# Online Learning with Gated Linear Networks

**Joel Veness**[†]                                                       AIXI@GOOGLE.COM
**Tor Lattimore**[†]                                        LATTIMORE@GOOGLE.COM
**Avishkar Bhoopchand**[†]                                AVISHKAR@GOOGLE.COM
**Agnieszka Grabska-Barwinska**                                AGNIGB@GOOGLE.COM
**Christopher Mattern**                                     CMATTERN@GOOGLE.COM
**Peter Toth**                                               PETERTOTH@GOOGLE.COM
*Google, DeepMind*

## Abstract

This paper describes a family of probabilistic architectures designed for online learning under the logarithmic loss. Rather than relying on non-linear transfer functions, our method gains representational power by the use of data conditioning. We state under general conditions a learnable capacity theorem that shows this approach can in principle learn any bounded Borel-measurable function on a compact subset of euclidean space; the result is stronger than many universality results for connectionist architectures because we provide both the model and the learning procedure for which convergence is guaranteed.

**Keywords:** Compression, online learning, geometric mixing, logarithmic loss.

## 1. Introduction

This paper explores the use of techniques from the online learning and data compression communities for the purpose of high dimensional density modeling, with a particular emphasis on image density modeling. The current state of the art is almost exclusively dominated by various deep learning based approaches from the machine learning community. These methods are trained offline, can generate plausible looking samples, and generally offer excellent empirical performance. But they also have some drawbacks. Notably, they require multiple passes through the data (not online) and do not come with any kind of meaningful theoretical guarantees. Of course density modeling is not just a machine learning problem; in particular, it has received considerable study from the statistical data compression community. Here the emphasis is typically on *online* density estimation, as this in conjunction with an adaptive arithmetic encoder (Witten et al., 1987) obviates the need to encode the model parameters along with the encoded data when performing compression. Our main contribution is to show that a certain family of neural networks, composed of techniques originating from the data compression and online learning communities, can in some circumstances match the performance of deep learning based approaches in just a single pass through the data, while also enjoying universal source coding guarantees.

---

. † denotes joint first authorship.

## 1.1 Related Work

### 1.1.1 ONLINE, NEURAL, DENSITY ESTIMATION

Although neural networks have been around for a long time, their application to data compression has been restricted until relatively recently due to hardware limitations. To the best of our knowledge, the first promising empirical results were presented by Schmidhuber and Heil (1996), who showed that an offline trained 3-layer neural network could outperform Lempel-Ziv based approaches on text prediction; interestingly online prediction was noted as a promising future direction if the computational performance issues could be addressed.

Building on the work of Schmidhuber and Heil (1996), Mahoney (2000) introduced a number of key algorithmic and architectural changes tailored towards computationally efficient online density estimation. Rather than use a standard fully connected MLP, neurons within a layer were partitioned into disjoint sets; given a single data point, hashing was used to select only a single neuron from each set in each layer, implementing a kind of hard gating mechanism. Since only a small subset of the total weights were used for any given data point, a speed up of multiple orders of magnitude was achieved. Also, rather than using backpropagation, the weights in each neuron were adjusted using a local learning rule, which was found to work better in an online setting.

The next major enhancement was the introduction of *context-mixing*, culminating in the PAQ family of algorithms (Mahoney, 2013). At a high level, the main idea is to combine the predictions of multiple history-based models using a network architecture similar in spirit to the one described above. Common choices of models to combine include $n$-grams, skip-grams (Guthrie et al., 2006), match models (Mahoney, 2013), Dynamic Markov Coding (Cormack and Horspool, 1987), as well as specialized models for commonly occurring data sources. Many variants along this theme were developed, including improvements to the network architecture and weight update mechanism, with a full history given in the book by Mahoney (2013). At the time of writing, context mixing variants achieve the best compression ratios on widely used benchmarks such as the Calgary Corpus (Bell et al., 1990), Canterbury Corpus (Bell and Arnold, 1997), and Wikipedia (Hutter, 2017).

The excellent empirical performance of these methods has recently motivated further investigation into understanding why these techniques seem to perform so well and whether they have applicability beyond data compression. Knoll and de Freitas (2012) investigated the PAQ8 variant (Mahoney, 2013) from a machine learning perspective, providing a useful interpretation of the as a kind of *mixture of experts* (Jacobs et al., 1991) architecture, as well as providing an up-to-date summary of the algorithmic developments since the tech report of Mahoney (2005). They also explored the performance of PAQ8 in a variety of standard machine learning settings including text categorization, shape recognition and classification, showing that competitive performance was possible. One of the key enhancements in PAQ7 was the introduction of logistic mixing, which applied a logit based non-linearity to a neurons input; Mattern (2012, 2013) generalized this technique to non-binary alphabets and justified this particular form of input combination via a KL-divergence minimization argument. Furthermore he showed that due to the convexity of the instantaneous loss (in the single neuron case), the local gradient based weight update used within PAQ7 and later versions is justified in a regret-minimizing sense with respect to piece-wise stationary sources. Subsequent work (Mattern, 2016) further analyzed the regret properties of many

of the core building blocks used to construct many of the base models used within PAQ models, but a theoretical understanding of the mixing network with multiple interacting geometric mixing neurons remained open.

### 1.1.2 Batch Neural Density Estimation

Generative image modeling is a core topic in unsupervised learning and has been studied extensively over the last decade from a deep learning perspective. The topic is sufficiently developed that there are now practical techniques for generating realistic looking natural images (Larochelle and Murray, 2011b; Gregor et al., 2015; Van Den Oord et al., 2016; Chen et al., 2016; Gulrajani et al., 2016, and others). Modeling images is necessarily a high dimensional problem, which precludes the use of many traditional density estimation techniques on both computational and statistical grounds. Although image data is high dimensional, it is also quite structured, and is well suited to deep learning approaches which can easily incorporate appropriate notions of locality and translation invariance into the architecture.

The aforementioned methods all require multiple passes over the dataset to work well, which makes it significantly more challenging to incorporate them into naturally online settings such as reinforcement learning, where density modeling is just starting to see promising applications (Veness et al., 2015; Ostrovski et al., 2017).

### 1.1.3 Miscellaneous

Our work shares some similarity to that of Balduzzi (2016), who used techniques from game theory and online convex programming to analyze some convergence properties of deep convolutional ReLU neural networks by introducing a notion of gated regret. While extremely general, this notion of gated regret is defined with respect to an oracle that knows which ReLU units will fire for any given input; while this is sufficient to guarantee convergence, it does not provide any guarantees on the quality of the obtained solution. Our work fundamentally differs in that we explicitly construct a family of gated neural architectures where it is easy to locally bound the usual notion of regret, allowing us to provide far more meaningful theoretical guarantees in our restricted setting.

Foerster et al. (2017) recently introduced a novel recurrent neural architecture whose modeling power was derived from using a data-dependent affine transformation as opposed to a non-linear activation function. As we shall see later, this is similar in spirit to our approach; we use a product of data-dependent weight matrices to provide representation power instead of a non-linear activation function. Our work differs in that we consider an online setting, and use a local learning rule instead of backpropagation to adjust the weights.

### 1.2 Contribution

Our main contributions are:

- to introduce a family of neural models, Gated Linear Networks, which consist of a sequence of data dependent linear networks coupled with an appropriate choice of

gating function; we show that the high performance PAQ models are a special case of this framework, opening up many potential avenues for improvement;

- to theoretically justify the local weight learning mechanism in these architectures, and highlight how the structure of the loss allows many other techniques to give similar guarantees. Interestingly, while the original motivation for introducing gating was computational (Mahoney, 2000; Knoll and de Freitas, 2012), we show that this technique also adds meaningful representational power as well;

- introduce an adaptive regularization technique which allows a Gated Linear Network to have competitive loss guarantees with respect to all possible sub-networks obtained via pruning the original network;

- to provide an *effective* capacity theorem showing that with an appropriate choice of gating function, a sufficiently large network size and a Hannan consistent/no-regret (Cesa-Bianchi and Lugosi, 2006) local learning method, Gated Linear Networks *will* learn any continuous density function to an arbitrary level of accuracy.
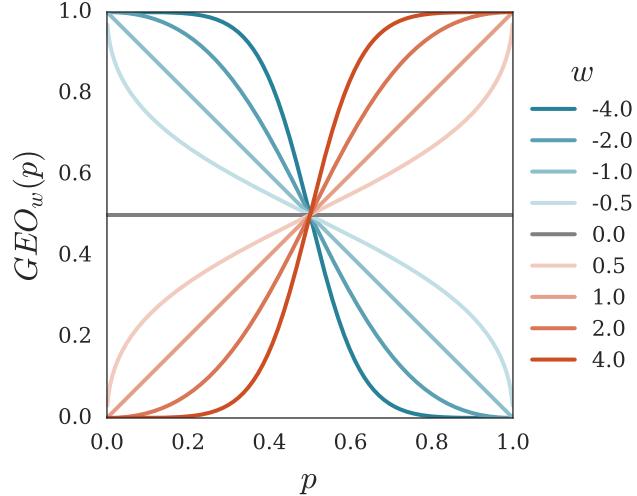
### 1.3 Paper Outline

After introducing the required notation, we proceed by describing a notion of "neuron" suitable for our purposes, which we will call a *gated geometric mixer*; this will form the elementary building block for all of our subsequent work. We then proceed to describe *gated linear networks*, which are feed-forward architectures composed of multiple layers of such gated geometric mixing neurons. Our theoretical results are then presented, which include a capacity analysis of this family of models. A technique for adaptive regularization, *sub-network switching*, is then introduced, along with its associated theoretical properties. Next we describe an architecture built from these principles, reporting empirical results on MNIST, both classification and density modelling.

## 2. Gated Geometric Mixing

This section introduces the basic building block of our work, the *gated geometric mixer*. We start with some necessary notation, before describing (ungated) *geometric mixing* and some of its properties. We then generalize this technique to the gated case, giving rise to the gated geometric mixer.

### 2.1 Notation

Let $\Delta_d = \{x \in [0,1]^d : \|x\|_1 = 1\}$ be the $d$ dimensional probability simplex embedded in $\mathbb{R}^{d+1}$ and $\mathcal{B} = \{0,1\}$ the set of binary elements. The indicator function for set $A$ is $\mathbb{1}_A$ and satisfies $\mathbb{1}_A(x) = 1$ if $x \in A$ and $\mathbb{1}_A(x) = 0$ otherwise. For predicate $P$ we also write $\mathbb{1}[P]$, which evaluates to 1 if $P$ is true and 0 otherwise. The scalar element located at position $(i,j)$ of a matrix $A$ is $A_{ij}$, with the $i$th row and $j$ column denoted by $A_{i*}$ and $A_{*j}$ respectively. For functions $f : \mathbb{R} \to \mathbb{R}$ and vectors $x \in \mathbb{R}^d$ we adopt the convention of writing $f(x) \in \mathbb{R}^d$ for the coordinate-wise image of $x$ under $f$ so that $f(x) = (f(x_1), \ldots, f(x_d))$. If $p, q \in [0,1]$, then $\mathcal{D}(p,q) = p \log p/q + (1-p) \log(1-p)/(1-q)$ is the Kullback-Leibler (KL)

Figure 1: Role of $w$ in a one-dimensional geometric mixer with input probability $p$.

divergence between Bernoulli distributions with parameters $p$ and $q$ respectively. Let $\mathcal{X}$ be a finite, non-empty set of symbols, which we call the alphabet. A string of length $n$ over $\mathcal{X}$ is a finite sequence $x_{1:n} = x_1 x_2 \ldots x_n \in \mathcal{X}^n$ with $x_t \in \mathcal{X}$ for all $t$. For $t \leq n$ we introduce the shorthands $x_{<t} = x_{1:t-1}$ and $x_{\leq t} = x_{1:t}$. The string of length zero is $\epsilon$ and the set of all finite strings is $\mathcal{X}^* = \{\epsilon\} \cup \bigcup_{i=1}^{\infty} \mathcal{X}^i$. The concatenation of two strings $s, r \in \mathcal{X}^*$ is denoted by $sr$. A sequential, probabilistic model $\rho$ is a probability mass function $\rho : \mathcal{X}^* \to [0, 1]$, satisfying the constraint that $\rho(x_{1:n}) = \sum_{y \in \mathcal{X}} \rho(x_{1:n}y)$ for all $n \in \mathbb{N}$, $x_{1:n} \in \mathcal{X}^n$, with $\rho(\epsilon) = 1$. Under this definition, the conditional probability of a symbol $x_n$ given previous data $x_{<n}$ is defined as $\rho(x_n \,|\, x_{<n}) = \rho(x_{1:n})/\rho(x_{<n})$ provided $\rho(x_{<n}) > 0$, with the familiar chain rules $\rho(x_{1:n}) = \prod_{i=1}^{n} \rho(x_i \,|\, x_{<i})$ and $\rho(x_{i:j} \,|\, x_{<i}) = \prod_{k=i}^{j} \rho(x_k \,|\, x_{<k})$ applying as usual.

## 2.2 Geometric Mixing

We now describe geometric mixing, an adaptive online ensemble technique for obtaining a single conditional probability estimate from the output of multiple models. We provide a brief overview here, restricting our attention to the case of binary sequence prediction for simplicity of exposition; for more information see the work of Mattern (2012, 2013, 2016).

### 2.2.1 MODEL

Given $m$ sequential, probabilistic, binary models $\rho_1, \ldots, \rho_m$, Geometric Mixing provides a principled way of combining the $m$ associated conditional probability distributions into a single conditional probability distribution, giving rise to a probability measure on binary sequences that has a number of desirable properties. Let $x_t \in \{0, 1\}$ denote the Boolean target at time $t$. Furthermore, let $p_t = (\rho_1(x_t = 1|x_{<t}), \ldots, \rho_m(x_t = 1|x_{<t}))$. Given a convex set $\mathcal{W} \subset \mathbb{R}^m$ and parameter vector $w \in \mathcal{W}$, the Geometric Mixture is defined by

$$\text{GEO}_w(x_t = 1; p_t) = \frac{\prod_{i=1}^{m} p_{t,i}^{w_i}}{\prod_{i=1}^{m} p_{t,i}^{w_i} + \prod_{i=1}^{m} (1 - p_{t,i})^{w_i}}, \tag{1}$$

with $\text{GEO}_w(x_t = 0; p_t) = 1 - \text{GEO}_w(x_t = 1; p_t)$.

**Properties.** A few properties of this formulation are worth discussing. Setting $w_i = 1/m$ for $i \in [1, m]$ is equivalent to taking the geometric mean of the $m$ input probabilities. As

illustrated in Figure 1, higher absolute values of $w_i$ translate into an increased belief into model $i$'s prediction; for negative values of $w_i$, the prediction needs to be reversed (see blue lines in Fig. 1). If $w = 0$ then $\text{GEO}_w(x_t = 1; p_t) = 1/2$; and in the case where $w_i = 0$ for $i \in \mathcal{S}$ where $\mathcal{S}$ is a proper subset of $[1, m]$, the contributions of the models in $\mathcal{S}$ are essentially ignored (taking $0^0$ to be 1). Due to the product formulation, every model also has "the right of veto", in the sense that a single $p_{t,i}$ close to 0 coupled with a $w_i > 0$ drives $\text{GEO}_w(x_t = 1; p_t)$ close to zero. These properties are graphically depicted in Figure 1.

**Alternate form.** Via simple algebraic manipulation, one can also express Equation 1 as

$$\text{GEO}_w(x_t = 1; p_t) = \sigma\left(w \cdot \text{logit}(p_t)\right), \tag{2}$$

where $\sigma(x) = 1/(1 + e^{-x})$ denotes the sigmoid function, and $\text{logit}(x) = \log(x/(1 - x))$ is its inverse. This form is best suited for numerical implementation. Furthermore, the property of having an input non-linearity that is the inverse of the output non-linearity is the reason why a linear network is obtained when layers of geometric mixers are stacked on top of each other.

**Remarks** This form of combining probabilities is certainly not new; for example it is discussed by Genest and Zidek (1986) under the name *logarithmic opinion pooling* and is closely related to the Product of Experts model of Hinton (2002). While one can attempt to justify this choice of probability combination via appealing to notions such as its *externally Bayesian* properties (Genest and Zidek, 1986), or as the solution to various kinds of weighted redundancy minimization problems (Mattern, 2013, 2016), the strongest case for this form of model mixing is simply its superior empirical performance when $w$ is adapted over time via online gradient descent (Zinkevich, 2003) compared with other similarly complex alternates such as linear mixing (Mattern, 2012; Mahoney, 2013).

### 2.2.2 Properties under the Logarithmic Loss

We assume a standard online learning setting, whereby at each round $t \in \mathbb{N}$ a predictor outputs a binary distribution $q_t : \mathcal{B} \to [0, 1]$, with the environment responding with an observation $x_t \in \mathcal{B}$. The predictor then suffers the logarithmic loss

$$\ell_t(q_t, x_t) = -\log q_t(x_t),$$

before moving onto round $t + 1$. The loss will be close to 0 when the predictor assigns high probability to $x_t$, and large when low probability is assigned to $x_t$; in the extreme cases, a zero loss is obtained when $q_t(x_t) = 1$, and an infinite loss is suffered when $q_t(x_t) = 0$. In the case of geometric mixing, which depends on both the $m$ dimensional input predictions $p_t$ and the parameter vector $w \in \mathcal{W}$, we abbreviate the loss by defining

$$\ell_t^{\text{GEO}}(w) = \ell_t(\text{GEO}_w(\cdot; p_t), x_t). \tag{3}$$

The following proposition, proven in Appendix A, establishes some useful properties about the logarithmic loss when applied to Geometric Mixing.

**Proposition 1** *For all $t \in \mathbb{N}$, $x_t \in \mathcal{B}$, $p_t \in (0, 1)^m$ and $w \in \mathcal{W}$ we have:*

6

1. $\nabla \ell_t^{\text{GEO}}(w) = (\text{GEO}_w(1; p_t) - x_t) \, \text{logit}(p_t)$.

2. $\|\nabla \ell_t^{\text{GEO}}(w)\|_2 \leq \|\text{logit}(p_t)\|_2$.

3. $\ell_t^{\text{GEO}}(w)$ *is a convex function of* $w$.

4. *If* $p_t \in [\varepsilon, 1 - \varepsilon]^m$ *for some* $\varepsilon \in (0, 1/2)$, *then:*

   (a) $\ell_t^{\text{GEO}} : \mathcal{W} \to \mathbb{R}$ *is* $\alpha$-*exp-concave with* $\alpha = \sigma \left( \log \left( \frac{\varepsilon}{1-\varepsilon} \right)^{\max_{w \in \mathcal{W}} \|w\|_1} \right)$;

   (b) $\|\nabla \ell_t^{\text{GEO}}(w)\|_2 \leq \sqrt{m} \log \left( \frac{1}{\varepsilon} \right)$.

Note also that Part 4.(a) of Proposition 1 implies that $\alpha = \Theta \left( (\varepsilon/(1 - \varepsilon))^{\max_{w \in \mathcal{W}} \|w\|_1} \right)$.

**Remarks** The above properties of the sequence of loss functions make it straightforward to apply one of the many different online convex programming techniques to adapt $w$ at the end of each round. The simplest to implement and most computationally efficient is Online Gradient Descent (Zinkevich, 2003) with $\mathcal{W}$ equal to some choice of hypercube. The convexity of $\ell_t^{\text{GEO}}(w)$ allows one to derive an $O(\sqrt{T})$ regret bound using standard techniques with respect to the best $w^* \in \mathcal{W}$ chosen in hindsight provided an appropriate schedule of decaying learning rates is used, where $T$ denotes the total number of rounds. Furthermore, when a fixed learning rate is used one can show $O(s\sqrt{T})$ regret guarantees with respect to data sequences composed of $s$ pieces (Mattern, 2013). More refined algorithms based on coin betting allow the $s$ to be moved inside the square root as shown by Jun et al. (2017). The $\alpha$-exp-concavity also allows second order techniques such as Online Newton Step (Hazan et al., 2007) and its sketched variants (Luo et al., 2016) to be applied. These techniques allow for $O(\log T)$ regret guarantees with respect to the best $w^* \in \mathcal{W}$ chosen in hindsight at the price of a more computationally demanding weight update procedure and further dependence on $\varepsilon$ as given in Part 4.(a) of Proposition 1.

### 2.3 Gated Geometric Mixing

We are now in a position to define our notion of a single neuron, a *gated geometric mixer*, which we obtain by adding a contextual gating procedure to geometric mixing. Here, contextual gating has the intuitive meaning of mapping particular examples to particular sets of weights. More precisely, associated with each neuron is a context function $c : \mathcal{Z} \to \mathcal{C}$, where $\mathcal{Z}$ is the set of possible *side information* and $\mathcal{C} = \{0, \ldots, k - 1\}$ for some $k \in \mathbb{N}$ is the *context space*. Given a convex set $\mathcal{W} \subset \mathbb{R}^d$, each neuron is parametrized by a matrix

$$W = \begin{bmatrix} w_0 \\ \vdots \\ w_{k-1} \end{bmatrix}$$

with each row vector $w_i \in \mathcal{W}$ for $0 \leq i < k$. The context function $c$ is responsible for mapping a given piece of side information $z_t \in \mathcal{Z}$ to a particular row $w_{c(z_t)}$ of $W$, which we then use with standard geometric mixing. More formally, we define the gated geometric mixture prediction as

$$\text{GEO}_W^c(x_t = 1 \, ; \, p_t, z_t) = \text{GEO}_{w_{c(z_t)}}(x_t = 1 \, ; \, p_t), \tag{4}$$

with $\text{GEO}_W^c(x_t = 0\,;\,p_t, z_t) := 1 - \text{GEO}_W^c(x_t = 1\,;\,p_t, z_t)$. Once again we have the following equivalent form

$$\text{GEO}_W^c(x_t = 1; p_t, z_t) = \sigma\left(w_{c(z_t)} \cdot \text{logit}(p_t)\right). \tag{5}$$

The key idea is that our neuron can now specialize its weighting of the input predictions based on some property of the side information $z_t$. The side information can be arbitrary, for example it could be some additional input features, or even functions of $p_t$. Ideally the choice of context function should be informative in the sense that it simplifies the probability combination task.

### 2.3.1 CONTEXT FUNCTIONS

Here we introduce several classes of general purpose context functions that have proven useful empirically, theoretically, or both. All of these context functions take the form of an indicator function $\mathbb{1}_{\mathcal{S}}(z) : \mathcal{Z} \to \mathcal{B}$ on a particular choice of set $\mathcal{S} \subseteq \mathcal{Z}$, with $\mathbb{1}_{\mathcal{S}}(z) := 1$ if $z \in \mathcal{S}$ and 0 otherwise. This list, while sufficient for understanding the content in this paper, is by no means exhaustive; practitioners can and should choose context functions that make sense for the given domain of interest.

**Half-space contexts** This choice of context function is useful for real-valued side information. Given a normal $v \in \mathbb{R}^d$ and offset $b \in \mathbb{R}$, consider the associated affine hyperplane $\{x \in \mathbb{R}^d : x \cdot v = b\}$. This divides $\mathbb{R}^d$ in two, giving rise to two half-spaces, one of which we denote $H_{v,b} = \{x \in \mathbb{R}^d : x \cdot v \geq b\}$. The associated half-space context is then given by $\mathbb{1}_{H_{v,b}}(z)$.
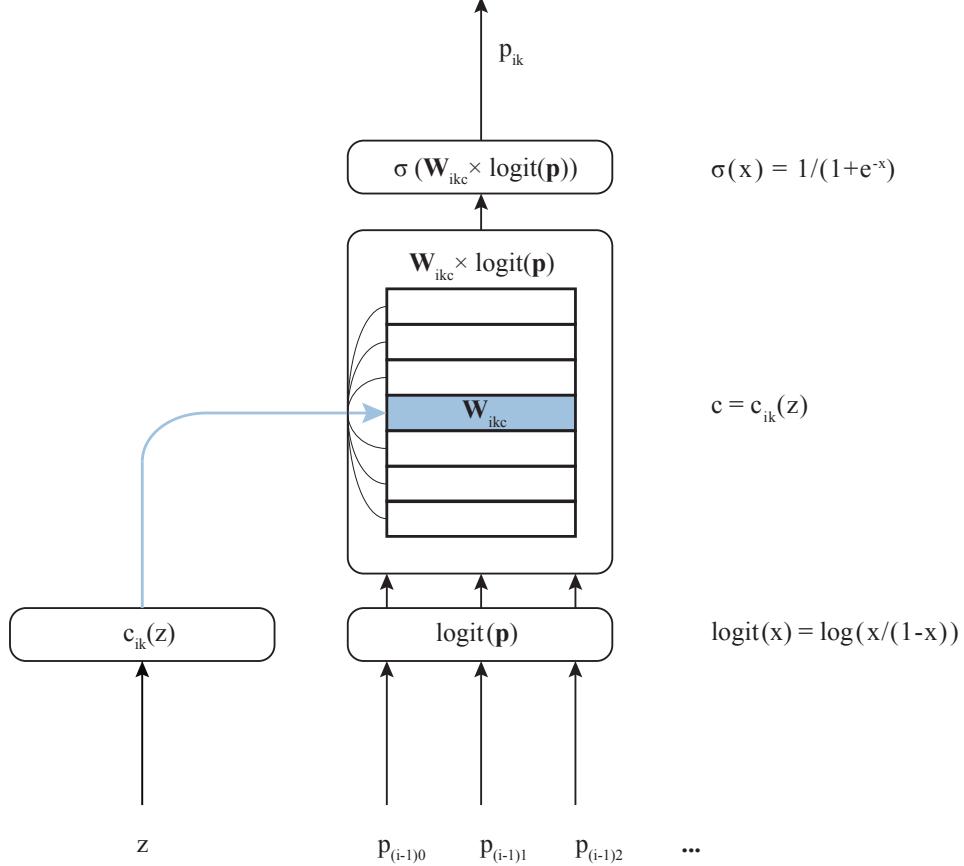
**Skip-gram contexts** The following type of context function is useful when we have multi-dimensional binary side information and can expect single components of $\mathcal{Z}$ to be informative. If $\mathcal{Z} = \mathcal{B}^d$, given an index $i \in [1, d]$, a skip-gram context is given by the function $\mathbb{1}_{\mathcal{S}_i}(z)$ where $\mathcal{S}_i := \{z \in \mathcal{Z} : z_i = 1\}$. One can also naturally extend this notion to categorical multi-dimensional side information or real valued side information by using thresholding.

### 2.3.2 CONTEXT FUNCTION COMPOSITION

Richer notions of context can be created from composition. In particular, any finite set of $d$ context functions $\{c_i : \mathcal{Z} \to \mathcal{C}_i\}_{i=1}^d$ with associated context spaces $\mathcal{C}_1, \ldots, \mathcal{C}_d$ can be composed into a single higher order context function $c : \mathcal{Z} \to \mathcal{C}$, where $\mathcal{C} = \left\{0, 1, \ldots, -1 + \prod_{i=1}^d |\mathcal{C}_i|\right\}$ by defining

$$c(z) = \sum_{i=1}^d c_i(z) \left(\prod_{j=i+1}^d |\mathcal{C}_j|\right).$$

For example, we can combine four different skip-gram contexts into a single context function with a context space containing 16 elements. The combined context function partitions the side information based on the values of the four different binary components of the side information.

$$p_{ik}$$

$$\sigma\,(\mathbf{W}_{ikc}\times \mathrm{logit}(\mathbf{p}))$$

$$\sigma(x) = 1/(1+e^{-x})$$

$$\mathbf{W}_{ikc}\times \mathrm{logit}(\mathbf{p})$$

$$\mathbf{W}_{ikc}$$

$$c = c_{ik}(z)$$

$$c_{ik}(z)$$

$$\mathrm{logit}(\mathbf{p})$$

$$\mathrm{logit}(x) = \log(x/(1-x))$$

$$z \qquad p_{(i-1)0} \qquad p_{(i-1)1} \qquad p_{(i-1)2} \qquad \cdots$$

Figure 2: Behaviour of the $k$th neuron in layer $i > 0$ of a Gated Linear Network.

## 3. Gated Linear Networks

We now introduce *gated linear networks* (GLNs), which are feed-forward networks composed of many layers of gated geometric mixing neurons. Each neuron in a given layer outputs a gated geometric mixture over the predictions from the previous layer, with the final layer consisting of just a single neuron that determines the output of the entire network.

### 3.1 Model

Once again let $\mathcal{Z}$ denote the set of possible side information and $\mathcal{C} \subset \mathbb{N}$ be a finite set called the context space. A GLN is a network of sequential, probabilistic models organized in $L + 1$ layers indexed by $i \in \{0, \ldots, L\}$, with $K_i$ models in each layer. Models are indexed by their position in the network when laid out on a grid; for example, $\rho_{ik}$ will refer to the $k$th model in the $i$th layer. The zeroth layer of the network is called the *base layer* and is constructed from $K_0$ probabilistic base models $\{\rho_{0k}\}_{k=0}^{K_0-1}$ of the form given in Section 2.1. Since each of their predictions is assumed to be a function of the given side information and all previously seen examples, these base models can essentially be arbitrary. The nonzero layers are composed of gated geometric mixing neurons. Associated to each of these will be a fixed context function $c_{ik} : \mathcal{Z} \to \mathcal{C}$ that determines the behavior of the gating. In addition
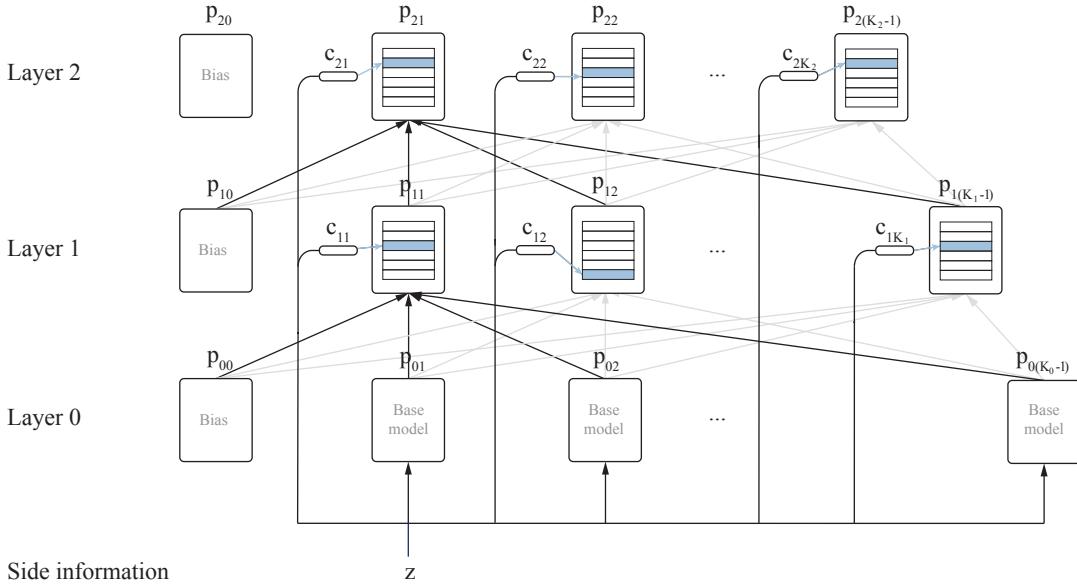
9

Figure 3: A graphical illustration of a Gated Linear Network.

to the context function, for each context $c \in \mathcal{C}$ and each neuron $(i, k)$ there is an associated weight vector $w_{ikc} \in \mathbb{R}^{K_{i-1}}$ which is used to geometrically mix the inputs. We also enforce the constraint of having a non-adaptive bias model on every layer, which will be denoted by $\rho_{i0}$ for each layer $i$. Each of these bias models will correspond to a Bernoulli Process with parameter $\beta$. These bias models play a similar role to the bias inputs in MLPs.

**Network Output** Given a $z \in \mathcal{Z}$, a weight vector for each neuron is determined by evaluating its associated context function. The output of each neuron can now be described inductively in terms of the outputs of the previous layer. To simplify the notation, we assume an implicit dependence on $x_{<t}$ and let $p_{ij}(z) = \rho_{ij}(x_t = 1 \,|\, x_{<t}; z)$ denote the output of the $j$th neuron in the $i$th layer, and $p_i(z) = (p_{i0}(z), p_{i1}(z), \ldots, p_{iK_{i-1}}(z))$ the output of the $i$th layer. The bias output for each layer is defined to be $p_{i0}(z) = \beta$ for all $z \in \mathcal{Z}$, for all $0 \leq i \leq L + 1$, where $\beta \in (0, 1) \setminus \{1/2\}$.[1] From here onwards we adopt the convention of setting $\beta = e/(e+1)$ so that $\mathrm{logit}(\beta) = 1$. For layers $i \geq 1$, the $k$th node in the $i$th layer receives as input the vector of dimension $K_{i-1}$ of predictions of the preceding layer (see Figure 3). The output of a single neuron is the geometric mixture of the inputs with respect to a set of weights that depend on its context, namely

$$p_{ik}(z) = \sigma \left( w_{ikc_{ik}(z)} \cdot \mathrm{logit} \left( p_{i-1}(z) \right) \right),$$

as illustrated by Figure 2. The output of layer $i$ can be re-written in matrix form as

$$p_i(z) = \sigma(W_i(z) \, \mathrm{logit}(p_{i-1}(z))), \tag{6}$$

---

1. The seemingly odd technical constraint $\beta \neq 0.5$ is required to stop the partial derivatives of the loss with respect to the bias weight being 0 under geometric mixing.

where $W_i(z) \in \mathbb{R}^{K_i \times K_{i-1}}$ is the matrix with $k$th row equal to $w_{ik}(z) = w_{ikc_{ik}(z)}$. Iterating Equation 6 once gives

$$p_i(z) = \sigma\left(W_i(z)\operatorname{logit}\left(\sigma\left(W_{i-1}\operatorname{logit}(p_{i-2}(z))\right)\right)\right).$$

Since logit is the inverse of $\sigma$, the $i$th iteration of Equation 6 simplifies to

$$p_i(z) = \sigma\Big(W_i(z)W_{i-1}(z)\ldots W_1(z)\operatorname{logit}(p_0(z))\Big). \tag{7}$$

Equation (7) shows the network behaves like a linear network (Baldi and Hornik, 1989; Saxe et al., 2013), but with weight matrices that are data-dependent. Without the data dependent gating, the product of matrices would collapse to single linear mapping, giving the network no additional modeling power over a single neuron (Minsky and Papert, 1969).

### 3.2 Learning in Gated Linear Networks

We now describe how the weights are learnt in a Gated Linear Network. While architecturally a GLN appears superficially similar to the well-known multilayer perception (MLP), what and how it learns is very different. The key difference is that every neuron in a GLN probabilistically predicts the target. This allows us to associate a loss function to each neuron. This loss function will be defined in terms of just the parameters of the neuron itself; thus, unlike backpropagation, *learning will be local*. Furthermore, this loss function will be convex, which will allow us to avoid many of the difficulties associated with training typical deep architectures. For example, we can get away with simple deterministic weight initializations, which aids the reproducibility of empirical results. The convexity allows us to learn from correlated inputs in an online fashion without suffering significant degradations in performance. And as we shall see later, GLNs are extremely data efficient, and can produce state of the art results in a single pass through the data. One should think of each layer as being responsible for trying to directly improve the predictions of the previous layer, rather than a form of implicit non-linear feature/filter construction as is the case with MLPs trained offline with back-propagation (Rumelhart et al., 1988).

**Weight space and initialization** For our subsequent theoretical analysis, we will require that the weights satisfy the following mild technical constraints:

1. $w_{ijc0} \in [a, b] \subset \mathbb{R}$ for some real $a < 0$ and $b > 0$;

2. $w_{ijc} \in \mathcal{S} \subset \mathbb{R}^{K_{i-1}}$ where $\mathcal{S}$ is a compact, convex set such that $\Delta_{K_{i-1}-1} \subset \mathcal{S}$.

One natural way to simultaneously meet these constraints is to restrict each neurons contextual weight vectors to lie within some (scaled) hypercube: $w_{ijc} \in [-b, b]^{K_{i-1}}$, where $b \geq 1$. For what follows we will need notation to talk about weight vectors at specific times. From here onwards we will use $w_{ijc}^{(t)}$ to denote the weight vector $w_{ijc}$ at time $t$. As each neuron will be solving an online convex programming problem, initialization of the weights is straightforward and is non-essential to the theoretical analysis. Choices found to work well in practice are:

- *zero initialization:* $w_{ikc}^{(1)} = 0$     for all $i, k, c$.

- *geometric average initialization:* $w_{ikc}^{(1)} = 1/K_{i-1}$     for all $i, k, c$.

The zero initialization can be seen as a kind of sparsity prior, where each input model is considered a-priori to be unimportant, which has the effect of making the geometric mixture rapidly adapt to incorporate the predictions of the best performing models. The geometric average initialization forces the geometric mixer to (unsurprisingly) initially behave like a geometric average of its inputs, which makes sense if one believes that the predictions of each input model are reasonable. One could also use small random weights, as is typically done in MLPs, but we recommend against this choice because it makes little practical difference and has a negative impact on reproducibility.

**Weight update**    Learning in GLNs is straightforward in principle; as each neuron probabilistically predicts the target, one can treat the current input to any neuron as a set of expert predictions and apply a single step of local online learning using one of the no-regret methods discussed in Section 2.2.2. The particular choice of optimization method will not in any way affect the capacity results we present later. For practical reasons however we focus our attention on Online Gradient Descent (Zinkevich, 2003) with $\mathcal{W}_{ik}$ a hypercube. This allows the weight update for any neuron at layer $i$ to be done in time complexity $O(K_{i-1})$, which permits the construction of large networks.

More precisely, let $\ell_t^{ij}(w_{ijc})$ denote the loss of the $j$th neuron in layer $i$. Using Equation (3) we have

$$\ell_t^{ij}(w_{ijc}) = \ell_t(\text{GEO}_w(\,\cdot\,; p_{i-1}(z_t), x_t). \tag{8}$$

Now, for all $i \in [1, L]$, $j \in K_i$, and for all $c = c_{ij}(z_t)$, we set

$$w_{ijc}^{(t+1)} = \Pi_i \left( w_{ijc}^{(t)} - \eta_t \nabla \ell_t^{ij}(w_{ijc}^{(t)}) \right), \tag{9}$$

where $\Pi_i$ is the projection operation onto hypercube $[-b, b]^{K_{i-1}}$:

$$\Pi_i(x) = \operatorname*{arg\,min}_{y \in [-b,b]^{K_{i-1}}} \|y - x\|_2 \ .$$

The projection is efficiently implemented by clipping every component of $w_{ijc}^{(t)}$ to the interval $[-b, b]$. The learning rate $\eta_t \in \mathbb{R}_+$ can depend on time; we will revisit the topic of appropriately setting the learning rate in subsequent sections.

**Performance guarantees for individual neurons**    Let $(i, k)$ refer to a non-bias neuron. We now state a performance guarantee for a single neuron when $\mathcal{W}_{ik} = [-b, b]^{K_{i-1}}$. Given $a \in \mathcal{C}$ consider the index set $N_{ika}(n) = \{t \leq n : c_{ik}(z_t) = a\}$, the set of rounds when context $a$ is observed by neuron $(i, k)$. The regret experienced by the neuron on $N_{ika}(n)$ is defined as the difference between the losses suffered by the neuron and the losses it would have suffered using the best choice of weights in hindsight.

$$R_{ika}(n) = \sum_{t \in N_{ika}(n)} \ell_t^{ik}(w_{ika}^{(t)}) - \min_{w \in \mathcal{W}_{ik}} \sum_{t \in N_{ika}(n)} \ell_t^{ik}(w) \,. \tag{10}$$

Zinkevich (2003) showed that if the learning rates are set to $\eta_t = \frac{D}{G\sqrt{t}}$, then the regret of gradient descent is at most

$$R_{ika}(n) \leq \frac{3DG}{2} \sqrt{|N_{ika}(n)|} \,, \tag{11}$$

12

where $D = \max_{x,y \in \mathcal{W}} \|x - y\|_2$ is the diameter of $\mathcal{W}_{ik}$ and $G$ is an upper bound on the 2-norm of the gradient of the loss. In our case we have $D = 2b\sqrt{K_{i-1}}$, with a bound $G \leq \sqrt{K_{i-1}} \log(\frac{1}{\varepsilon})$ following from Proposition 1 Part 4.(b) by ensuring the input to each neuron is within $[\varepsilon, 1 - \varepsilon]$. The regret definition given in Eq. (10) is on the subsequence of rounds when $c_{ik}(z_t) = a$. The total regret for neuron $(i, k)$ is

$$R_{ik}(n) = \sum_{a \in \mathcal{C}} R_{ika}(n) \,. \tag{12}$$

Combining the above and simplifying, we see that

$$
\begin{aligned}
R_{ik}(n) &\leq \frac{3DG}{2} \sum_{a \in \mathcal{C}} \sqrt{|N_{ika}(n)|} = \frac{3DG|\mathcal{C}|}{2} \sum_{a \in \mathcal{C}} \frac{1}{|\mathcal{C}|} \sqrt{|N_{ika}(n)|} \\
&\leq \frac{3DG|\mathcal{C}|}{2} \sqrt{\sum_{a \in \mathcal{C}} \frac{1}{|\mathcal{C}|} |N_{ika}(n)|} = \frac{3DG}{2} \sqrt{|\mathcal{C}| \sum_{a \in \mathcal{C}} |N_{ika}(n)|} \,,
\end{aligned}
$$

where the first and second inequalities follow from Eq. (11) and Jensen's inequality respectively. Noting that $\bigcup_{a \in \mathcal{C}} N_{ika}(n) = \{1, \ldots, n\}$ allows us to further simplify the bound to

$$R_{ik}(n) \leq \frac{3DG}{2} \sqrt{|\mathcal{C}|n} \,.$$

In the particular case where each input probability is bounded between $[\varepsilon, 1 - \varepsilon]$ and the weights reside in the hypercube $[-b, b]^{K_{i-1}}$, the above bound becomes

$$R_{ik}(n) \leq 3bK_{i-1} \sqrt{|\mathcal{C}|n} \, \log\left(\frac{1}{\varepsilon}\right) . \tag{13}$$

The cumulative loss suffered by the omniscient algorithm is:

$$\sum_{a \in \mathcal{C}} \min_{w \in \mathcal{W}_{ik}} \sum_{t \in N_{ika}(n)} \ell_t^{ik}(w) \,.$$

The regret is the difference between the losses suffered by the learning procedure and the above quantity. Since the latter usually grows like $\Omega(n)$, the $O(\sqrt{n})$ additional penalty suffered by the algorithm is asymptotically negligible.

How should we interpret Equation 13? First of all, note the linear dependence on the number of inputs; this is the price a single neuron pays in the worst case for increasing the width of the preceding layer. There is also a linear dependence on $b$, which is expected as when the competitor set gets larger it should be more difficult to compete with the optimal solution in hindsight. Clipping inputs is common implementation trick for many deep learning systems; here this has the effect of stopping the 2-norm of the gradient exploding, which would adversely affect the worst case regret. Furthermore, by assuming the unit basis vector $e_j \in \mathcal{W}_{ik}$ we guarantee that the cumulative loss of neuron $(i, k)$ is never much larger than the cumulative loss of neuron $(i-1, j)$; to see this, note that if the weights were set to $e_j$ then geometric mixing simply replicates the output of the $j$th neuron in the preceding

layer. This is the bare minimum we should expect, but the guarantee actually provides much more than this by ensuring that the cumulative loss of neuron $(i, k)$ is never much worse than the best geometric mixture of its inputs. As we will observe in Section 4, if the context functions are sufficiently rich, then the layering of geometric mixtures combined with Eq. (11) leads to a network that learns to approximate arbitrary functions.

### 3.3 Computational properties

Some computational properties of Gated Linear Networks are now discussed.

**Complexity of a single online learning step**   Generating a prediction requires computing the contexts from the given side information for each neuron, and then performing $L$ matrix-vector products. Under the assumption that multiplying a $m \times n$ by $n \times 1$ pair of matrices takes $O(mn)$ work, the total time complexity to generate a single prediction is $O\left(\sum_{i=1}^{L} K_i K_{i-1}\right)$ for the matrix-vector products, which in typical cases will dominate the overall runtime. Using online gradient descent just requires updating the rows of the weight matrices using Equation 9; this again takes time $O\left(\sum_{i=1}^{L} K_i K_{i-1}\right)$.

**Parallelism**   When generating a prediction, parallelism can occur within a layer, similar to an MLP. The local training rule however enables all the neurons to be updated simultaneously, as they have no need to communicate information to each other. This compares favorably to back-propagation and significantly simplifies any possible distributed implementation. Furthermore, as the bulk of the computation is primarily matrix multiplication, large speedups can be obtained straightforwardly using GPUs.

**Static prediction with pre-trained models**   In the case where no online updating is desired, prediction can be implemented potentially more efficiently depending on the exact shape of the network architecture. Here one should implement Equation 7 by first solving a Matrix Chain Ordering problem (Hu and Shing, 1982) to determine the optimal way to group the matrix multiplications.

## 4. Effective Capacity of Gated Linear Networks

Neural networks have long been known to be capable of approximating arbitrary continuous functions with almost any reasonable activation function (Hornik, 1991, and others). We will show that provided the contexts are chosen sufficiently richly, then GLNs also have the capacity to approximate large classes of functions. More than this, the capacity is *effective* in the sense that gradient descent will eventually find the approximation. In contrast, similar results for neural networks show the existence of a choice of weights for which the neural network will approximate some function, but do not show that gradient descent (or any other single algorithm) will converge to these weights. Of course we do not claim that gated linear networks are the only model with an effective capacity result. For example, $k$-nearest-neighbour with appropriately chosen $k$ is also universal for large function classes. We view universality as a good first step towards a more practical theoretical understanding of a model, but the ultimate goal is to provide meaningful finite-time theoretical guarantees for *interesting* and *useful* function classes. Gated linear networks have some advantages

over other architectures in the sense that they are constructed from small pieces that are well-understood in isolation and the nature of the training rule eases the analysis relative to neural networks.

Our main theorem is the following abstract result that connects the richness of the context functions used in each layer to the ability of a network with no base predictors to make constant improvements. After the theorem statement we provide intuition on simple synthetic examples and analyse the richness of certain classes of context functions.

**Theorem 1 (Capacity)** *Let $\mu$ be a measure on $(\mathbb{R}^d, \mathcal{F})$ with $\mathcal{F}$ the Lebesgue $\sigma$-algebra and let $z_t \in \mathbb{R}^d$ be a sequence of independent random variables sampled from $\mu$. Furthermore, let $x_1, x_2, \ldots$ be a sequence of independent Bernoulli random variables with $\mathbb{P}(x_t = 1|z_t) = f(z_t)$ for some $\mathcal{F}$-measurable function $f : \mathbb{R}^d \to [0, 1]$. Consider a gated linear network and for each layer $i$ let $\mathcal{G}_i = \{c_{ik} : 1 \le k \le K_i\}$ be the set of context functions in that layer. Assume there exists a $\delta > 0$ such that for each non-bias neuron $(i, k)$ the weight-space $\mathcal{W}_{ik}$ is compact and $(-\delta, \delta) \times \Delta^{K_{i-1}-2} \subseteq \mathcal{W}_{ik}$. Provided the weights are learned using a no-regret algorithm, then the following hold with probability one:*

1. *For each non-bias neuron $(i, k)$ there exists a non-random function $p_{ik}^* : \mathbb{R}^d \to (0, 1)$ such that*

$$\lim_{n \to \infty} \frac{1}{n} \sum_{t=1}^{n} \sup_{z \in \text{Supp}(\mu)} \left| p_{ik}(z; w^{(t)}) - p_{ik}^*(z) \right| = 0.$$

2. *The average of $p_{ik}^*$ converges to the average of $f$ on the partitions of $\mathbb{R}^d$ induced by the contexts:*

$$\sum_{i=1}^{\infty} \max_k \max_{c \in \mathcal{G}_{i+1}} \sum_{a \in \mathcal{C}} \left( \int_{c^{-1}(a)} (f(z) - p_{ik}^*(z)) \, d\mu(z) \right)^2 \le \max \left\{ 2, \frac{4}{\delta} \right\} \log \left( \frac{1}{\beta} \right),$$

   *where $\beta = e/(1+e)$ is the output of bias neurons $p_{i0}(z)$.*

3. *There exists a non-random $\mathcal{F}$-measurable function $p_\infty^* : \mathbb{R}^d \to (0, 1)$ such that*

$$\lim_{i \to \infty} \max_k \int_{\mathbb{R}^d} (p_{ik}^*(z) - p_\infty^*(z))^2 d\mu(z) = 0.$$

The assumption that the sequence of side information is independent and identically distributed may be relaxed significantly. A sufficient criteria is that $z_1, z_2, \ldots$ follows a Markov processes that mixes to stationary measure $\mu$. Unfortunately this requires a benign assumption on the stability of the learning algorithm, which is satisfied by most online learning algorithms including online gradient descent with appropriately decreasing learning rates.

**Definition 2** *A learning procedure for GLNs is stable if there exist constants $r < 0$ and $C > 0$ such that for all data sequences $(x_t)_t$ and $(z_t)_t$ it holds that:*

$$\left\| w_{ika}^{(t)} - w_{ika}^{(t+1)} \right\|_\infty \le \mathbb{1}_{c_{ik}^{-1}(a)}(z_t) C t^r \qquad \text{for all } t.$$

The proofs of Theorems 1 and 3 are given in Appendices B and C respectively.

**Theorem 3** *Suppose the learning procedure is stable and that all conditions of Theorem 1 hold except that $(z_t)_t$ is sampled from an aperiodic and $\phi$-irreducible Markov chain with stationary distribution $\mu$. Then (1), (2) and (3) from Theorem 1 hold.*

Theorem 3 suggests that GLNs could be well suited reinforcement learning applications. For example, one could combine a GLN-based density model with the policy evaluation approach of Veness et al. (2015).

**Interpretation of Theorems 1 and 3**  The theorem has three parts. The first says the output of each neuron converges almost surely in Cesaro average to some nonrandom function on the support of $\mu$. Outside the support of $\mu$ the algorithm will never observe data, so provided there is no covariate shift we should not be concerned about this. The second part gives a partial characterisation of what this function is. Let

$$\varepsilon_i = \max_k \max_{c \in \mathcal{G}_{i+1}} \left( \int_{c^{-1}(a)} (f(z) - p_{ik}^*(z)) d\mu(z) \right)^2.$$

The theorem says that $\sum_{i=1}^\infty \varepsilon_i = O(1)$, which since $\varepsilon_i$ is clearly positive implies that $\lim_{i \to \infty} \varepsilon_i \to 0$. The convergence is fast in the sense that $\varepsilon_i$ is approximately subharmonic. Suppose that $\mathcal{G}_i = \mathcal{G}$ is the same for all layers and $\varepsilon_i$ is small, then for all $c \in \mathcal{G}$ and $a \in \mathcal{C}$ we have

$$\int_{c^{-1}(a)} f(z) d\mu(z) \approx \int_{c^{-1}(a)} p_{ik}^*(z) d\mu(z),$$

which means that $p_{ik}^*$ is approximately equal to $f$ on average on all sets $c^{-1}(a) \subset \mathbb{R}^d$. We will shortly see that if $\{c^{-1}(a) : c \in \mathcal{G}, a \in \mathcal{C}\}$ is sufficiently rich, then $f(z) \approx p_{ik}^*(z)$ also holds. Finally, the last part of the theorem shows that as the number of layers tends to infinity, the output of neurons in later layers converges to some single function. Note that all results are about what happens in the limit of infinite data. In principle one could apply the regret guarantees from gradient descent to calculate a rate of convergence, a task which we leave for the future.

**Intuition on synthetic data**  We now illustrate the behaviour of a small network with various contexts in the simple case of $d = 1$. Figure 5 shows the chosen architecture and the problem setup. Each box relates to a non-bias neuron of the network, comparing the ground truth, $f(z)$ (black line) to the output of each neuron, $p_{ik}(x|z)$ (coloured line). Axes are identical for all boxes, and labelled on bottom left: The single-dimensional side information $z \in [-3, 3]$ is used to derive the Bernoulli parameters according to a Gaussian transformation $f(z) = e^{-z^2/2}$. Thus, the training data consist of samples $z_t$ drawn uniformly over the range $[-3, 3]$, and their labels drawn from Bernoulli distributions parametrised by $f(z_t)$.

The GLN has six layers, with 3/2/2/2/2/1 non-bias neurons in each layer reading from bottom to top (Fig. 5), and a single base predictor $p_{00}(z) = \alpha$ for all $z$ (not shown). All bias neurons are also set to $p_{0k} = \beta$ (not shown). The half-space contexts of the non-bias

neurons are parametrized by an offset, chosen randomly (coloured dashed lines in Fig. 5). The bottom row depicts the output of the network in the first layer after training. Each neuron learns to predict the average of $f$ on the partitions of $[-3, 3]$ induced by $c_{1k}^{-1}(a)$ for each $a \in \{0, 1\}$. Already in the second layer, the neurons' outputs gain in complexity, combining information not only from the neuron's own contexts, but also from the estimates of the neurons below. In effect, every output has 4 discontinuities (albeit two of them are very close, due to two contexts in the first layer having similar offsets). The top panel shows the output of the non-bias neuron in the sixth layer, which even for this small network is predicting with reasonable accuracy.

The performance of the network is easily improved by increasing the number of neurons in the first layer, as shown in Figure 4. If $K_1 = 100$, then the output of the non-bias neuron in the second layer is a near-perfect predictor. As it happens, a two-layer network with half-space contexts and sufficiently many neurons in the first layer can learn to approximate any continuous function, which explains the performance of such a shallow network on this problem. Note that the base predictor is completely uninformative in this case and the half-space contexts are not tuned to the data.

**Applications of Theorem 1**  We now see several applications of Theorem 1 for different contexts. If we assume that $\mathcal{G}_i = \mathcal{G}$ for all $i$, then the result shows that in the limit of infinite data a sufficiently deep network predicts like the average of $f$ on $c^{-1}(a)$ for all context functions $c \in \mathcal{G}$ and all $a \in \mathcal{C}$. When the space of context functions is sufficiently rich this implies that the network indeed predicts like $f$. Let $\mathcal{G}$ be a (possibly infinite) set of context functions such that

$$\|h\|_{\mathcal{G}} = \sup_{c \in \mathcal{G}} \sum_{a \in \mathcal{C}} \left| \int_{c^{-1}(a)} h(z) d\mu(z) \right|$$

is a norm on the space of bounded measurable functions. Note that homogeneity, non-negativity and the triangle inequality are trivial, which leaves one only to check that $\|h\|_{\mathcal{G}} = 0$ implies that $h = 0$ $\mu$-almost-everywhere. Now suppose that $\mathcal{G}_1 \subset \mathcal{G}_2 \subset \mathcal{G}_3 \subset \cdots$ is a sequence of finite subsets of $\mathcal{G}$ such that for any bounded measurable $h$

$$\lim_{i \to \infty} \|h\|_{\mathcal{G}_i} = \|h\|_{\mathcal{G}} \ , \tag{14}$$
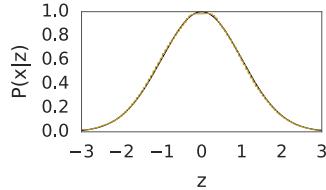


Figure 4: Output of two-layer network with evenly spaced half-space contexts and $K_1 = 100$
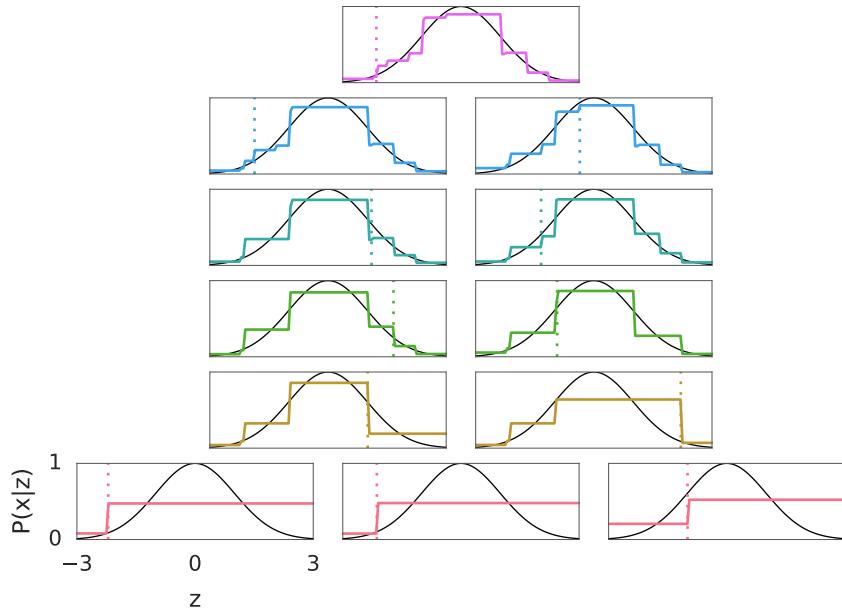
Figure 5: Output of a six-layer network with half-space contexts. Each box represents a non-bias neuron in the network, the function to fit is shown in black, and the output distribution learnt by each neuron is shown in colour (for example, red for the first layer and purple for the top-most neuron). All axes are identical, as labeled on bottom left. The dashed coloured lines represent the choice of hyperplane for each neuron.

which is true in the special case that $\mathcal{G}$ is countable and $\mathcal{G}_i \uparrow \mathcal{G}$. Then by Parts (2) and (3) of Theorem 1 we have

$$0 = \lim_{i \to \infty} \max_{c \in \mathcal{G}_i} \sum_{a \in \mathcal{C}} \left| \int_{c^{-1}(a)} (f(z) - p_\infty^*) d\mu(z) \right| = \lim_{i \to \infty} \|f - p_\infty^*\|_{\mathcal{G}_i} = \|f - p_\infty^*\|_{\mathcal{G}} \ .$$

Therefore $\mu(z : f(z) - p_\infty(z) = 0) = 1$, which means that almost surely the network asymptotically predicts like $f$ $\mu$-almost-everywhere. We now give some example countable classes of contexts, which by Eq. (14) can be used to build a universal network. All the claims in the following enumeration are proven in Appendix D.

- (*Covers*) Let $B_r(x) = \{y : \|x - y\|_2 < r\}$ be the open set of radius $r > 0$. If $\mathcal{G} = \{\mathbb{1}_{B_r(x)} : x \in \mathbb{Q}^d, r \in \mathbb{Q} \cap (0, \infty)\}$, then $\|\cdot\|_{\mathcal{G}}$ is a norm on the space of bounded measurable functions.

- (*Covers cont.*) The above result can easily be generalised to the situation where $\mathcal{G}$ is the set of indicator functions on any countable base of $\mathbb{R}^d$.

- (*Hyperplanes*) Now assume that $\mu$ is absolutely continuous with respect to the Lebesgue measure. If $\mathcal{G}$ is the space of context functions that are indicators on half-spaces $\mathcal{G} = \{\mathbb{1}_{H_{v,c}} : v \in \mathbb{Q}^d, \|v\| = 1, c \in \mathbb{Q}\}$, then $\|\cdot\|_{\mathcal{G}}$ is a norm on the spaces of bounded $\mathcal{F}$-measurable functions.
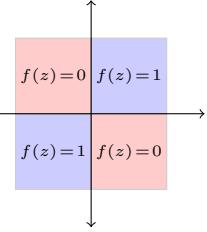
18

We conclude this section with three remarks.

**Remark 4 (Asymptotic non-diversity)** *Part (3) of Theorem 1 shows that asymptotically in the depth of the network there is vanishingly little diversity in the predictors in the bottom layer. The intuition is quite simple. If the contexts are sufficient that the network learn the true function, then eventually all neurons will predict perfectly (and hence be the same). On the other hand, if the class of contexts is not rich enough for perfect learning, then there comes a time when all neurons in a given layer cannot do better than copying the output of their best parent. Once this happens, all further neurons make the same predictions.*

**Remark 5 (Universality of shallow half-spaces)** *If $\mu$ is supported on a compact set and $f$ is continuous, then a two-layer network is sufficient to approximate $f$ to arbitrary precision with half-space contexts. A discussion of this curiosity is deferred to Appendix E. In practice we observe that except when $d = 1$ it is beneficial to use more layers. And note that arbitrary precision is only possible by using arbitrarily large weights and wide networks.*

**Remark 6 (Effective capacity $\neq$ capacity)** *Theorem 1 shows that if the context functions are sufficiently rich and the network is deep, then the network can learn any well-behaved function using online gradient descent. This is an effective capacity result because we have provided a model (the network) and a learning procedure (OGD) that work together to learn any well-behaved function. For a fixed architecture, the capacity of a GLN can be much larger than the effective capacity. In particular, there exist functions that can be modelled by a particular choice of weights that online gradient descent will never find. (This does not contradict Theorem 1; one can always construct a larger GLN whose effective capacity is sufficiently rich to model such functions.)*

*To construct an example demonstrating this failure we exploit the fact that training each neuron against the target means that neurons do not coordinate to solve the following 'exclusive-or' problem. Let $d = 2$ and $\mu$ be uniform on $[-1, 1]^2$ and $f : \mathbb{R}^2 \to [0, 1]$ be given by $f(z) = \mathbb{1}\{z_1 z_2 \geq 0\}$. Next suppose $\mathcal{G} = \{c_1, c_2\}$ where $c_i(z) = \mathbb{1}\{z_i \geq 0\}$. Notice that $z_1$ or $z_2$ alone is insufficient to predict $f(z)$. If a gated linear network with $K_i = 2$ for all $i \geq 2$ and contexts from $\mathcal{G}$ is trained in the iid setting of Theorem 1, then it is easy to check that $\lim_{i \to \infty} p_{ik}^*(z) = 1/2$ almost surely. And yet there exists a choice of weights such that $p_{ik}(z; w) = f(z)$ for all $z$. The problem is the optimal weights have neurons in the same layer coordinating so that a single neuron in the next layer has access to $z_1$ and $z_2$ as inputs. But online gradient descent has neurons selfishly optimizing their own performance, which does not lead to coordination in this case.*

## 5. Adaptive Regularization via Sub-network Switching

While GLNs can have almost arbitrary capacity in principle, large networks are susceptible to a form of the *catch-up phenomenon* (van Erven et al., 2007). During the initial stages of learning neurons in the lower layers typically have better predictive performance than neurons in the higher layers. This section presents a solution to this problem based on
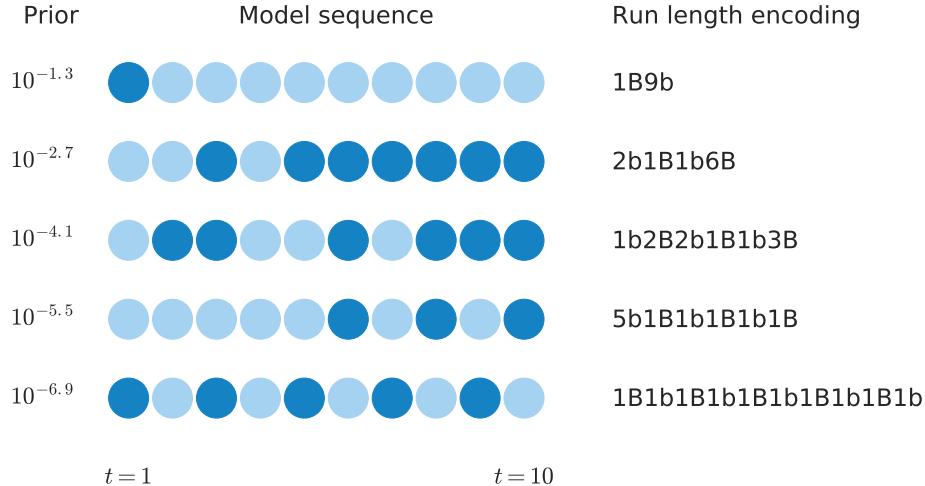
Figure 6: A visual depiction of the run-length encoding prior.

*switching* (Veness et al., 2012), a fixed share (Herbster and Warmuth, 1998) variant tailored to the logarithmic loss. The main idea is simple: as each neuron predicts the target, one can construct a switching ensemble across all neurons predictions. This guarantees that the predictions made by the ensemble are not much worse than the predictions made by the best sparsely changing sequence of neurons. We now describe this process in detail.

### 5.1 Model

Let $\mathcal{M} = \{\rho_{ij} : i \in [1, L], j \in [0, K_i - 1]\}$ denote the model class consisting of all neurons that make up a particular GLN with $L$ layers and $K_i$ neurons in each layer. Now for all $n \in \mathbb{N}$, for all $x_{1:n} \in \mathcal{X}^n$, consider a Bayesian (non-parametric) mixture that puts a prior $w_\tau(\cdot)$ over all *sequences* of neurons inside the index set $\mathcal{I}_n(\mathcal{M}) = \mathcal{M}^n$, namely

$$\tau(x_{1:n}) = \sum_{\nu_{1:n} \in \mathcal{I}_n(\mathcal{M})} w_\tau(\nu_{1:n}) \, \nu_{1:n}(x_{1:n}) \tag{15}$$

where $\nu_{1:n}(x_{1:n}) = \prod_{k=1}^n \nu_k(x_k \mid x_{<k})$. As $w_\tau(\cdot)$ is a prior, it needs to be non-negative and satisfy $\sum_{\nu_{1:n}} w_\tau(\nu_{1:n}) = 1$ for all $n \in \mathbb{N}$. From the dominance property of Bayesian mixtures it immediately follows that for any $\nu_{1:n}^* \in \mathcal{I}_n(\mathcal{M})$ we have

$$-\log \tau(x_{1:n}) \leq -\log \left( w_\tau \left( \nu_{1:n}^* \right) \, \nu_{1:n}^* (x_{1:n}) \right) \leq -\log w_\tau(\nu_{1:n}^*) - \log \nu_{1:n}^*(x_{1:n}). \tag{16}$$

Thus the regret $\mathcal{R}_n = -\log \left( \frac{\tau(x_{1:n})}{\nu_{1:n}^*(x_{1:n})} \right)$ with respect to a sequence of models $\nu_{1:n}^*$ is upper bounded by $-\log w_\tau(\nu_{1:n}^*)$. Putting a uniform prior over all neuron sequences would lead to a vacuous regret of $n \log(|\mathcal{M}|)$, so we are forced to concentrate our prior mass on a smaller set of neuron sequences we think a-priori are likely to predict well.

Empirically we found that when the number of training examples is small neurons in the lower layers usually predict better than those in higher layers, but this reverses as more data becomes available. Viewing the sequence of best-predicting neurons over time as a

string, we see that a run-length encoding gives rise to a highly compressed representation with a length linear in the number of times the best-predicting neuron changes. Run-length encoding can be implemented probabilistically by using an arithmetic encoder with the following recursively defined prior:

$$
w_\tau(\nu_{1:n}) = \begin{cases} 1 & \text{if} \quad n = 0 \\ \frac{1}{|\mathcal{M}|} & \text{if} \quad n = 1 \\ w_\tau(\nu_{<n}) \times \left( \frac{n-1}{n} \mathbb{1}[\nu_n = \nu_{n-1}] + \frac{1}{n(|\mathcal{M}|-1)} \mathbb{1}[\nu_n \neq \nu_{n-1}] \right) & \text{otherwise} . \end{cases} \tag{17}
$$

A graphical depiction of this prior is shown in Figure 6. Multiple sequences of models are depicted by rows of light blue and dark blue circles, with each colour depicting a different choice of model. A run length encoding of the sequence is given by a string made up of pairs of symbols, the first being an integer representing the length of the run and the second a character ($B$ for dark blue, $b$ for light blue) depicting the choice of model. One can see from the figure that high prior weight is assigned to model sequences which have shorter run length encodings. More formally, one can show (see Appendix F) that

$$
- \log w_\tau(\nu_{1:n}) \leq (s(\nu_{1:n}) + 1)(\log |\mathcal{M}| + \log n), \tag{18}
$$

for all $\nu_{1:n} \in \mathcal{I}_n(\mathcal{M})$, where $s(\nu_{1:n}) := \sum_{t=2}^n \mathbb{I}[\nu_t \neq \nu_{t-1}]$ denotes the number of switches from one neuron to another in $\nu_{1:n}$. Combining Equations 16 and 18 allows us to upper bound the regret $\mathcal{R}_n$ with respect to an arbitrary $\nu_{1:n}^* \in \mathcal{I}_n(\mathcal{M})$ by

$$
(s(\nu_{1:n}^*) + 1)(\log |\mathcal{M}| + \log n) .
$$

Therefore, when there exists a sequence of neurons with a small number $s(\nu_{1:n}^*) \ll n$ of switches that performs well, only logarithmic regret is suffered, and one can expect the switching ensemble to predict almost as well as if we knew what the best performing sparsely changing sequence of neurons was in advance. Note too that the bound holds in the case where we just predict using a single choice of neuron; here only $\log |\mathcal{M}| + \log n$ is suffered, so there is essentially no downside to using this method if such architecture adaptivity is not needed.

## 5.2 Algorithm

A direct computation of Equation 15 would require $|\mathcal{M}|^n$ additions, which is clearly intractable. The switching algorithm of Veness et al. (2012) runs in time $O(n|\mathcal{M}|)$, but was originally presented in terms of log-marginal probabilities, which requires a log-space implementation to avoid numerical difficulties. An equivalent numerically robust formulation is described here, which incrementally maintains a weight vector that is used to compute a convex combination of model predictions at each time step.

Let $u_{ik}^{(t)} \in (0, 1]$ denote the switching weight associated with the neuron $(i, k)$ at time $t$. The weights will satisfy the invariant $\sum_{i=1}^{|L|} \sum_{k=0}^{K_i-1} u_{ik}^{(t)} = 1$ for all $t$. At each time step $t$, the switching mixture outputs the conditional probability

$$
\tau(x_t|x_{<t}) = \sum_{i=1}^{|L|} \sum_{k=0}^{K_i-1} u_{ik}^{(t)} \rho_{ik}(x_t|x_{<t}),
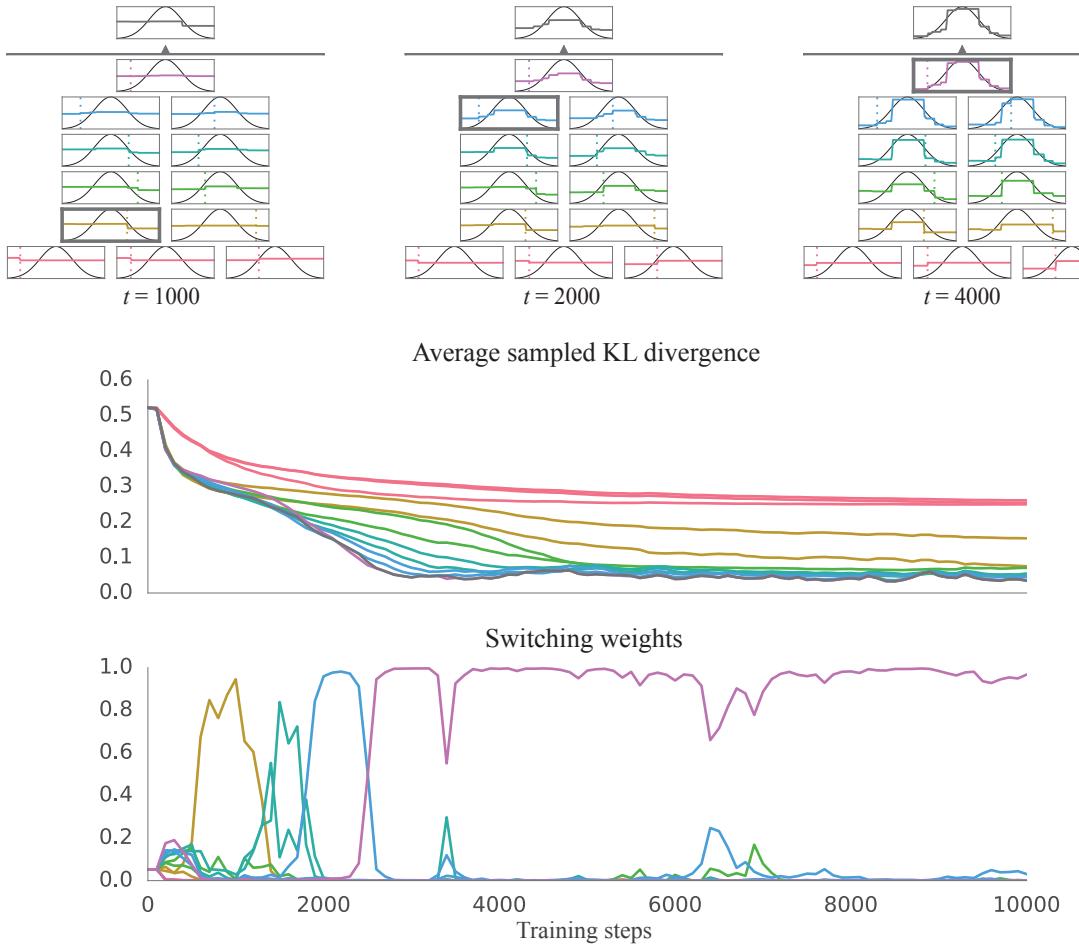$$

Figure 7: Evolution of model fit over time

with the weights defined, for all $1 \leq i \leq L$ and $0 \leq k < K_i$, by $u_{ik}^{(1)} = 1/|\mathcal{M}|$ and

$$u_{ik}^{(t+1)} = \frac{1}{(t+1)(|\mathcal{M}|-1)} + \left( \frac{t|\mathcal{M}| - t - 1}{(t+1)(|\mathcal{M}|-1)} \right) \frac{u_{ik}^{(t)} \, \rho_{ik}(x_t|x_{<t})}{\tau(x_t|x_{<t})}.$$

This weight update can be straightforwardly implemented in $O(|\mathcal{M}|)$ time per step. To avoid numerical issues, we recommend enforcing the weights to sum to 1 by explicitly dividing by $\sum_{i=1}^{|L|} \sum_{k=0}^{K_i-1} u_{ik}^{(t+1)}$ after each weight update.

## 5.3 Illustration

Figure 7 shows an example of a small 6 layer GLN fitting a family of Bernoulli distributions parametrized by a Gaussian transformation $g(z) = e^{-z^2/2}$ of a single dimensional side information value $z \in [-3, 3]$, which is the same setup depicted in Fig. 5 without switching. The GLN has 6 layers, with $3/2/2/2/2/1$ neurons on each layer reading from bottom to top. The output distribution learnt by each neuron is shown in colour (for example, red
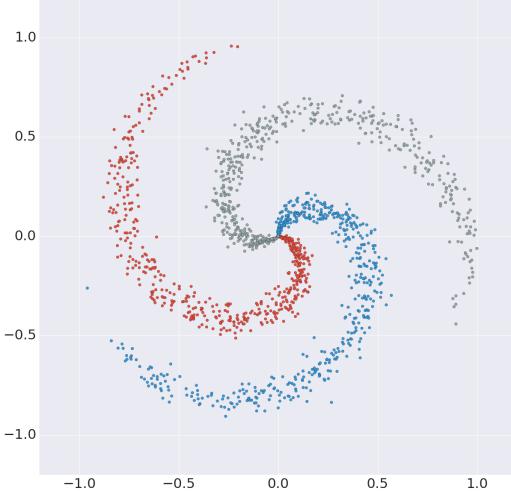
22

Figure 8: The spiral dataset. The horizontal and vertical axes describe the values of the two dimensional side information, with the class label being indicated by either a red, grey or blue point.

for the bottom layer and purple for the top neuron). The dashed vertical line in each distribution denotes the choice of half-space which controls the gating of input examples. Each example was generated by first sampling a $z$ uniformly distribution between $[-3, 3]$ and then sampling a label from a Bernoulli distribution parameterized by $g(z)$. The topmost box in each of the four subfigures shows the output distribution (in black) obtained using subnetwork switching. The top-left subfigure shows the model fit at $t = 1000$; the top-middle at $t = 2000$; and the top-right the fit at $t = 4000$. One can clearly see the model fit improving over time, and that with sufficient training the higher level neurons better model the target family of distributions. The bottom graph in the figure shows the evolution of the switching weights over time; here we see that initially a neuron on the 3rd layer dominates the predictions, then one in the 4th layer, then the 5th and subsequently settling down on the predictions made by the top-most neuron.

## 6. Empirical Results

We now present a short series of case studies to shed some insight into the capabilities of GLNs. When describing our particular choice of architecture, we adopt the convention of describing the number of neurons on each layer by a dashed list of natural numbers going from the most shallow layer to the deepest layer.

### 6.1 Non-linear decision boundaries with half-spaces

Our first experiment is on a synthetic 'spiral' ternary classification task, which we use to demonstrate the ability of GLNs to model non-linear decision boundaries. A visualization of the dataset is shown in Figure 8. The horizontal and vertical axes describe the values of the two dimensional side information, with the class label being indicated by either a red, grey or blue point. To address this task with a GLN, we constructed an ensemble of 3 GLNs to form a standard one-vs-all classifier. Each member of the ensemble was a 3 layer network
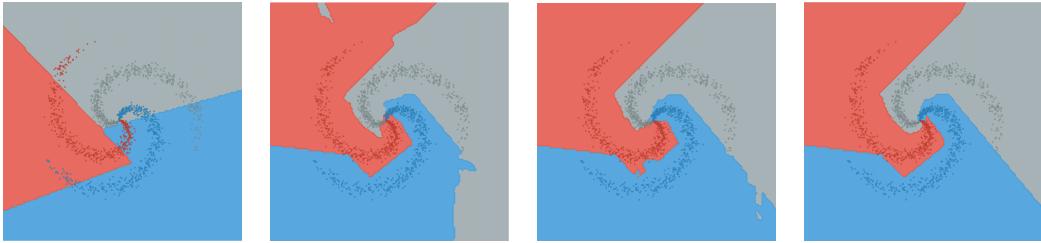
Figure 9: From left to right, the decision boundary determined by a neuron on the bottom layer, the middle layer, and the final layer, and the top-level switching mixture. The fit of the bottommost neuron is quite poor. The middle neuron improves significantly, making mistakes only at the center of the spiral. The topmost neuron has learnt to fit the data almost perfectly. The final decision boundary of the switching ensemble resembles a smoothed version of the deepest neuron.

consisting of 50-25-1 neurons which each used a half-space context for gating. The half-space contexts for each neuron were determined by sampling a 2 dimensional normal vector whose components were distributed according to $\mathcal{N}(0, 36)$, and a bias weight distributed according to $\mathcal{N}(0, 9)$. Each component of all weight vectors were constrained to lie within $[-200, 200]$. The side information was the 2-dimensional x/y values, and the input to the network was the component-wise sigmoid of these x/y values (as GLNs require the input to be within [0,1]). The learning rate was set to 0.01. Figure 9 shows a plot of the resultant decision boundaries for neurons of differing depths after training on a single pass through all the examples.

## 6.2 Online Half-space Classification: MNIST

Next we explore the use of GLNs for classification on the well known MNIST dataset (Lecun et al., 1998). Once again we used an ensemble of 10 GLNs to construct a one-vs-all classifier. Each member of the ensemble was a 3 layer network consisting of 1500-1500-1 neurons, each of which used 6 half-space context functions (using the context composition technnique of Section 2.3.2) as the gating procedure, meaning that each neuron contained 64 different possible weight vectors. The half-space contexts for each neuron were determined by sampling a 784 dimensional normal vector whose components were distributed according to $\mathcal{N}(0, 0.01)$, and a bias weight of 0. The learning rate for an example at time $t$ was set to $\min\{8000/t, 0.3\}$. These parameters were determined by a grid search across the training data. Each component of all weight vectors were constrained to lie within $[-200, 200]$. The input features were pre-processed by first applying mean-subtraction and a de-skewing operation (Ghosh and Wan, 2017).

Running the method purely online across a single pass of the data gives an accuracy on the test set of 98.3%. If weight updating during the testing set is disabled, the test set accuracy drops slightly to 98.1%. Without the de-skewing operation, the accuracy drops to 96.9%. It is worth nothing that our classifier contains no image specific domain knowledge; in fact, one could even apply a (fixed across all images) permutation to each image input and get exactly the same result.

### 6.3 Online Density Modeling: MNIST

Our final result is to use GLNs and image specific gating to construct an online image density model for the binarized MNIST dataset (Larochelle and Murray, 2011a), a standard benchmark for image density modeling. By exploiting the chain rule

$$\mathbb{P}(X_1, X_2, \ldots, X_d) = \prod_{i=1}^{d} \mathbb{P}(X_i \mid X_{<i})$$

of probability, we constructed a so-called autoregressive density model over the $28 \times 28$ dimensional binary space by using 784 GLNs to model the conditional distribution for each pixel; a row-major ordering was used to linearize the two dimensional pixel locations.

**Base Layer**  The base layer for each GLN was determined by taking the outputs of a number of skip-gram models (Guthrie et al., 2006). The context functions used by these skip-gram models included a number of specific geometric patterns that are known to be helpful for lossless image compression (Witten et al., 1999), plus a large number of randomly sampled pixel locations. Some example skip-grams are shown in Figure 10; the red box
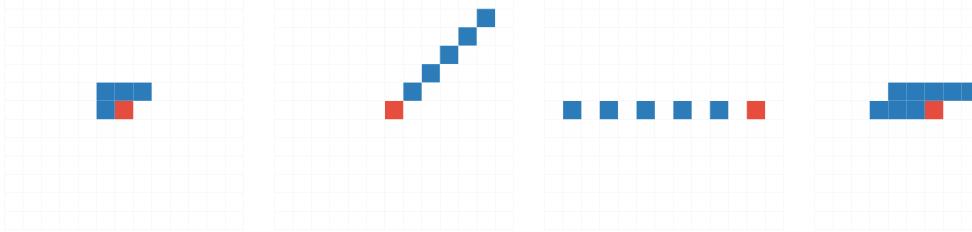


Figure 10: Example skip-gram context functions.

indicates the location of the pixel in the image, and the blue boxes indicate the pixels which are used to determine the skip-gram context. For each pixel model, the base layer would consist of up to 600 skip-gram predictions; the exact number of predictions depends on the location of the pixel in the linear ordering; the set of permitted skip-gram models for pixel $i$ are those whose context functions do not depend on pixels greater than or equal $i$ within the linear ordering. The Zero-Redundancy estimator (Willems et al., 1997) was used to estimate the per-context skip-gram probability estimates online.

**Image specific context functions**  Empirically we found that we could significantly boost the performance of our method by introducing two simple kinds of image specific context functions.

The first were *max-pool* context functions, some examples of which are depicted in Figure 11. Once again the red square denotes the location of the current pixel. For each shaded region of a particular colour, the maximum of all binary pixel values was computed. The max-pool context function returns the number represented by the binary representation obtained by concatenating these max-pooled values together (in some fixed order).

The next type of context we introduced were *distance* context functions, some examples of which are depicted in Figure 12. Once again the red square denotes the location of the current pixel. Each non-white/non-red location is labeled with an index between 1 and
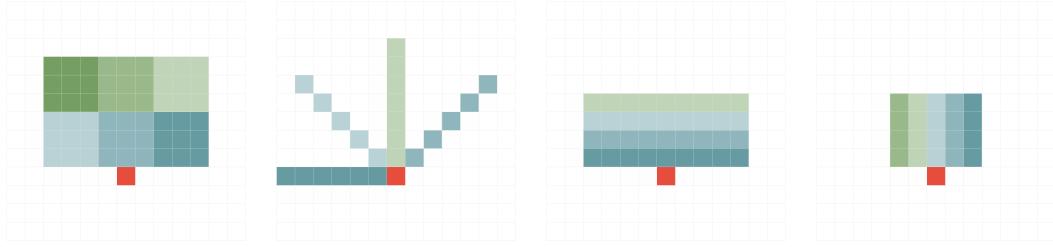
Figure 11: Example max-pool context functions.

the number of non-white/non-red locations; darker colours indicate smaller index values. The distance context function returns the smallest index value whose pixel value is equal to 1. This class of context allows one to measure distance to the nearest active pixel under various orderings of local pixel locations.
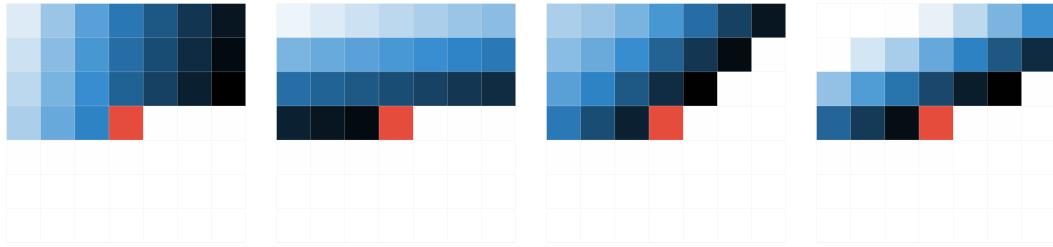


Figure 12: Example distance context functions.

**Network construction**    Each GLN used 4-layers, with the gating function for each neuron determined by a choice of skip-gram/max-pool/distance context function. A set of 200 different context functions were defined, and these were randomly distributed across a network architecture whose shape was 35-60-35-70. The learning rate for an example seen at time $t$ was set to $\min\{25/t, 0.005\}$.

**Results**    Running the method purely online across a single pass of the data (we concatenated the training, validation and test sets into one contiguous stream of data) gave an average loss of 79.0 nats per image across the test data. To the best of our knowledge, this result matches the state of the art (Van Den Oord et al., 2016) of any batch trained density model which outputs exact probabilities, and is close to the estimated upper bounds reported for state of the art variational approaches (Bachman, 2016).

## 7. Relationship to the PAQ family compressors

One of the main contributions of this paper is to justify and explain the empirical success of the PAQ family of compression programs. At the time of writing, the best performing PAQ implementation in terms of general purpose compression performance is CMIX, an open-source project whose details are described by Knoll (2017) on his personal website. CMIX uses a mixing network comprised of many gated geometric mixers, each of which use different skip-gram context functions to define the gating operations. While many of the core building blocks have been analyzed previously by Mattern (2016), the reason

for the empirical success of such locally trained mixing networks had until now remained somewhat of a mystery. Our work shows that such architectures are a special cases of Gated Linear Networks, and that future improvements may result by exploring alternate no-regret learning methods, by using a broader class of context functions or by using our gated linear network formulation to enable efficient implementation on GPUs via vectorized matrix operations. In principle the universality of GLNs suggests that the performance of PAQ-like approaches will continue to scale as hardware improves. Furthermore, our MNIST density modeling results suggest that such methods can be competitive with state of the art deep learning approaches on domains beyond general purpose file compression.

## 8. Conclusion

We have introduced Gated Linear Networks, a family of architectures designed for online learning under the logarithmic loss. Under appropriate conditions, such architectures are universal in the sense that they can can model any bounded Borel-measurable function; more significantly, they are guaranteed to find this solution using any appropriate choice of no-regret online optimization algorithm to locally adapt the weights of each neuron. Initial experimental results suggest that the method warrants further investigation.

## Acknowledgments

## References

Philip Bachman. An architecture for deep, hierarchical generative models. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 4826–4834. Curran Associates, Inc., 2016.

P. Baldi and K. Hornik. Neural networks and principal component analysis: Learning from examples without local minima. *Neural Networks*, 2(1):53–58, January 1989. ISSN 0893-6080. doi: 10.1016/0893-6080(89)90014-2.

David Balduzzi. Deep online convex optimization with gated games. *CoRR*, abs/1604.01952, 2016. URL http://arxiv.org/abs/1604.01952.

T. C. Bell, J. G. Cleary, and I. H. Witten. *Text Compression*. Prentice Hall, Englewood Cliffs, NJ, 1990.

Tim Bell and Ross Arnold. A corpus for the evaluation of lossless compression algorithms. *Data Compression Conference*, 00:201, 1997. ISSN 1068-0314. doi: doi. ieeecomputersociety.org/10.1109/DCC.1997.582019.

Nicolo Cesa-Bianchi and Gabor Lugosi. *Prediction, Learning, and Games*. Cambridge University Press, New York, NY, USA, 2006. ISBN 0521841089.

Xi Chen, Diederik P. Kingma, Tim Salimans, Yan Duan, Prafulla Dhariwal, John Schulman, Ilya Sutskever, and Pieter Abbeel. Variational lossy autoencoder. *CoRR*, abs/1611.02731, 2016.

G. V. Cormack and R. N. S. Horspool. Data compression using dynamic markov modelling. *The Computer Journal*, 30(6), 1987.

Jakob N. Foerster, Justin Gilmer, Jascha Sohl-Dickstein, Jan Chorowski, and David Sussillo. Input switched affine networks: An RNN architecture designed for interpretability. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 1136–1145, International Convention Centre, Sydney, Australia, 06–11 Aug 2017. PMLR.

Christian Genest and James V. Zidek. Combining probability distributions: A critique and an annotated bibliography. *Statistical Science*, 1(1):114–135, 1986.

Dibya Ghosh and Alvin Wan, 2017. URL `https://fsix.github.io/mnist/Deskewing.html`.

Karol Gregor, Ivo Danihelka, Alex Graves, Danilo Rezende, and Daan Wierstra. Draw: A recurrent neural network for image generation. In David Blei and Francis Bach, editors, *Proceedings of the 32nd International Conference on Machine Learning (ICML-15)*, pages 1462–1471. JMLR Workshop and Conference Proceedings, 2015.

Ishaan Gulrajani, Kundan Kumar, Faruk Ahmed, Adrien Ali Taiga, Francesco Visin, David Vázquez, and Aaron C. Courville. Pixelvae: A latent variable model for natural images. *CoRR*, abs/1611.05013, 2016.

David Guthrie, Ben Allison, W. Liu, Louise Guthrie, and Yorick Wilks. A closer look at skip-gram modelling. In *Proceedings of the Fifth international Conference on Language Resources and Evaluation (LREC-2006)*, Genoa, Italy, 2006.

Elad Hazan. Introduction to online convex optimization. *Foundations and Trends in Optimization*, 2(3-4):157–325, 2016. doi: 10.1561/2400000013.

Elad Hazan, Amit Agarwal, and Satyen Kale. Logarithmic regret algorithms for online convex optimization. *Machine Learning*, 69(2-3):169–192, 2007.

Sigurdur Helgason. The radon transform on r n. In *Integral Geometry and Radon Transforms*, pages 1–62. Springer, 2011.

Mark Herbster and Manfred K. Warmuth. Tracking the best expert. *Machine Learning*, 32 (2):151–178, August 1998. ISSN 0885-6125.

Geoffrey E. Hinton. Training products of experts by minimizing contrastive divergence. *Neural Computation*, 14(8):1771–1800, August 2002. ISSN 0899-7667.

Kurt Hornik. Approximation capabilities of multilayer feedforward networks. *Neural networks*, 4(2):251–257, 1991.

T. C. Hu and M. T. Shing. Computation of matrix chain products. part i. *SIAM Journal on Computing*, 11(2):362–373, 1982. doi: 10.1137/0211028.

Marcus Hutter. Hutter prize, 2017. URL `http://prize.hutter1.net/`.

Robert A. Jacobs, Michael I. Jordan, Steven J. Nowlan, and Geoffrey E. Hinton. Adaptive mixtures of local experts. *Neural Comput.*, 3(1):79–87, 1991. ISSN 0899-7667. doi: 10.1162/neco.1991.3.1.79.

Kwang-Sung Jun, Francesco Orabona, Stephen Wright, and Rebecca Willett. Improved strongly adaptive online learning using coin betting. In *Artificial Intelligence and Statistics*, pages 943–951, 2017.

B. Knoll and N. de Freitas. A machine learning perspective on predictive coding with paq8. In *Data Compression Conference (DCC)*, pages 377–386, April 2012.

Byron Knoll, 2017. URL `http://www.byronknoll.com/cmix.html`.

Hugo Larochelle and Iain Murray. The neural autoregressive distribution estimator. In Geoffrey Gordon, David Dunson, and Miroslav Dudk, editors, *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15 of *Proceedings of Machine Learning Research*, pages 29–37, Fort Lauderdale, FL, USA, 11–13 Apr 2011a. PMLR.

Hugo Larochelle and Iain Murray. The neural autoregressive distribution estimator. *Journal of Machine Learning Research (JMLR)*, 15:29–37, 2011b.

Yann Lecun, Lon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. In *Proceedings of the IEEE*, pages 2278–2324, 1998.

Haipeng Luo, Alekh Agarwal, Nicolò Cesa-Bianchi, and John Langford. Efficient second order online learning by sketching. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 902–910, 2016.

Matthew Mahoney. Fast text compression with neural networks. *AAAI*, 2000.

Matthew Mahoney. Adaptive weighing of context models for lossless data compression. *Technical Report, Florida Institute of Technology CS*, 2005.

Matthew Mahoney. *Data Compression Explained*. 2013.

Christopher Mattern. Mixing strategies in data compression. In *2012 Data Compression Conference, Snowbird, UT, USA, April 10-12*, pages 337–346, 2012.

Christopher Mattern. Linear and geometric mixtures - analysis. In *2013 Data Compression Conference, DCC 2013, Snowbird, UT, USA, March 20-22, 2013*, pages 301–310, 2013.

Christopher Mattern. *On Statistical Data Compression*. PhD thesis, Technische Universität Ilmenau, Germany, 2016.

Sean P Meyn and Richard L Tweedie. *Markov chains and stochastic stability*. Springer Science & Business Media, 2012.

Marvin Minsky and Seymour Papert. *Perceptrons: An Introduction to Computational Geometry*. MIT Press, Cambridge, MA, USA, 1969.

Georg Ostrovski, Marc G. Bellemare, Aäron van den Oord, and Rémi Munos. Count-based exploration with neural density models. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, pages 2721–2730, 2017.

David E. Rumelhart, Geoffrey E. Hinton, and Ronald J. Williams. Neurocomputing: Foundations of research. chapter Learning Representations by Back-propagating Errors, pages 696–699. MIT Press, Cambridge, MA, USA, 1988. ISBN 0-262-01097-6.

Andrew M. Saxe, James L. McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *CoRR*, abs/1312.6120, 2013.

J. Schmidhuber and S. Heil. Sequential neural text compression. *IEEE Transactions on Neural Networks*, 7(1):142–146, Jan 1996. ISSN 1045-9227. doi: 10.1109/72.478398.

Aäron Van Den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48*, ICML'16, pages 1747–1756. JMLR.org, 2016.

Tim van Erven, Peter Grunwald, and Steven de Rooij. Catching up faster in bayesian model selection and model averaging. In *Advances in Neural Information Processing Systems 20, Proceedings of the Twenty-First Annual Conference on Neural Information Processing Systems, Vancouver, British Columbia, Canada, December 3-6, 2007*, pages 417–424, 2007.

Joel Veness, Kee Siong Ng, Marcus Hutter, and Michael H. Bowling. Context tree switching. In *2012 Data Compression Conference, Snowbird, UT, USA, April 10-12, 2012*, pages 327–336, 2012.

Joel Veness, Marc G. Bellemare, Marcus Hutter, Alvin Chua, and Guillaume Desjardins. Compress and control. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence, January 25-30, 2015, Austin, Texas, USA.*, pages 3016–3023, 2015.

Frans Willems, Yuri Shtarkov, and Tjalling Tjalkens. Reflections on "the context-tree weighting method: Basic properties". In *IEEE Information Theory Society Newsletter*, volume 47, 1997.

Ian H. Witten, Radford M. Neal, and John G. Cleary. Arithmetic coding for data compression. *Communications of the ACM*, 30(6):520–540, June 1987. ISSN 0001-0782. doi: 10.1145/214762.214771.

Ian H. Witten, Alistair Moffat, and Timothy C. Bell. *Managing Gigabytes (2nd Ed.): Compressing and Indexing Documents and Images*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1999. ISBN 1-55860-570-3.

Martin Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *Machine Learning, Proceedings of the Twentieth International Conference (ICML 2003), August 21-24, 2003, Washington, DC, USA*, pages 928–936, 2003.

## Appendix A. Proof of Proposition 1

Part (1): If $x_t = 1$ then $\ell_t^{\text{GEO}}(w) = -\log \sigma(\text{logit}(p_t) \cdot w)$ from Equation 2, hence using $\sigma'(x) = \sigma(x)[1 - \sigma(x)]$ and Equation 2 once more gives

$$\frac{\partial \ell_t^{\text{GEO}}}{\partial w_i} = [\sigma(\text{logit}(p_t) \cdot w) - 1] \, \text{logit}(p_{t,i}) = (\text{GEO}_w(1 \; ; \; p_t) - x_t) \, \text{logit}(p_{t,i}). \qquad (19)$$

Similarly, if $x_t = 0$ then $\ell_t^{\text{GEO}}(w) = -\log(1 - \sigma(\text{logit}(p_t) \cdot w))$, and so

$$\frac{\partial \ell_t^{\text{GEO}}}{\partial w_i} = \sigma(\text{logit}(p_t) \cdot w) \, \text{logit}(p_{t,i}) = (\text{GEO}_w(1 \; ; \; p_t) - x_t) \, \text{logit}(p_{t,i}). \qquad (20)$$

Hence $\nabla \ell_t^{\text{GEO}}(w) = (\text{GEO}_w(1 \; ; \; p_t) - x_t) \, \text{logit}(p_t)$.

Part (2): As $|(\text{GEO}_w(1 \; ; \; p_t) - x_t)| \leq 1$ it holds that $\|\nabla \ell_t^{\text{GEO}}(w)\|_2 \leq \|\text{logit}(p)\|_2$ for all $w \in \mathcal{C}$.

Part (3): the convexity of $\ell_t^{\text{GEO}}$ has been established already by Mattern (2013).

Part (4a): We make use of Lemma 4.1 in Hazan (2016), which states that a twice differentiable function $f : \mathbb{R}^m \to \mathbb{R}$ is $\alpha$-exp-concave at $x \in \mathbb{R}^m$ if and only if there exists a scalar $\alpha > 0$ such that $\nabla^2 f(x) - \alpha \nabla f(x) \nabla f(x)^\top$ is positive semi-definite.

We can calculate the Hessian of the loss directly from Equations 19 and 20 by observing that

$$\frac{\partial^2 \ell_t^{\text{GEO}}}{\partial w_j \, \partial w_i} = \text{logit}(p_{t,i}) \, \text{logit}(p_{t,j}) \, \sigma(\text{logit}(p_t) \cdot w)[1 - \sigma(\text{logit}(p_t) \cdot w)],$$

and so

$$\nabla^2 \ell_t^{\text{GEO}}(w) = \sigma(\text{logit}(p_t) \cdot w)[1 - \sigma(\text{logit}(p_t) \cdot w)] \, \text{logit}(p_t) \, \text{logit}(p_t)^\top \qquad (21)$$
$$= \text{GEO}_w(1 \; ; \; p_t)[1 - \text{GEO}_w(1 \; ; \; p_t)] \, \text{logit}(p_t) \, \text{logit}(p_t)^\top.$$

Furthermore, from Part (1), we have

$$\nabla \ell_t^{\text{GEO}}(w) \, \nabla \ell_t^{\text{GEO}}(w)^\top = (\text{GEO}_w(1 \; ; \; p_t) - x_t)^2 \, \text{logit}(p_t) \, \text{logit}(p_t)^\top.$$

Therefore, letting $q = \text{GEO}_w(1 \; ; \; p_t)$, we have that

$$A := \nabla^2 \ell_t^{\text{GEO}}(w) - \alpha \nabla \ell_t^{\text{GEO}}(w) \, \nabla \ell_t^{\text{GEO}}(w)^\top$$
$$= \text{logit}(p_t) \, \text{logit}(p_t)^\top [q(1 - q) - \alpha(q - x_t)^2]$$

Now if we could show $k := q(1 - q) - \alpha(q - x_t)^2 \geq 0$, then we would have that $A$ is positive-semidefinite as $x^\top A x = k \, x^\top \text{logit}(p_t) \, \text{logit}(p_t)^\top x = k \, (x \cdot \text{logit}(p_t))^2 \geq 0$ for any non-zero real vector $x$. Considering the two cases $x_t = 0$ and $x_t = 1$ separately, for $k \geq 0$ to hold, we

must have either $\alpha \leq (1-q)/q$ and $\alpha \leq q/(1-q)$; these conditions can be met independently of $x_t$ by choosing $\alpha = \varepsilon_0/(1-\varepsilon_0)$, where $\varepsilon_0$ is the smallest possible value of $\text{GEO}_w(1 \; ; \; p_t)$ for any $p_t$ or $w$. Given the assumption that $p_t \in [\varepsilon, 1-\varepsilon]^d$ for some $\varepsilon \in (0, 1/2)$, we conclude that $\ell_t^{\text{GEO}}(w)$ is $\alpha$-exp-concave with

$$\alpha = \min_{w \in \mathcal{W}} \sigma(w \cdot \text{logit}(\varepsilon)) = \min_{w \in \mathcal{W}} \sigma\left(\prod_{i=1}^{m} \log\left(\frac{\varepsilon}{1-\varepsilon}\right)^{w_i}\right) = \sigma\left(\log\left(\frac{\varepsilon}{1-\varepsilon}\right)^{\max_{w \in \mathcal{W}} \|w\|_1}\right).$$

Also, notice that Eq. (21) is clearly positive-semidefinite, confirming Part (3).

Part (4b): Applying Part (2), and using the assumption that $p_t \in [\varepsilon, 1-\varepsilon]^m$ we have

$$\|\nabla \ell(w)\|_2 \leq \|\text{logit}(p)\|_2 = \sqrt{\sum_{i=1}^{m} \text{logit}^2(p_{t,i})} \leq \sqrt{m \log^2\left(\frac{1-\varepsilon}{\varepsilon}\right)}$$

$$= \sqrt{m} \left(\log(1-\varepsilon) - \log(\varepsilon)\right) \leq \sqrt{m} \log\left(\frac{1}{\varepsilon}\right).$$

## Appendix B. Proof of Theorem 1

The proof of Theorem 1 relies on several simple technical lemmas. Recall for $p, q \in (0, 1)$ that $\mathcal{D}(p, q) = p \log(p/q) + (1-p) \log((1-p)/(1-q))$ is the relative entropy (or Kullback-Leibler divergence) between Bernoulli distributions with biases $p$ and $q$ respectively.

**Lemma 7** Let $p \in [0, 1]$ and $g_p(y) = \mathcal{D}(p, \sigma(y))$, then:

1. $g_p'(\text{logit}(q)) = q - p$.

2. $g_p''(y) = \exp(y)/(1 + \exp(y))^2 \in (0, 1/4]$.

3. $g_p(\Delta + \text{logit}(q)) \leq \mathcal{D}(p, q) + \Delta(q - p) + \Delta^2/8$.

**Proof** The first and second parts are trivial. The third follows from the fact that $g_p$ is everywhere twice differentiable for all $p \in [0, 1]$, which implies that

$$g_p(x + \Delta) \leq g_p(x) + \Delta g_p'(x) + \max_y g_p''(y) \frac{\Delta^2}{2} \leq g_p(x) + \Delta g_p'(x) + \frac{\Delta^2}{8}.$$

Then choose $x = \text{logit}(q)$ and note that $g_p(\text{logit}(q)) = \mathcal{D}(p, q)$. ∎

**Lemma 8** Given $f, p, q \in (0, 1)$,

$$\mathcal{D}(f, q) - \mathcal{D}(f, p) = \mathcal{D}(p, q) + \frac{\partial}{\partial \alpha} \mathcal{D}(f, \sigma((1-\alpha) \text{logit}(p) + \alpha \text{logit}(q)))\bigg|_{\alpha=0}$$

**Proof** Note that

$$
\frac{\partial}{\partial \alpha} \mathcal{D}(f, \sigma((1-\alpha)\operatorname{logit}(p) + \alpha \operatorname{logit}(q))) \Big|_{\alpha=0} = (f - p)(\operatorname{logit}(q) - \operatorname{logit}(p)) .
$$

The proof follows from the definition of $\mathcal{D}(\cdot, \cdot)$. ∎

**Proof** [of Theorem 1]

**Part (1)** Recall that for all layers $i$ (including $i = 0$) the output of the bias neuron in the $i$th layer is $p_{i0}^*(z) = \beta$. Given non-bias neuron $(i, k)$ and context $a \in \mathcal{C}$ let

$$
N_{ika}(n) = \{t \le n : c_{ik}(z_t) = a\}
$$

be the set of rounds when the context for neuron $(i, k)$ is $a$. Now define

$$
w_{ika}^* = \underset{w \in \mathcal{W}_{ik}}{\arg\min} \int_{c_{ik}^{-1}(a)} \mathcal{D}(f(z), \sigma(w \cdot \operatorname{logit}(p_{i-1}^*(z)))) d\mu(z)
$$
$$
p_{ik}^*(z) = \sigma(w_{ika}^* \cdot \operatorname{logit}(p_{i-1}^*(z))) .
$$

Note that $w_{ika}^*$ need not be unique, but the strict convexity of $\mathcal{D}(f(z), \sigma(\cdot))$ for all $z$ ensures that $p_{ik}^*$ does not depend on this choice. We prove by induction that the following holds almost surely for all non-bias neurons $(i, k)$

$$
\lim_{n \to \infty} \frac{1}{n} \sum_{t=1}^{n} \sup_{z \in \operatorname{Supp}(\mu)} \left| p_{ik}(z; w^{(t)}) - p_{ik}^*(z) \right| = 0 . \tag{22}
$$

Let $(i, k)$ be a non-bias neuron and $c = c_{ik}$ be its context function and assume the above holds for all preceding layers. Fix $a \in \mathcal{C}$ and abbreviate $N(n) = N_{ika}(n)$ and $q_t = \operatorname{logit}(p_{i-1}(z_t; w^{(t)}))$ be the input to the $i$th layer in the $t$th round $q^*(z) = \operatorname{logit}(p_{i-1}^*(z))$ and

$$
\ell_t(w) = x_t \log\left(\frac{1}{\sigma(w \cdot q_t))}\right) + (1 - x_t) \log\left(\frac{1}{1 - \sigma(w \cdot q_t)}\right) .
$$
$$
\ell_t^*(w) = x_t \log\left(\frac{1}{\sigma(w \cdot q^*(z_t))}\right) + (1 - x_t) \log\left(\frac{1}{1 - \sigma(w \cdot q^*(z_t))}\right) .
$$

Let $\ell^*(w) = \mathbb{E}[\ell_t^*(w) \mathbb{1}_{N(n)}(t)]$, which does not depend on $t$ by the assumption that $(z_t, x_t)$ is independent and identically distributed. Furthermore,

$$
\ell^*(w) = \int_{c^{-1}(a)} \mathcal{D}(f(z), \sigma(w \cdot q^*(z))) d\mu(z) + C , \tag{23}
$$

where $C$ is a constant depending only on $\mu$ and $f$ and given by:

$$
C = -\int_{c^{-1}(a)} (f(z) \log f(z) + (1 - f(z)) \log(1 - f(z))) d\mu(z) .
$$

33

Since $\mathcal{W}_{ik}$ is bounded it follows that $q_t$ is bounded. Therefore $\ell_t(u)$ is bounded and continuous in $u$ and $q_t$, which by the induction hypothesis means that

$$\lim_{n\to\infty} \frac{1}{n} \sum_{t\in N(n)} \sup_{u\in\mathcal{W}_{ik}} |\ell_t(u) - \ell_t^*(u)| = 0\,.$$

Since the weights are learned using an algorithm with sublinear regret it follows that

$$0 \geq \lim_{n\to\infty} \min_{u\in\mathcal{W}_{ik}} \frac{1}{n} \sum_{t\in N(n)} \left(\ell_t(w^{(t)}) - \ell_t(u)\right) = \lim_{n\to\infty} \min_{u\in\mathcal{W}_{ik}} \frac{1}{n} \sum_{t\in N(n)} \left(\ell_t^*(w^{(t)}) - \ell_t^*(u)\right)\,.$$

For each $u \in \mathcal{W}_{ik}$ define martingale

$$M_n(u) = \sum_{t\in N(n)} \left(\ell_t^*(w^{(t)}) - \ell_t^*(u) - (\ell^*(w^{(t)}) - \ell^*(u))\right)\,. \tag{24}$$

By the martingale law of large numbers $\mathbb{P}\left(\lim_{n\to\infty} M_n(u)/n = 0\right) = 1$. Since $\ell_t^*(u)$ is uniformly Lipschitz in $u$ for all $t$, a union bound on a dense subset of $\mathcal{W}_{ik}$ is enough to ensure that with probability one it holds for all $u \in \mathcal{W}_{ik}$ that

$$0 \geq \lim_{n\to\infty} \frac{1}{n} \sum_{t\in N(n)} \ell_t^*(w^{(t)}) - \ell_t^*(u) = \lim_{n\to\infty} \frac{1}{n} \sum_{t\in N(n)} \left(\ell^*(w^{(t)}) - \ell^*(u)\right)\,.$$

The conclusion follows from the definition of $\ell^*$ in Eq. (23).

**Part (2)**  We begin by noting that for any non-bias neuron $(i,k)$ it holds that

$$\begin{aligned}
D_{ik} &= \int_{\mathbb{R}^d} \mathcal{D}(f(z), p_{ik}^*(z))d\mu(z) \\
&= \sum_{a\in\mathcal{C}} \int_{c_{ik}^{-1}(a)} \mathcal{D}(f(z), p_{ik}^*(z))d\mu(z) \\
&= \sum_{a\in\mathcal{C}} \min_{w\in\mathcal{W}_{ik}} \int_{c_{ik}^{-1}(a)} \mathcal{D}(f(z), \sigma(w\cdot\mathrm{logit}(p_{i-1}^*(z))))d\mu(z) \\
&\leq \min_{w\in\mathcal{W}_{ik}} \int_{\mathbb{R}^d} \mathcal{D}(f(z), \sigma(w\cdot\mathrm{logit}(p_{i-1}^*(z))))d\mu(z) \\
&\leq \min_j \int_{\mathbb{R}^d} \mathcal{D}(f(z), p_{i-1,j}^*(z))d\mu(z) = D_{i-1}^*\,,
\end{aligned}$$

where the first and last equalities serve as definitions and in the last line we made use of the assumption that $e_j \in \mathcal{W}_{ik}$ for each neuron $(i-1,j)$. Therefore the expected loss of any neuron in the $i$th layer is smaller than that of the best neuron in the $(i-1)$th layer. Since $p_{i0}^*(z) = \beta$ for all layers $i$ and $z \in \mathbb{R}^d$ we have

$$\int_{\mathbb{R}^d} \mathcal{D}(f(z), p_{i0}^*(z))d\mu(z) = \int_{\mathbb{R}^d} \mathcal{D}(f(z), \beta)d\mu(z) \leq \mathcal{D}(0, \beta) = -\log(\beta)\,. \tag{25}$$

Next fix a layer $i$ and let $(i, k)$ and $(i-1, j)$ be neurons maximising

$$B_i = \sum_{a \in \mathcal{C}} \left( \int_{c_{ik}^{-1}(a)} \left( f(z) - p_{i-1,j}^*(z) \right) d\mu(z) \right)^2 \in [0, 1].$$

Let $\Delta_a = \int_{c_{ik}^{-1}(a)} (f(z) - p_{i-1,j}^*(z)) d\mu(z)$ and $\tilde{\Delta}_a = \text{sign}(\Delta_a) \min\{|\Delta_a|, \delta/4\}$. Then by Lemma 7 we have

$$\begin{aligned}
D_i^* \leq D_{ik} &= \sum_{a \in \mathcal{C}} \min_{w \in \mathcal{W}_{ik}} \int_{c_{ik}^{-1}(a)} \mathcal{D}(f(z), \sigma(w \cdot p_{i-1}^*(z))) d\mu(z) \\
&\leq \sum_{a \in \mathcal{C}} \int_{c_{ik}^{-1}(a)} \mathcal{D}\left( f(z), \sigma\left( 4\tilde{\Delta}_a + \text{logit}(p_{i-1,j}^*(z)) \right) \right) d\mu(z) \\
&\leq \sum_{a \in \mathcal{C}} \int_{c_{ik}^{-1}(a)} \left( \mathcal{D}(f(z), p_{i-1,k_{i-1}}^*(z)) + 4\tilde{\Delta}_a(p_{i-1,j}^*(z) - f(z)) + 2\tilde{\Delta}_a^2 \right) d\mu(z) \\
&\leq \int_{\mathbb{R}^d} \mathcal{D}(f(z), p_{i-1,j}^*(z)) d\mu(z) - 4 \sum_{a \in \mathcal{C}} \Delta_a \tilde{\Delta}_a + 2 \sum_{a \in \mathcal{C}} \tilde{\Delta}_a^2 \\
&\leq \int_{\mathbb{R}^d} \mathcal{D}(f(z), p_{i-1,j}^*(z)) d\mu(z) - 2 \sum_{a \in \mathcal{C}} \Delta_a \tilde{\Delta}_a \\
&\leq \int_{\mathbb{R}^d} \mathcal{D}(f(z), p_{i-1,j}^*(z)) d\mu(z) - 2 \min\left\{ 1, \frac{\delta}{4} \right\} \sum_{a \in \mathcal{C}} \Delta_a^2 \\
&= D_{i-1,j} - \min\left\{ 2, \frac{\delta}{2} \right\} B_i = D_{i-2}^* - \min\left\{ 2, \frac{\delta}{2} \right\} B_i.
\end{aligned}$$

Therefore by telescoping the sum over odd and even natural numbers and Eq. (25) we have

$$\sum_{i=1}^{\infty} B_i \leq \frac{2}{\min\left\{2, \frac{\delta}{2}\right\}} \log\left(\frac{1}{\beta}\right) = \max\left\{1, \frac{4}{\delta}\right\} \log\left(\frac{1}{\beta}\right),$$

which implies the result.

**Part (3)** Let $(i, k)$ and $(i-1, j)$ be two non-bias neurons and abbreviate $p(z) = p_{ik}^*(z) \in \mathbb{R}$ and $q(z) = p_{i-1,j}^*(z) \in \mathbb{R}^{K_{i-1}}$. Let

$$\varepsilon = \int_{\mathbb{R}^d} \left( \mathcal{D}(f(z), q(z)) - \mathcal{D}(f(z), p(z)) \right) d\mu(z).$$

Then by Lemma 8 and Pinsker's inequality

$$\begin{aligned}
\varepsilon &= \int_{\mathbb{R}^d} \left( \mathcal{D}(p(z), q(z)) + \frac{\partial}{\partial \alpha} \mathcal{D}(f(z), \sigma((1-\alpha)\text{logit}(p(z)) + \alpha\,\text{logit}(q(z)))) \Big|_{\alpha=0} \right) d\mu(z) \\
&\geq \int_{\mathbb{R}^d} \mathcal{D}(p(z), q(z)) d\mu(z) \geq 2 \int_{\mathbb{R}^d} (p(z) - q(z))^2 d\mu(z),
\end{aligned}$$

35

where the first inequality follows from the definition of $p = p_{ik}^*$, which ensures the derivative is nonnegative and the second follows from Pinsker's inequality. Therefore for all layers $i$ and non-bias neurons $(i, k)$ and $(i-1, j)$ it holds that

$$\int_{\mathbb{R}^d} \left( p_{ik}^*(z) - p_{i-1,j}^*(z) \right)^2 d\mu(z) \leq \frac{1}{2} \int_{\mathbb{R}^d} \left( \mathcal{D}(f(z), p_{i-1,j}^*(z)) - \mathcal{D}(f(z), p_{i,k}^*(z)) \right) d\mu(z)$$

$$\leq \frac{1}{2} (D_{i-2}^* - D_i^*) .$$

By telescoping the sum over odd and even $i$ there exists a $p_\infty^*$ such that

$$\lim_{i \to \infty} \max_{k > 1} \int_{\mathbb{R}^d} (p_{ik}^*(z) - p_\infty^*(z))^2 d\mu(z) = 0 .$$

Since $p_{ik}^*$ is $\mathcal{F}$-measurable for all neurons $(i, k)$ and pointwise convergence preserves measurability it follows that $p_\infty^*$ is also $\mathcal{F}$-measurable. ∎

## Appendix C. Proof of Theorem 3

We use the same notation as the proof of Theorem 1 in the previous section. Only the first part is different. Since the side information is no longer independent and identically distributed we cannot claim anymore that $\mathbb{E}[\ell_t^*(u)] = \ell^*(u)$. The idea is to use the stability of the learning procedure to partition the data into chunks of increasing size on which the weights are slowly changing, but where the mixing of the Markov chain will eventually ensure that the empirical distribution converges to stationary. Let $(S_i)_i$ be a random contiguous partition of $N(\infty) = \{t \in \mathbb{N} : c_{ik}(z_t) = a\}$, which means that $\max S_i < \min S_{i+1}$ for all $i$ and $\bigcup_{i=1}^\infty S_i = N(\infty)$. By the stability assumption we may choose $(S_i)_i$ such that $\lim_{i \to \infty} |S_i| = \infty$ and $\max_{s,t \in S_i} \left\| w^{(s)} - w^{(t)} \right\|_\infty \leq \varepsilon_i$ with $\lim_{i \to \infty} \varepsilon_i = 0$ for all possible data sequences. Then abusing notation by letting $w^{(i)} = w^{(\min S_i)}$ and letting $n_I = \sum_{i=1}^I |S_i|$ we have the following holding almost surely:

$$0 \geq \lim_{n \to \infty} \frac{1}{n} \sum_{t=1}^n \left( \ell_t^*(w^{(t)}) - \ell_t^*(u) \right) = \lim_{I \to \infty} \frac{1}{n_I} \sum_{i=1}^I \sum_{t \in S_i} \left( \ell_t^*(w^{(t)}) - \ell_t^*(u) \right)$$

$$= \lim_{I \to \infty} \frac{1}{n_I} \sum_{i=1}^I \sum_{t \in S_i} \left( \ell_t^*(w^{(i)}) - \ell_t^*(u) + \ell_t^*(w^{(t)}) - \ell_t^*(w^{(i)}) \right)$$

$$\geq \lim_{I \to \infty} \frac{1}{n_I} \sum_{i=1}^I \sum_{t \in S_i} \left( \ell_t^*(w^{(i)}) - \ell_t^*(u) - \varepsilon_i \right)$$

$$= \lim_{I \to \infty} \frac{1}{n_I} \sum_{i=1}^I \sum_{t \in S_i} \left( \ell_t^*(w^{(i)}) - \ell_t^*(u) \right) = \lim_{I \to \infty} \frac{1}{n_I} \sum_{i=1}^I |S_i| \left( \ell^*(w^{(i)}) - \ell^*(u) \right)$$

$$= \lim_{I \to \infty} \frac{1}{n_I} \sum_{i=1}^I \sum_{t \in S_i} \left( \ell^*(w^{(t)}) - \ell^*(u) \right) = \lim_{n \to \infty} \frac{1}{n} \sum_{t=1}^n \left( \ell^*(w^{(t)}) - \ell^*(u) \right) ,$$

where we used the convergence theorem for aperiodic $\phi$-irreducible Markov chains (Meyn and Tweedie, 2012, Chap. 10). The proof is concluded in the same way as the proof of Theorem 1.

## Appendix D. Norms and topological arguments

Recall that $(\mathbb{R}^d, \mathcal{F}, \mu)$ is a probability space with $\mathcal{F}$ the Lebesgue $\sigma$-algebra and for $\mathcal{G}$ a set of $\mathcal{F}$-measurable functions from $\mathbb{R}^d$ to $\mathcal{C}$ we define

$$\|h\|_{\mathcal{G}} = \sup_{c \in \mathcal{G}} \sum_{a \in \mathcal{C}} \left| \int_{c^{-1}(a)} h(x) d\mu(x) \right|.$$

Of course $\|\cdot\|_{\mathcal{G}}$ is not a norm for all $\mathcal{G}$. In the following we establish that it is a norm for a few natural choices. First we note that $\|\cdot\|_{\mathcal{G}}$ is 'almost' a norm for all $\mathcal{G}$ as summarised in the following trivial lemma.

**Lemma 9** *Let $g, h : \mathbb{R}^d \to \mathbb{R}$ be $\mu$-integrable and $\mathcal{G} \subset \mathcal{C} \to \mathbb{R}^d$. Then*

1. $\|g + h\|_{\mathcal{G}} \leq \|g\|_{\mathcal{G}} + \|h\|_{\mathcal{G}}$

2. $\|ag\|_{\mathcal{G}} = a \|g\|_{\mathcal{G}}$

3. $\|0\|_{\mathcal{G}} = 0$.

Thus in order to show that $\|\cdot\|_{\mathcal{G}}$ is a norm it is only necessary to prove that $\|h\|_{\mathcal{G}} = 0$ implies that $h = 0$ $\mu$-almost-everywhere. We are especially interested in countable $\mathcal{G}$, since we can construct (infinite) networks that eventually use all contexts in such sets. The following lemma allows us to connect Theorem 1 to a claim about the asymptotic universality of a network.

**Lemma 10** *If $\mathcal{G} = \{c_1, c_2, \ldots\}$ is countable and $\|\cdot\|_{\mathcal{G}}$ is a norm on the space of $\mu$-integrable $h$, then for $\mathcal{G}_i = \{c_1, c_2, \ldots, c_i\}$ it holds that $\lim_{i \to \infty} \|h\|_{\mathcal{G}_i} = \|h\|_{\mathcal{G}}$ for all $h$.*

For the remainder we assume that $\mu$ is absolutely continuous with respect to the Lebesgue measure.

**Lemma 11** *Let $\mathcal{G} = \{\mathbb{1}_{B_r(x)} : x \in \mathbb{Q}^d, r \in (0, \infty) \cap \mathbb{Q}\}$ and suppose that $h$ is bounded and measurable with $\|h\|_{\mathcal{G}} = 0$. Then $h = 0$ $\mu$-almost-everywhere.*

**Proof** Assume without loss of generality that $\sup_{x \in \mathbb{R}^d} |h(x)| \leq 1$. Suppose that $|h|$ does not vanish $\mu$-almost everywhere. Then there exists an $\varepsilon > 0$ and $\delta > 0$ such that $\mu(|h| > \varepsilon) > \delta$. Let $A = \{x : |h(x)| > \varepsilon\}$, which is measurable because $|h|$ is measurable. Therefore by outer-regularity there exists an open $U \subseteq \mathbb{R}^d$ with $A \subset U$ and $\mu(U) \leq \mu(A) + \varepsilon \delta / 2$. Therefore $\int_U h(x) d\mu(x) = \int_A h(x) d\mu(x) + \int_{U-A} h(x) d\mu(x) \geq \varepsilon \delta / 2$. Next write $U = \bigcup_i B_i$ where $(B_i)$ are disjoint open balls with rational-valued centers and radii. Then by assumption $\int_U h(x) d\mu(x) = \sum_{i=1}^{\infty} \int_{B_i} h(x) d\mu(x) = 0$, which is a contradiction. ∎

**Lemma 12** *Let $\mathcal{B}$ be a countable base for $\mathbb{R}^d$ and $\mathcal{G} = \{\mathbb{1}_U : U \in \mathcal{B}\}$ and $\|h\|_{\mathcal{G}} = 0$ for some bounded and $\mathcal{F}$-measurable $h$. Then $h = 0$ $\mu$-almost-everywhere.*

**Proof** As above. ∎

**Lemma 13** *Let $\mathcal{G} = \{\mathbb{1}_{H_{\nu,c}} : \nu \in \mathbb{Q}^d, \|\nu\| = 1, c \in \mathbb{Q}\}$. If $h : \mathbb{R}^d \to \mathbb{R}$ is measurable and $\|h\|_{\mathcal{G}} = 0$, then $h = 0$ almost-everywhere.*

**Proof** Since $\mathbb{Q}$ is dense in $\mathbb{R}$ it follows from absolutely continuity of $\mu$ that if $\nu \in \mathbb{R}^d$ and $c \in \mathbb{R}$ and $\|h\|_{\mathcal{G}} = 0$, then $\int_H h d\mu = 0$ for all half-spaces $H$. Therefore if $\lambda$ is the Lebesgue measure, then

$$\hat{h}(\partial H) = \int_{\partial H} h \frac{d\mu}{d\lambda} d\lambda = 0$$

for all hyperplanes $\partial H$. But $\hat{h}$ is the Radon transform of $h d\mu/d\lambda$, which is zero everywhere only if $h d\mu/d\lambda$ is zero almost everywhere (Helgason, 2011). Therefore $h\frac{d\mu}{d\lambda} \equiv 0$ and so $h \equiv 0$ almost everywhere. ∎

## Appendix E. Capacity for Half-Space Contexts with Two Layers

We now argue that isolation networks can approximate continuous functions to arbitrary precision with just two layers where the first layer is sufficiently wide and consists of half-space contexts and the second layer (output) has just one neuron and no context. For dimension one we give an explicit construction, while for higher dimensions we sketch the argument via the Radon transform inversion formula. The interestingness of these results is tempered by the fact that the precision of the approximation can only be made arbitrarily small by having an arbitrarily wide network *and* that the weights must also be permitted to be arbitrarily large. In practice we did not find two-layer half-space contexts to be effective beyond the one-dimensional case.

**Theorem 14** *Let $f : \mathbb{R} \to [0, 1]$ be continuous and $\mu$ have compact support and be absolutely continuous with respect to the Lebesgue measure. Then for any $\varepsilon > 0$ there exists a two-layer network with half-space contexts in the first layer and a single non-bias neuron in the second layer with a trivial context function such that $\|p_{21}^* - f\|_\infty < \varepsilon$, where $p_{21}^*(z)$ is given in Theorem 1.*

**Proof** By applying a mollifier, any continuous $f$ may be approximated to arbitrary precision by a differentiable function with range a subset of $(\varepsilon, 1 - \varepsilon)$ for small $\varepsilon > 0$. From now on we simply assume that $f$ is differentiable. For $b \in (0, 1)$ let

$$\ell(b) = \text{logit}\left(\frac{1}{\mu([b,1])}\int_b^1 f(z)d\mu(z)\right) \qquad \bar{\ell}(b) = \text{logit}\left(\frac{1}{\mu([0,b])}\int_0^b f(z)d\mu(z)\right)$$

and $\ell = \int_0^1 f(z)d\mu(z)$. Because $f(z) \in (\varepsilon, 1-\varepsilon)$ it follows that $\ell(b)$ and $\bar{\ell}(b)$ are also bounded. By the first part of Theorem 1, if neuron $(1, k)$ uses context function $c_{1k}(z) = \mathbb{1}_{z \geq b_k}(z)$, then

$$\text{logit}(p_{1k}^*(z)) = \begin{cases} \ell(b_k) & \text{if } z \geq b_k \\ \bar{\ell}(b_k) & \text{otherwise.} \end{cases}$$

Since neuron $(2, 1)$ has a trivial context function there exists an $a \in \mathcal{C}$ such that $c_{21}(z) = a$ for all $z \in \mathbb{R}$. Abbreviating $w_k = w_{21ak}$ we can write the output of neuron $(2, 1)$ as

$$p_{21}^*(z) = \sigma \left( \sum_{k=0}^{K_1} w_k \, \text{logit}(p_{1k}^*(z)) \right)$$

$$= \sigma \left( w_0 + \sum_{k: b_k \leq z} w_k \ell(b_k) + \sum_{k: b_k > z} w_k \bar{\ell}(b_k) \right)$$

$$= \sigma \left( w_0 + \sum_{k: b_k \leq z} w_k (\ell(b_k) - \bar{\ell}(b_k)) + \sum_k w_k \bar{\ell}(b_k) \right) . \tag{26}$$

The main idea is to write an approximation of $f$ as an integral that in turn may be approximated by the sum above for sufficiently large $K_1$ and well chosen $w_k$. Define

$$v(x) = \frac{f'(x)}{f(x)(1 - f(x))} \left( \frac{1}{\ell(x) - \bar{\ell}(x)} \right) .$$

By differentiating the logit function and the fundamental theorem of calculus and the assumption that $f$ is differentiable and $f(z) \in (\varepsilon, 1-\varepsilon)$ we have

$$f(z) = \sigma \left( \text{logit}(f(0)) + \int_0^z \frac{f'(x)}{f(x)(1 - f(x))} dx \right)$$

$$= \sigma \left( \text{logit}(f(0)) + \int_0^z v(x)(\ell(x) - \bar{\ell}(x)) dx \right) .$$

The result follows by approximating the integral above with the sum in Eq. (26). $\blacksquare$

The situation for higher dimensions is more complicated, but the same theorem may be proven for $d > 1$ by appealing to the Radon transform inversion formula (Helgason, 2011, Chap 1), which says that any infinitely differentiable and compactly supported function $g : \mathbb{R}^d \to \mathbb{R}$ may be represented by

$$g(z) = \int_S \tilde{g}(\eta, \langle z, \eta \rangle) d\eta ,$$

where $S = \{x \in \mathbb{R}^d : \|x\|_2 = 1\}$ is the sphere and the integral is with respect to the uniform measure and $\tilde{g} : S \times \mathbb{R} \to \mathbb{R}$ is carefully chosen. The precise form of $\tilde{g}$ is quite complicated, but only depends on the Radon transform of $g$, which is the operator $(Rg) : S \times \mathbb{R} \to \mathbb{R}$

given by $(Rg)(\eta, x) = \int_{\partial H_{\eta,x}} f(z)dz$. This result means that sufficiently regular functions are entirely determined by their integrals on hyperplanes. It remains to note that for wide two-layer networks with arbitrarily large weights and half-space contexts, the output of neuron $(2, 1)$ can approximate the above integral to arbitrary precision. We leave the details for the interested reader.

## Appendix F. Proof of Equation 18

This result is the same as that of Lemma 1 in (Veness et al., 2012). We provide the proof here for the sake of completeness.

**Lemma 15** *If* $s(\nu_{1:n}) = \sum_{k=2}^{n} \mathbb{1}[\nu_k \neq \nu_{k-1}]$, *then for all* $\nu_{1:n} \in \mathcal{I}_n(\mathcal{M})$, *we have*

$$- \log w_\tau(\nu_{1:n}) \leq (s(\nu_{1:n}) + 1) \left( \log |\mathcal{M}| + \log_2 n \right).$$

**Proof** Consider an arbitrary $\nu_{1:n} \in \mathcal{I}_n(\mathcal{M})$. Letting $m = s(\nu_{1:n})$, from Equation 17 we have that

$$
\begin{aligned}
- \log w_\tau(\nu_{1:n}) &= \log |\mathcal{M}| - \log \prod_{t=2}^{n} \left( \tfrac{t-1}{t} \mathbb{1}[\nu_t = \nu_{t-1}] + \tfrac{1}{t(|\mathcal{M}|-1)} \mathbb{1}[\nu_t \neq \nu_{t-1}] \right) \\
&\leq \log |\mathcal{M}| - \log \prod_{t=2}^{n} \left( \tfrac{t-1}{t} \mathbb{1}[\nu_t = \nu_{t-1}] + \tfrac{1}{n(|\mathcal{M}|-1)} \mathbb{1}[\nu_t \neq \nu_{t-1}] \right) \\
&\leq \log |\mathcal{M}| - \log \left( n^{-m}(|\mathcal{M}|-1)^{-m} \prod_{t=2}^{n-m} \tfrac{t-1}{t} \right) \\
&= \log |\mathcal{M}| + m \log n + m \log(|\mathcal{M}|-1) + \log(n-m) \\
&\leq (m+1)[\log |\mathcal{M}| + \log n].
\end{aligned}
$$

■