
Biologically Plausible Deep Learning: A Critical Review


Robert T. Lange *

Einstein Center for Neurosciences Berlin

robert.lange17@imperial.ac.uk

www.rob-lange.com

Todo-List for Rob

Figure: Learning performance comparison.	7
Figure: Plot of weight change dynamics/convergence.	7
 Rob: Add BO classic reference	7
Figure: Table of hyperparameters search space for different models	8

1 Introduction

Backpropagation provides a biologically implausible solution to the synaptic credit assignment problem in Deep Learning. While computational graphs and the chain rule successfully provide approximate gradients in deep layered structures, the mere empirical success does not imply that the brain is capable of implementing such a procedure. In this report we review different proposed solutions to such problems that render backpropagation biologically implausible. More specifically, we focus on an approach which implements local plasticity rules in a neural network architecture with dendritic compartments [2]. Previously it has been argued that such an architecture overcomes multiple points of critique while accomplishing similar strong results. Our robustness checks reveal that such a claim is not justified. Deep Learning (DL) has rightfully been the poster child of Machine Learning success in the 21st century. It has dominated competitions and research across all domains (computer vision, natural language processing, robotics as well as computational neuroscience).

*This progress report was submitted as part of the final project of the "Models of Neural Systems" (Winter Term 2018/2019) computer practical course taught and organized by Prof. Richard Kempter (Bernstein Center for Computational Neuroscience, Berlin).

2 Credit Assignment in Deep Layered Structures

In the following section, we will briefly set the ground for our following investigations. More specifically, we introduce all required notation as well as background and problematic considerations regarding backpropagation.² Finally, we briefly discuss the physiological characteristics of pyramidal neurons found in sensory cortices.

Arguably, Deep Learning's most simple layered architecture is the Multi-Layer Perceptron (MLP). A MLP composes multiple layers $\{h_l\}_{l=1}^L$ of non-linear and affine transformations:

$$h_l := f(h_{l-1}; \theta) = \sigma_l(W_l h_{l-1} + b_l)$$

In a classification task the final output layer h_L represents the output distribution over the possible labels. In order to train such a composition one has to define a loss function. A standard classification loss function is given by the cross-entropy between the actual labels distribution, $q(y|x)$, and the output distribution of the network, $p(y|h_L)$:

$$\mathcal{L}(h_L) = - \sum_y q(y|x) \log p(y|h_L)$$

Backpropagation: A Successful Deep Learning and Critical Neuroscience Perspective

In order to train the parameters $\{\theta\}_{l=1}^L$ of a network one makes of powerful auto-differentiation tools and stochastic or batch gradient descent methods.

$$\frac{\partial \mathcal{L}}{\partial \theta_l} = \left(\frac{dh_l}{d\theta_l} \right)^T \frac{\partial \mathcal{L}}{\partial h_l} \quad (1)$$

$$\frac{\partial \mathcal{L}}{\partial h_l} = \left(\frac{dh_{l+1}}{dh_l} \right)^T \frac{\partial \mathcal{L}}{\partial h_{l+1}} \quad (2)$$

$$\frac{dh_{l+1}}{dh_l} = W_{l+1} \text{diag}(\sigma'_{l+1}(W_{l+1}h_l + b_{l+1})) \quad (3)$$

In order to compute the gradient with respect to the parameters of a specific layer l , one has to first compute a forward pass of the network to obtain predictions $\hat{y} := h_L$. Afterwards, one is able to compute a loss-based error signal.

Synaptic Integration via Compartmental Dendrites

²In terms of notation we follow Bartunov et al. [1] who provide a wonderful outline as well as review.

3 Literature Review

Feedback Alignment - Lillicrap et al. [4]

General Content: Introduce a first feedback alignment approach to solve the weight transport problem of backpropagation. Forward and backward weights are modeled separately - backward weights align with weight matrix transpose through learning process. Argument follows from positive definiteness of weight and random matrix product and a rotation line of thought.

Keypoints:

- * Weight transport Problem: downstream errors are fed back to upstream neurons via exact symmetric copy of downstream synaptic weight matrix - neuron "deep" within network has to have precise knowledge of all downstream synapses!

- * Possible solutions: 1. Retrograde transmission of info along axons - problem of slow timescale 2. Feedback of errors via second network - problem of symmetry assumption of feedforward and feedback connections 3. Here: Show that even fixed random connections can allow for learning - symmetry not required! Instead implicit dynamics lead to soft alignment between forward and backward weights

- * Observations: * Feedback weights does not have to be exact: $B \approx W^T$ with $e^T W B e > 0$. rotation within 90 degrees of backprop signal. Learning speed depends on degree! * Alignment of B and W^T via adjustment of W (and B) possible

- * Feedback alignment: * Modulator signal (error-FA) does not impact forward pass post-synaptic activity but acts to alter plasticity at the forward synapses. * FA may encourage W to align with Moore-Penrose pseudoinverse of B - approximate functional symmetry * Inference vs learning - towards bayesian approaches

- * Experiments: * Learns linear function with single hidden layer - learning not slower than backprop * Sigmoid nonlinearity and classification task - altered function of post-synaptic activity - learned also to communicate info when 50% More layers 3 hidden layers - as well as backprop and making use of depth - froze layers and trained alternatingly - positive/negative phase? * Neurons that integrate activity over time and spike stochastically - synchronous pathways

- * Possible Extensions: * Fixed spike thresholds/refractory period * Dropout/stochasticity

Questions:

- * Still signed error signal has to be transferred which remains illusive - see target propagation. * Is result related to Johnson-Lindenstrauss concentration ineq ideas? * Usage of intricate/more complex architectures of communication of backward error - relation to multi-agent RL * Relationship to predictive coding

Target Propagation - Lee et al. [3], Bartunov et al. [1]

Bartunov et al. [1] - Simplified Difference

General Content: Extend the target propagation algorithm to not require exact gradient at penultimate layer. Test alternative learning rules in more complicated settings (CIFAR/ImageNet) and differentiate between locally and fully connected architectures. Very good review but not much additional innovation. Behavioral + Physiological Realism

Keypoints:

- * Problems with backpropagation * Feedback connections require exact copy of feedforward connections = Weight transport * Info propagation does not influence "neural activity" - does not conform to any known biological mechanism

- * Feedback alignment: Use random weights in backward pass to deliver info to earlier layers * Still requires delivery of signed error via distinct pathway * Direct/Broadcast FA - connect feedback from output layer directly to all previous ones

* Contrastive Hebbian Learning/Generalized Recirculation: Use top-down feedback connections to influence neural activity and differences to locally approx gradients * Positive/negative phase - need settling process - Likely to slow for brain to compute in real time

* Target Propagation: Trains distinct set of feedback connections defining backward activity propagation * Connections trained to approximately invert feedforward connections to compute target activities for each layer by successive inversion - decoders * Reconstruction + Forward loss * Different target constructions * Vanilla TP: Target computation via propagation from higher layers' targets backwards through layer-wise inverses * Difference TP: Standard delta rule with additional stabilization from prev reconstruction error. Still needs explicit grad comp at final layer * Not tested on data more complex than MNIST

* Simplified Difference Target Propagation: Computation also for penultimate layer with help of correct label distribution - removes implausible gradient communication * Need diversity in targets - problem of low entropy of classification targets * Need precision in targets - poor inverse learned * Combat both problems/weakness of targets with help of auxiliary output resembling random features from penultimate hidden layer * Parallel vs alternating inverse training - simultaneous more plausible

* Weight-Sharing is not plausible - regularizes by reducing number of free parameters

* Experiments - Mostly negative results: 1. None of existing algos is able to scale up - Good performance MNIST/Somewhat reasonable on CIFAR/Horrible on ImageNet - Seems like weight-sharing is not key to success 2. Need for behavioral realism - judged by performance on difficult tasks 3. Hyperparameter Sensitivity * First fix "good" architecture and then optimize * Use hyperbolic tanh instead of ReLu - work better

Questions:

* How could the brain do weight sharing - is approximate again satisfactory/functional approx? * Think more about communication: MARL agents learning communication channels

4 Local Synaptic Learning Rules with Dendritic Integration

Guerguiev et al. [2] - A Plausible Alternative?

A Hidden layer is described by a set of m three compartmental neurons:

$$\begin{aligned}\text{Apical: } \mathbf{V}^{0a}(t) &= [V_1^{0a}(t), \dots, V_m^{0a}(t)] \\ \text{Basal: } \mathbf{V}^{0b}(t) &= [V_1^{0b}(t), \dots, V_m^{0b}(t)] \\ \text{Somatic: } \mathbf{V}^0(t) &= [V_1^0(t), \dots, V_m^0(t)]\end{aligned}$$

The dynamics of the somatic membrane potential are given by

$$\tau \frac{dV_i^0(t)}{dt} = -V_i^0(t) + \frac{g_b}{g_l} (V_i^{0b}(t) - V_i^0(t)) + \frac{g_a}{g_l} (V_i^{0a}(t) - V_i^0(t))$$

The voltages of the dendritic compartments, on the other hand, are given by the weighted sum of incoming spike trains:

Sacramento et al. [5] - Dendritic Microcircuits

General Content: MLP with simplified dendritic compartments learned in local PE plasticity fashion. No separate phases needed. Errors represent mismatch between pre input from lateral interneurons and top-down feedback. First cortical microcircuit approach. Analytically derive that such a setup/learning rule approximates backprop weight updates and proof basic performance on MNIST.

Keypoints:

- * Hypothesis: Pred errors are encoded at distal dendrites of pyramidal neurons - receive input from downstream neurons - in model: error arise from mismatch of lateral local interneuron inputs (SST - somatostatin) - Learning via local plasticity

- * 3 Compartment Neuron: * Soma + Integration zones: Basal/Apical - convergence of top-down/bottom-up synapses on different compartments - Larkum (2013): Preferred connectivity patterns of cortico-cortical projections

- * 2nd Population within hidden layer - Interneurons = lateral + cross-layer connectivity: cancel t-d input - only backprop errors remain as apical dendrite activity * Predominantly driven by same layer but cross-layer feedback provides weak nudge for interneurons = modeled as conduc-based somatic input current * Modeled as one-to-one between layer interneuron and corresponding upper-layer neuron * Empirically justified by monosynaptic input mapping experiments: weak interneuron teaching signal

- * Neuron/network Model: - Simplifications: 1. Membrane capacity to 1 and resting potential 0; Background activity is white noise 2. Modeling of layer dynamics - where vectors represent units 3. No apical compartment in pyramidal output neurons - 3 compartments seem to suffice as comparison mechanism - Qualitative dynamics: error = apical voltage deflection -> propagates down soma -> modulates somatic firing rate -> plasticity at bottom-up synapses - Somatic conductance acts as nudging conductance - Lateral dendritic projections: interneuron is nudged to follow corresponding next layer pyramidal neuron

- * Synaptic learning rules = Dendritic Predictive Plasticity Rules - Originally: reduction of somatic spiking error - conductance based normalization of lateral projections based on dendritic attenuation factors of different compartments - Implementation requires subdivision of apical compartment into two distal parts (t-d input and lateral input from interneurons)

- * Prev work: * Guerguiev: View apical dendrites as integration zones - temp difference between activity of apical dendrite in presence/absence of teaching input = error inducing plasticity at forward synapses. Used directly for learning b-u synapses without influencing somatic activity. HERE: apical dendrite has explicit error representation by sim integration of t-d excitation and lateral inhibition - No need for separate temporal phases - continuous operation with plasticity always turned on * PC

based work - Whittington and Bogacz: Only plastic synapses are those connecting prediction and error neurons. HERE: all connections plastic - errors are directly encoded in dendritic compartments

* Main Results/Experiments: * Analytic derivation: Somatic MP at layer k integrate feedforward predictions (basal dendritic potentials) and backprop errors (apical dendritic potentials) * Analytic derivation: Plasticity rule converges to backprop weight change with weak feedback limit * Random/Fixed t-d weights = FA * Learned t-d weights minimizing inverse reconstruction loss = TP * Experiments: * Non-Linear regression task: Use soft rectifying nonlinearity as transfer fct - Tons of hyperparameters - injected noise current (dropout/regularization effect?) * MNIST - Deeper architectures: Use convex combination of learning/nudging

* General notes: * Kriegeskorte/DiCarlo/RSA - DNNs outperform alternative frameworks in accurately reproducing activity patterns in cortex - What does this mean? Is DL just extremely flexible/expressive? * bottom-up = feedforward, top-down = feedback

Questions:

* Neural transfer fct = Activation fct! * Again tons of hyperparameters to be chosen - How? * Think of learning (accurate gradient approx) vs architecture (depth, number of hyperparameters) complexity * Different interneuron types (PV = parvalbumin-positive) - different types of errors (generative)

	Backprop (Rummelhart et al., 1986)	Feedback Alignment (Lillicrap et al., 2016)	Target Propagation (LeCun, 1986)	Difference TP (Lee et al., 2015)	Simplified DTP (Bartunov et al., 2018)	Segregated Compartments (Guergiev et al., 2017)	Microcircuits (Sacramento et al., 2018)
Exact Gradients	✓	✗	✗	✗	✗	✗	✓ (In Limit)
No Weight Transport	✗	✓	✓/✗ (Final Layer)	✓/✗ (Final Layer)	✓	✓	✓
No Separate Pathways	✓	✗	✗	✗	✗	✓	✓
Dendritic Integration	-	✗	✗	✗	✗	✓	✓
Separate Weights Learned	-	✗	✓	✓	✓	✗	✓/✗
Linear Stabilization	-	-	✗	✓	✓	-	-
Explicit Error Representation	✓	✓	✓	✓	✓	✗	✓

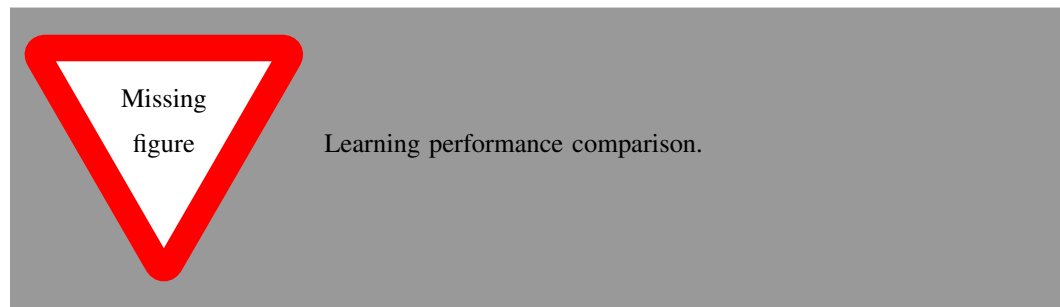
Figure 1: Literature Review.

5 Empirical Investigations

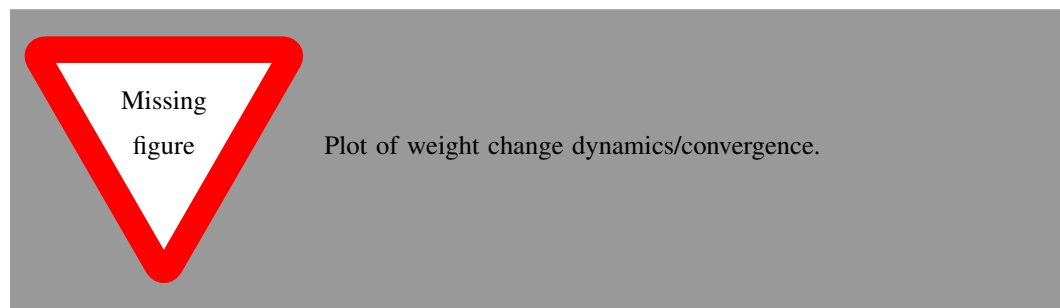
Scalability Across Datasets



Figure 2: Illustration of the 10 different classes/labels of the analyzed datasets. **Top Row:** MNIST dataset. Data format: $70000 \times 1 \times 28 \times 28$. **Middle Row:** Fashion-MNIST dataset. Data format: $70000 \times 1 \times 28 \times 28$. **Bottom Row:** CIFAR-10 dataset. Data format: $60000 \times 3 \times 32 \times 32$. From top to bottom the intra-class variability/entropy increases significantly. We normalize the pixel values to lie within $[0, 1]$ and reshape the images into vector format (e.g. $X \in [0, 1]^{784}$) before training the classifiers. This helps dealing with erratic gradient behavior.



Learning Dynamics

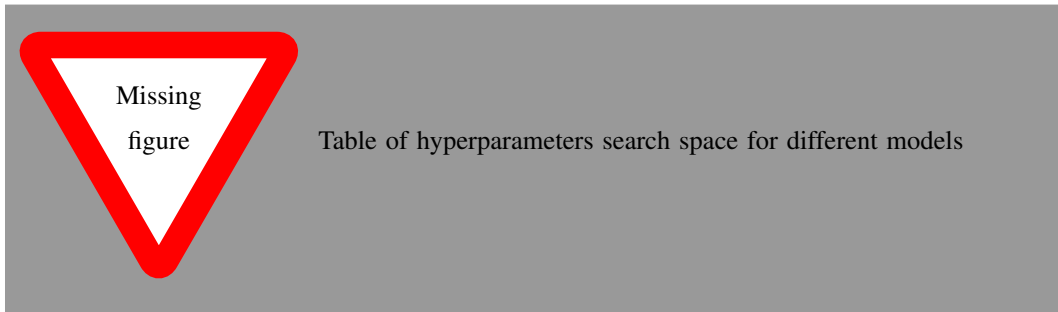


Hyperparameter Robustness

Bayesian Optimization is a probabilistic technique commonly used to optimize the parameters of complex functions which are costly to evaluate. Computing cross-validated test accuracies of deep

Rob: Add
BO classic
reference

networks is one such costly function. Instead of randomly searching through the hyperparameter space, one approximates the loss function $\mathcal{L}^{k-fold}(\theta|X, y)$ with the help of a Gaussian Process (GP).



MLP Hyperparameter Search Space:

1. Layer 1

CNN Hyperparameter Search Space:

1. Layer 1

Guerguiev et al. [2] Hyperparameter Search Space:

1. Layer 1

6 Outlook and Conclusion

In this report we have empirically investigated the robustness and learning dynamics of an alternative learning rule in deep layered structures. We first reviewed and formalized the classical backpropagation algorithm. Afterwards, we put on computational neuroscience goggles and highlighted several short-comings such as the weight transport problem as well the necessity to propagate signed errors. In Section 3 of this report we then introduced the methodology outlined by Guerguiev et al. [2] which intends to overcome such limitations. Inspired by dendritic compartments and information integration at different sites, the algorithm solves the weight transport problem. In Section 4 we reviewed more current approaches and compared their benefits and limitations. Thereby, we highlight the difference between behavioral and neurophysiological realism. Furthermore, we discuss the differences between learning and architecture complexity across the different approaches. Afterwards, we implement the approach by Guerguiev et al. [2] and compare model selection as well as hyperparameter robustness across different popular datasets. Our experiments reveal major performance decreases. This brings up the following question: Why should the brain implement a suboptimal **and** non-robust learning rule on a neurophysiological level? A simple answer to this is the flexibility that such an alternative architecture comes with.

References

- [1] BARTUNOV, S., A. SANTORO, B. RICHARDS, L. MARRIS, G. E. HINTON, AND T. LILLICRAP (2018): “Assessing the scalability of biologically-motivated deep learning algorithms and architectures,” in *Advances in Neural Information Processing Systems*, 9389–9399.
- [2] GUERGUIEV, J., T. P. LILLICRAP, AND B. A. RICHARDS (2017): “Towards deep learning with segregated dendrites,” *ELife*, 6, e22901.
- [3] LEE, D., S. ZHANG, A. BIARD, AND Y. BENGIO (2014): “Target Propagation,” *CoRR*, abs/1412.7525.
- [4] LILLICRAP, T. P., D. COWNDEN, D. B. TWEED, AND C. J. AKERMAN (2016): “Random synaptic feedback weights support error backpropagation for deep learning,” *Nature communications*, 7, 13276.
- [5] SACRAMENTO, J., R. P. COSTA, Y. BENGIO, AND W. SENN (2018): “Dendritic cortical microcircuits approximate the backpropagation algorithm,” in *Advances in Neural Information Processing Systems*, 8735–8746.

Supplementary Material

Biological Plausible Deep Learning

Author: Robert Tjarko Lange | December 2018

This project analyzes different learning rules in deep layered structures. More specifically, we explore alternatives to backpropagation (aka the chain rule). Weight transport (access to all weights at every layer of the backward pass) renders backpropagation biologically implausible. Recent alternatives explore local learning rules and draw inspiration from the compartmental design of pyramidal neurons.

Repository Structure

```
Bio-Plausible-DeepLearning
+- workspace.ipynb: Main workspace notebook - Execute for replication
```

How to use this code

1. Clone the repo.

```
git clone https://github.com/RobertTLange/Bio-Plausible-DeepLearning
cd Bio-Plausible-DeepLearning
```

2. Create a virtual environment (optional but recommended).

```
virtualenv -p python BPDFL
```

Activate the env (the following command works on Linux, other operating systems might differ):

```
source BPDFL/bin/activate
```

3. Install all dependencies:

```
pip install -r requirements.txt
```

4. Run the main notebook:

```
jupyter notebook workspace.ipynb
```