{\varepsilon} \in \mathbb{R}^{\varepsilon \times \mathcal{V}_{\text{gen}}}$ (resp. $I'_{\mathcal{V}} \in \mathbb{R}^{\varepsilon \times \mathcal{V}_{\text{gen}}}$) as the weight matrix ξ (resp. $I_{\mathcal{V}}$) without the rows of the input layer $\mathcal{V}^{(0)}$, the *predictive coding network* is defined by

$$\theta = \begin{pmatrix} I_{\mathcal{V}} & -\xi^T + I'_{\mathcal{V}}^T \\ -I'_{\mathcal{V}} + \xi & (I_{\varepsilon} - \Sigma)^T \end{pmatrix} \in \mathbb{R}^{\mathcal{V} \times \mathcal{V}} \quad (31)$$

together with the unit functions

$$g_v(s \mid \theta_{\rightarrow v}) := \theta_{\rightarrow v}^T(\mathcal{V}) s(\mathcal{V}) + \theta_{\rightarrow v}^T(\varepsilon) s(\varepsilon) f'(N(v)) \quad (32)$$

$$g_{\varepsilon_v}(s \mid \theta_{\rightarrow \varepsilon_v}) := \theta_{\rightarrow \varepsilon_v}^T(\mathcal{V}) f(s(\mathcal{V})) - f(s(v)) + s(v) + \theta_{\rightarrow \varepsilon_v}^T(\varepsilon) s(\varepsilon) \quad (33)$$

2.9 REMARK. By expressing (32,33) by ξ and Σ instead of θ , their definition and the similarity to (22,23) becomes clear:

$$g_v(s|\theta_{\rightarrow v}) = s(v) - s(\varepsilon) + f'(s(v))\xi_{v \rightarrow}^\varepsilon s(\varepsilon) \quad (34)$$

$$g_v(s|\theta_{\rightarrow \varepsilon_v}) = s(v) - (\xi_{\rightarrow v}^\varepsilon)^T s(\mathcal{V}) + (I_\varepsilon - \Sigma)s(\varepsilon) \quad (35)$$

2.10 THEOREM. The predictive coding network has the local maxima of F as fixed states.

Proof. It is sufficient to prove that the fixed states of θ are identical to the fixed states of the step function defined in (22,23). We use the equations from remark 2.9 where (34) implies (22). For (35), we observe that for some fixed state t ,

$$I_\varepsilon t(\varepsilon) = t(\varepsilon) = t(v) - (\xi_{\rightarrow v}^\varepsilon)^T f(t(\mathcal{V})) + (I_\varepsilon - \Sigma)t(\varepsilon)$$

which is, again, equivalent to the fixed state equation corresponding to 23. \square

2.11 (PREDICTIVE CODING MODEL III: LEARNING PERSPECTIVES) Theorem 2.10 can be taken as evidence that the predictive coding network is a reasonable implementation of the generative model in the brain. We will now turn towards the question how such a network may learn its parameters. With respect to this goal, we have two options: we may consider learning ξ and Σ or we may consider learning θ . While every learning rule on ξ and Σ can be transformed into a learning rule on θ this transformation does not preserve the Hebbian property because of the two applications of ξ , a problem that has been coined *weight symmetry* (see e. g. [4, p. 1254]. On the other hand, not every learning rule on θ can be transformed into a learning rule on ξ and Σ but if there exists such a transformation, it preserves the Hebbian property.

2.12 (PREDICTIVE CODING MODEL IV: LEARNING ξ AND Σ) In this framework, we begin by treating some fixed state of the network as a constant value and try to improve ξ and Σ , again by gradient descent. This corresponds to the approach in [4] and is based on the expectation maximization algorithm [2] as presented in [3]. In this framework, prediction corresponds to the estimation of the expected value of all units and therefore the states of the generative model if we fix the states of certain units. On the other hand, learning ξ and Σ corresponds to changes in the generative model the brain builds to describe its environment (see [3]).

2.13 THEOREM (LEARNING ξ). If

$$\frac{\partial F}{\partial \xi^\varepsilon} \equiv \left(\frac{\partial F}{\partial \xi_{v \rightarrow w}} \right)_{(v,w) \in \mathcal{V} \times (\mathcal{V} \setminus \mathcal{V}^{(0)})} \quad (36)$$

we obtain

$$\frac{\partial F}{\partial \xi^\varepsilon} = N(\varepsilon) f(N(\mathcal{V}))^T \quad (37)$$

Proof.

$$\begin{aligned} \frac{\partial F}{\partial \xi_{v \rightarrow w}} &= -\varepsilon^T \Sigma \frac{\partial \varepsilon}{\partial \xi_{v \rightarrow w}} \\ \frac{\partial \varepsilon}{\partial \xi_{v \rightarrow w}} &= -\Sigma^{-1} \frac{\partial}{\partial \xi_{v \rightarrow w}} \xi^T f(N) \end{aligned}$$

As $\xi^T f(N)$ contains $\xi_{v \rightarrow w}$ exactly once, we have a vector with one non-zero entry. As $\xi_{v \rightarrow w}$ is in the w -row of ξ^T , it is in the w -row of $\xi^T f(N)$. Evidently, the value in the w -row is $f(N(v))$ and therefore

$$\frac{\partial F}{\partial \xi_{v \rightarrow w}} = \varepsilon_w f(N(v)) \quad (38)$$

which, in the notation of (36), yields (37). \square

2.14 (HEBBIAN LEARNING OF ξ) We note that ξ can be learned in a Hebbian way that can be directly motivated from gradient ascent. We also note that the formula resembles the formula for backpropagation but is generally not identical because the network functions may yield different results.

2.15 PROPOSITION (LEARNING Σ I: GRADIENT ASCENT). If

$$\frac{\partial F}{\partial \Sigma} \equiv \left(\frac{\partial F}{\partial \Sigma_{v \rightarrow w}} \right)_{(v,w) \in (\mathcal{V} \setminus \mathcal{V}^{(0)}) \times (\mathcal{V} \setminus \mathcal{V}^{(0)})} \quad (39)$$

we obtain

$$\frac{\partial F}{\partial \Sigma} = \frac{1}{2} \Sigma^{-2} \circ (\epsilon \epsilon^T - \Sigma) \quad (40)$$

Proof. Let $v, w \in \mathcal{V} \setminus \mathcal{V}^{(0)}$.

$$\frac{\partial F}{\partial \Sigma_{v \rightarrow w}} = -\frac{1}{2} \left(\frac{\partial}{\partial \Sigma_{v \rightarrow w}} \ln \det \Sigma + \frac{\partial}{\partial \Sigma_{v \rightarrow w}} \epsilon^T \Sigma^{-1} \epsilon \right) \quad (41)$$

We obtain

$$\frac{\partial}{\partial \Sigma_{v \rightarrow w}} \ln \det \Sigma = \Sigma_{w \rightarrow v}^{-1} = \Sigma_{v \rightarrow w}^{-1} \quad (42)$$

3 OPEN PROBLEM (JUSTIFICATION OF (42)). Reconstruct the [solution](#).

On the other hand,

$$\frac{\partial \epsilon^T \Sigma^{-1} \epsilon}{\partial \Sigma^{-1}} = \epsilon \epsilon^T \quad (43)$$

together with

$$\frac{\partial \Sigma^{-1}}{\partial \Sigma} = -\Sigma^{-2} \quad (44)$$

yields

$$\frac{\partial}{\partial \Sigma_{v \rightarrow w}} \epsilon^T \Sigma^{-1} \epsilon = \frac{\partial}{\partial \Sigma_{v \rightarrow w}^{-1}} \epsilon^T \Sigma^{-1} \epsilon \frac{\partial \Sigma_{v \rightarrow w}^{-1}}{\partial \Sigma_{v \rightarrow w}} = -\varepsilon_v \varepsilon_w (\Sigma^{-2})_{v \rightarrow w} \quad (45)$$

Combining (42,45) yields (40). \square

2.16 (LEARNING Σ I) Gradient ascent is clearly not Hebbian if Σ is not diagonal, that is, if the error terms are correlated. Note that, we cannot obtain ϵ either. However, the direction of the gradient is given by $\epsilon \epsilon^T - \Sigma$ and at least does not make computation of the gradient necessary. [1] proposes a solution for a less general situation.

4 OPEN PROBLEM (HEBBIAN Σ). How can we learn Σ in a Hebbian way without adding more PUs than necessary?

References

- [1] Rafal Bogacz. “A tutorial on the free-energy framework for modelling perception and learning”. In: *Journal of Mathematical Psychology* 76 (2017), pp. 198–211.
- [2] A. P. Dempster, N. M. Laird, and D. B. Rubin. “Maximum Likelihood from Incomplete Data via the EM Algorithm”. In: *Journal of the Royal Statistical Society. Series B (Methodological)* 39.1 (1977), pp. 1–38.
- [3] Karl Friston. “A theory of cortical responses”. In: *Philosophical Transactions of the Royal Society B* 360 (2005), pp. 815–836. DOI: [10.1098/rstb.2005.1622](https://doi.org/10.1098/rstb.2005.1622).
- [4] James C. R. Whittington and Rafal Bogacz. “An Approximation of the Error Backpropagation Algorithm in a Predictive Coding Network with Local Hebbian Synaptic Plasticity”. In: *Neural Computation* 29 (2017), pp. 1229–1262. ISSN: 1530888X. DOI: [10.1162/NECO](https://doi.org/10.1162/NECO.2017.02451). arXiv: [1706.02451](https://arxiv.org/abs/1706.02451). URL: <http://arxiv.org/abs/1706.02451>.