
Holographic Feature Representations of Deep Networks

Martin A. Zinkevich
Google, Inc.

Alex Davies
Google, Inc.

Dale Schuurmans*
University of Alberta

Abstract

It is often asserted that deep networks learn “features”, traditionally expressed by the activations of intermediate nodes. We explore an alternative concept by defining features as partial derivatives of model output with respect to model parameters—extending a simple yet powerful idea from generalized linear models. The resulting features are not equivalent to node activations, and we show that they can induce a *holographic* representation of the complete model: the network’s output on given data can be exactly replicated by a simple *linear* model over such features extracted from *any* ordered cut. We demonstrate useful advantages for this feature representation over standard representations based on node activations.

1 INTRODUCTION

Deep networks provide an effective foundation for machine learning, having delivered significant advances in computer vision [10, 26, 17, 19, 8], speech recognition [7], spoken dialogue systems [9], and machine translation [9, 22]. Such models are often said to learn useful “features” [4, 23], which are normally considered to be the activations of intermediate nodes. However, we consider a more fundamental concept than node activations, and show that this alternative notion of what constitutes a “feature” can provide more powerful properties.

The standard way to think about the features learned in a deep network—as node activations across a single layer [4, 25, 26, 23]—has led to some insight into what a trained network might be representing. For example, such a perspective has led to impressive visualizations

that illustrate how such models might “see” classes [20] or interpret images [15]. The layer-wise embeddings in deep networks trained on text have also been visualized and shown to encode semantically meaningful dimensions [14]. Attempts have also been made to “transfer” useful activation features between tasks [23], and to understand layer-wise representations by using them in linear classifiers [1]. In fact, this latter work bears some similarity to the present investigation, but without considering equally powerful feature representations.

Meanwhile, a key drawback of deep networks remains the stability and repeatability of training. It is well known that non-convex optimization problems can have suboptimal local minima separated from global minima, which also occurs in training deep networks [5]. It is true that local minima do not prevent good outcomes from being obtained in practice [3]; in fact, whenever one is attempting to produce a single model given a fixed training set (i.e. the “offline” scenario) any weak result can simply be discarded and training repeated. *However, the situation is quite different in the typical “online” scenario encountered in industry*, where data arrives continually and models must be constantly retrained or adjusted for user-facing applications. In such cases, training instability can at minimum be a nuisance and at worst be dangerous. The online scenario presents an important challenge for deep networks and raises some of the key questions we attempt to address in this paper. Can features be extracted from a deep network that can be reliably used in other models? Can extracted features be used to more efficiently and stably incorporate new data?

In this paper, we provide a different perspective on what a “feature” might represent about a deep network. We show how the specific features we extract can be used to recover a linear model that *exactly* reproduces the original network on the training data. These features incorporate both the “upstream” information from previous layers, and the “downstream” information from later layers, giving a complete “holographic” representation

* Work supported by Google.

of the network from the perspective of any given layer. An important advantage is that this representation can be rapidly re-trained to global optimality on new data, efficiently obtaining repeatable outcomes. Moreover, the features are “calibrated” in a manner that allows weaknesses in the model versus the chosen training objective to be disambiguated, as we discuss in the next section.

2 FEATURES AND CALIBRATION

To investigate whether the rich notion of “feature” from generalized linear models can be meaningfully applied to deep networks, we leverage the concept of *calibration* (which is a particular form of *moment matching*).

Consider linear least squares regression. If a feature is 1 in each example (e.g., the bias), the average label must equal the average prediction in any optimally trained model. More generally, for any $\{0, 1\}$ -valued (boolean) feature, the average label must equal the average prediction when the feature is 1 [16, Equation 3]. This calibration property provides a powerful diagnostic: Imagine one is trying to predict the number of minutes someone who searches for “cat videos” will watch `cat_12`. Having a feature that is 1 when the search is “cat videos” and the video shown is `cat_12` guarantees that the average prediction will equal the average label in this subset. Thus, if you wanted to rank videos by how many minutes will be watched, you can disambiguate between mistakes in the label (e.g. `cat_12` is watched more on average than `cat_13` but is a worse video by a different metric: do we need to change our metric?) and mistakes in the modeling (e.g. `cat_12` is watched more on average than `cat_13`, but predicted on average to be watched less).

A similar property holds for logistic regression after distinguishing the *predicted log odds* (a linear combination of the features) from the *predicted probability* (a function of the log odds). In an optimally trained model, the average predicted probability will equal the average label whenever a boolean feature is 1. That is, if we tried to predict the probability someone would click on video `cat_12` given that they search for “cat videos”, the average predicted probability would be the average label.

This notion of calibration provides a practical tool for disambiguating problems with the feature representation versus problems with the objective. We use it as a key defining concept in the technical development below.

3 FORMALISM

Our basic strategy is to relate deep networks to generalized linear models—such as linear, logistic and Poisson regression—by considering supervised learning prob-

lems defined in terms of exponential families [16, 21].

3.1 EXPONENTIAL FAMILIES

We base our development on full, minimal, 1-dimensional, standard exponential families [2] (see Appendix A for standard definitions), which will be used to formulate the training loss. We will denote the exponential family being used by p , from which one can derive the domain $\Omega_p \subseteq \mathbf{R}$ of possible outcomes, the mean function $\mu_p : \mathbf{R} \rightarrow \mathbf{R}$ of the outcome given the parameter, and the loss $\ell_p : \Omega_p \times \mathbf{R} \rightarrow \mathbf{R}$ given by the negative log likelihood of the outcome given the parameter.

For example, the Bernoulli family is specified by $\Omega_p = \{0, 1\}$, $\ell_p(\omega, \theta) = -\theta\omega + \ln(1 + \exp(\theta))$, and $\mu_p(\theta) = \frac{1}{1 + \exp(\theta)}$; this family forms the basis of logistic regression. Other examples, such as Gaussian with a fixed variance (least squares regression), and the Poisson family (Poisson regression) are given in Appendix A. The following lemma encapsulates the key facts we require:

Lemma 1 *The Bernoulli family, Gaussian family (fixed variance), and Poisson family are full, minimal, 1-dimensional standard exponential families. For a full, minimal, 1-dimensional standard exponential family p :*

1. for every $\theta \in \mathbf{R}$, $\omega \in \Omega_p$, $\ell_p(\omega, \theta)$ is differentiable with respect to θ and $\frac{\partial \ell_p(\omega, \theta)}{\partial \theta} = \mu_p(\theta) - \omega$; and
2. ℓ_p is **strictly convex in its second argument**: for every $\theta, \theta' \in \mathbf{R}$, if $\theta \neq \theta'$, then for all $\lambda \in (0, 1)$ we have $\ell_p(\omega, \lambda\theta + (1 - \lambda)\theta') < \lambda\ell_p(\omega, \theta) + (1 - \lambda)\ell_p(\omega, \theta')$.

Proofs for all lemmas and theorems stated in this paper are given in the appendices.

3.2 SUPERVISED LEARNING PROBLEM

Our focus in this paper will be on supervised learning problems. A **model** is defined to be a function $M : X \rightarrow \mathbf{R}$. A **model family** is given by a parameterized set of models $\mathcal{M} = \{M_w\}_{w \in \mathbf{R}^S}$ where S is a finite set of parameters. We define a **supervised learning problem** (or just a **problem** for short) to be a tuple $\mathcal{P} = (p, X, \{x_1 \dots x_m\}, \{y_1 \dots y_m\})$, where p is a full, minimal, 1-dimensional standard exponential family, X is the domain of instances, $x_1 \dots x_m \in X$ are the instances, and $y_1 \dots y_m \in \Omega_p$ are the labels. We define a **loss function** $L_{\mathcal{P}} : \mathbf{R}^X \rightarrow \mathbf{R}$ that maps a model M to a scalar loss value via:

$$L_{\mathcal{P}}(M) = \sum_{i=1}^m \ell_p(y_i, M(x_i)), \quad (1)$$

the negative log likelihood of the labels given instances.

3.3 CALIBRATION, FEATURES, OPTIMALITY

Definition 2 Given a supervised learning problem $\mathcal{P} = (p, X, \{x_1 \dots x_m\}, \{y_1 \dots y_m\})$,¹ a **feature** is a function $f : X \rightarrow \mathbf{R}$. We say that a model M is **calibrated with respect to feature f on \mathcal{P}** if:

$$\sum_{i=1}^m f(x_i) \mu_p(M(x_i)) = \sum_{i=1}^m f(x_i) y_i. \quad (2)$$

For example, if the feature $f : X \rightarrow \mathbf{R}$ is 1 if the instance comes from the United States, and 0 otherwise, then calibration of a model with respect to this feature would imply that the average prediction given an instance from the United States would equal the average label for an instance from the United States. Note that the concept of calibration with respect to a feature is independent of whether the feature is actually used by the model M . In the literature on regression, there are many results that show how various generalized linear regression problems find models that are calibrated with respect to the features [16, 21]. In this paper, we will consider starting with a model and generating calibrated features from it.

Definition 3 For a model family $\{M_w\}_{w \in \mathbf{R}^S}$ with $w \in \mathbf{R}^S$ and $s \in S$, we define the **feature generated from parameter s of model M_w** to be $f_s : X \rightarrow \mathbf{R}$ such that:

$$f_s(x) = \frac{\partial^+ M_w(x)}{\partial w_s} \quad \text{for all } x \in X; \quad (3)$$

in other words, a feature is specified by the right partial derivative of the output of model M_w acting on x with respect to some parameter s .²

Definition 4 Given a parameter subset $S' \subseteq S$, the features generated from S' of M_w are **total for $x \in X$** if:

$$\lim_{h \rightarrow 0} \frac{|M_{w+h}(x) - (M_w(x) + \sum_{s' \in S'} h_{s'} f_{s'}(x))|}{\|h\|} = 0 \quad (4)$$

(note $h \in \mathbf{R}^{S'}$), where $f_{s'}$ is the feature generated from s' of model M_w , and $w + h \in \mathbf{R}^S$ obeys $(w + h)_{s'} = w_{s'} + h_{s'}$ if $s' \in S'$, and $(w + h)_{s'} = w_{s'}$ if $s' \notin S'$. For brevity, we will also call such a set S' itself **total**.³

First note that the property of a feature set S' being total exhibits downward inheritance: if S' is total for some

¹ When we mention a feature in the context of a model or problem, we assume the feature shares the same domain.

² Note that this definition is related but not identical to the Fisher score, given by $\partial \log \ell_p(y, M_w(x)) / \partial w_s$. A key difference is that the Fisher score requires the label y , whereas a generated feature is only defined with respect to model output.

³ Appx. I.1 shows how totality relates to differentiability.

x , then so is S'' for any $S'' \subseteq S'$. Therefore, if $M_w(x)$ is differentiable at w , implying that the full feature set S is total, then any generated feature set must also be total. Note that it is not sufficient that $M_w(x)$ be partially differentiable with respect to w to ensure totality; specifically, if the features generated from S' are total, and the features generated from S'' are total, there is no guarantee that the features generated from $S' \cup S''$ are total—Appendix E provides a concrete example.

Definition 5 Given a problem $\mathcal{P} = (p, X, \{x_1 \dots x_m\}, \{y_1 \dots y_m\})$, two models M and M' are **equal on the training data** if for all $i \in \{1 \dots m\}$, $M(x_i) = M'(x_i)$.

Definition 6 Given a problem \mathcal{P} and a model family \mathcal{M} , we say that $M \in \mathcal{M}$ is **optimal for \mathcal{P} on \mathcal{M}** if, for all $M' \in \mathcal{M}$, $L_{\mathcal{P}}(M) \leq L_{\mathcal{P}}(M')$.

Notice that optimality is based upon both the problem and the model family: this will become important later when we talk about models outside of a family being equivalent to optimal models inside a family.

In what follows we will distinguish these generated features from the more common node activations. In some cases they are the same, as in the last layer, other times they are related, and sometimes they are completely different, as in the generated features from bias parameters.

4 GENERATED FEATURES

We first show that if the partial derivative of the loss with respect to a parameter is zero, then the feature generated by that parameter must be calibrated.

Theorem 7 For a problem \mathcal{P} and model family $\mathcal{M} = \{M_w\}_{w \in \mathbf{R}^S}$, given a $w \in \mathbf{R}^S$ and $s \in S$ such that f_s is the feature generated from parameter s of model M_w : if $\{s\}$ is total on the training data given M_w and $\frac{\partial L_{\mathcal{P}}(M_w)}{\partial w_s} = 0$, then M_w is calibrated with respect to f_s .

Note that we do not assume that the model family \mathcal{M} is linearly parameterized; far from it, below we will be considering model families defined by certain classes of deep neural networks. Nevertheless, the features generated in this way can be exported to an external family of linear models defined over the same feature set. Such an auxiliary family of linear models will possess several strong properties with respect to the original nonlinear model family, as we now show.

Definition 8 For a finite S , given a set of functions $\{f_s\}_{s \in S}$ where $f_s : X \rightarrow \mathbf{R}$, we define the **family of linear models induced by $\{f_s\}_{s \in S}$** , to be the model family

$\{N_w\}_{w \in \mathbf{R}^S}$ such that for all $w \in \mathbf{R}^S$ and $x \in X$:

$$N_w(x) = \sum_{s \in S} w_s f_s(x). \quad (5)$$

Let this family be denoted by $\mathcal{L}(\{f_s\}_{s \in S})$.

Note that if we were generating features from a model M_w that was already linearly parameterized, we would simply re-generate the original feature set. Next, we need some preliminary results for families of linear models.

Lemma 9 For finite S , given a set of features $\{f_s\}_{s \in S}$, the family of linear models $\mathcal{L}(\{f_s\}_{s \in S})$, and $N_w \in \mathcal{L}$: the set S is total for all $x \in X$ given N_w , and the feature generated from parameter s by model N_w is f_s .

Lemma 10 (c.f. [16, Equation 10]) Given a problem \mathcal{P} , a finite set S and a set of features $\{f_s\}_{s \in S}$: a model N in the family of linear models $\mathcal{L}(\{f_s\}_{s \in S})$ is optimal if and only if N is calibrated with respect to f_s for all $s \in S$.

Lemma 11 Given a problem \mathcal{P} , a finite set S , a subset $S' \subseteq S$, and a set of features $\{f_s\}_{s \in S}$: if a model N is optimal in the family of linear models $\mathcal{L}(\{f_s\}_{s \in S'})$ and N is calibrated with respect to f_s for all $s \in S - S'$, then N is also an optimal model in $\mathcal{L}(\{f_s\}_{s \in S})$.

Now consider a general model family \mathcal{M} . Even if the models in \mathcal{M} are not linear, our first key result is that a family of linear models defined over the features generated from \mathcal{M} (via Definition 3) will satisfy important properties.

Definition 12 Given a finite S , family $\{M_w\}_{w \in \mathbf{R}^S}$ of models, subset $S' \subseteq S$, and parameters $w \in \mathbf{R}^{S'}$: if f_s is the feature generated by s given M_w , then the **family of linear models associated with S' given M_w** will be denoted $\mathcal{L}(\{f_s\}_{s \in S'})$.

Lemma 13 Given a problem \mathcal{P} , a finite set S , a subset $S' \subseteq S$, a model family $\{M_w\}_{w \in \mathbf{R}^S}$, and a $w \in \mathbf{R}^{S'}$ such that $\frac{\partial L_{\mathcal{P}}(M_w)}{\partial w_s} = 0$ for all $s \in S'$: if \mathcal{L}' is the family of linear models associated with S' given M_w , S' is total on the training data given M_w , and M_w is equal on the training data to some $N' \in \mathcal{L}'$, then N' is optimal in \mathcal{L}' .

This next lemma goes one step further, by showing that if a model $M \in \mathcal{M}$ is stationary on a subset of parameters S' and if there is a linear model over features generated on S' that matches M , then there is an optimal linear model over features generated on S that also matches M .

Lemma 14 Given a problem \mathcal{P} , a finite set S , a subset $S' \subseteq S$, a model family $\{M_w\}_{w \in \mathbf{R}^S}$, and a $w \in \mathbf{R}^{S'}$ such

that $\frac{\partial L_{\mathcal{P}}(M_w)}{\partial w_s} = 0$ for all $s \in S'$: if \mathcal{L}' is the family of linear models associated with S' given M_w , S' is total, M_w is equal on the training data to some $N' \in \mathcal{L}'$, and \mathcal{L}'' is the family of linear models associated with S , then any optimal $N'' \in \mathcal{L}''$ equals M_w on the training data.

So a key question will be to understand when a model M_w that is stationary on a subset S' of parameters will also be equal on the training data to a stationary model that uses all the parameters S . To address deep neural networks we will introduce the concept of *homogeneity*, which will allow us to prove that, for certain families of deep models, if a model is stationary with respect to a particular subset of its parameters it will also be equal on the training data to a linear model over the same subset. We will use this construction to show how the generated features can capture *all* of the expressiveness of a given deep model with respect to the training data, while still allowing one to re-train in a linear setting.

5 FEEDFORWARD NETWORKS

We first need to formally define deep neural networks for subsequent analysis. For generality we follow the formalization given in [18] that extends the more conventional layered representation; see also Appendix B.

A **feedforward neural network** D is defined by a directed acyclic graph with objects attached to the vertices and edges: in particular, $D = (V, E, I, \{g_i\}_{i \in I}, o^*, A)$, where V is a set of vertices, $E \subseteq V \times V$ is a set of edges, $I = \{i_1 \dots i_m\} \subset V$ is a set of input vertices, $g_i : X \rightarrow \mathbf{R}$ is an input function connected to $i \in I$, $o^* \in V$ is the output vertex, and $A = \{a_v : v \in V\}$ is a set of activation functions, $a_v : \mathbf{R} \rightarrow \mathbf{R}$. The parameters are $w \in \mathbf{R}^E$.

We will need to refer to the partial ordering on vertices implied by $G = (V, E)$, where we assume G contains no cycles,⁴ the input vertices have no incoming edges (i.e. $(u, i) \notin E$ for all $i \in I, u \in V$), the output vertex is not an input (i.e. $o^* \notin I$) and the output vertex has no outgoing edges (i.e. $(o^*, v) \notin E$ for all $v \in V$). A directed acyclic graph defines a partial order \leq on vertices where $u \leq v$ if and only if there is a path from u to v .

Given a training input $x \in X$, the computation of the network D is specified by a circuit function $c_{v,w}$ that assigns values to each vertex based on the partial order:

$$c_{i,w}(x) = a_i(g_i(x)) \text{ for } i \in I; \quad (6)$$

$$c_{v,w}(x) = a_v \left(\sum_{u:(u,v) \in E} c_{u,w}(x) w_{(u,v)} \right) \text{ for } v \in V - I. \quad (7)$$

⁴ A **path** (v_1, \dots, v_k) is a sequence of vertices such that $(v_j, v_{j+1}) \in E$ for all j . A **cycle** is a path with $v_1 = v_k$.

Note that the activations on input and output nodes are usually the identity, i.e. $a_v(x) = x$ for $v \in I \cup \{o^*\}$. We can also add a bias vertex $b \in I$, with input function $g_b(x) = 1$ for all $x \in X$, so that adding an edge $(b, v) \in E$ ensures that vertex v receives an affine rather than a linear combination of its incoming circuit values.

Definition 15 Given $D = (V, E, I, \{g_i\}_{i \in I}, o^*, A)$, we will denote the **family of feedforward network models** by $\mathcal{D}(V, E, I, \{g_i\}_{i \in I}, o^*, A) = \{M_w\}_{w \in \mathbf{R}^E}$, where for all $w \in \mathbf{R}^E$ we let $M_w = c_{o^*, w}$.

6 HOMOGENEOUS MODEL FAMILIES

Our main theoretical development applies to feedforward neural networks with *homogeneous* activation functions.

Definition 16 A function $f: \mathbf{R}^p \rightarrow \mathbf{R}^q$ is **homogeneous** (of degree 1) [11] if for all $v \in \mathbf{R}^p$ and all $\lambda \geq 0$:

$$f(\lambda v) = \lambda f(v). \quad (8)$$

Homogeneity implies that a function is linear along any ray from the origin, and also that it can be decomposed as an inner product between inputs and partial derivatives.

Lemma 17 (Euler’s Homogeneous Function Theorem) (Degree 1 Case) If a homogeneous function $f: \mathbf{R}^p \rightarrow \mathbf{R}^q$ is differentiable at $x \in \mathbf{R}^p$, then for all $k \in \{1 \dots q\}$,

$$f_k(x) = \sum_{j=1}^p \frac{\partial f_k(x)}{\partial x_j} x_j. \quad (9)$$

An important example of a (degree 1) homogeneous function is the standard neural network activation function $\text{relu}: \mathbf{R} \rightarrow \mathbf{R}$, given by $\text{relu}(x) = \max(x, 0)$.

Fact 18 The *relu* and *leaky relu* [13] activation functions are homogeneous. Any linear function is homogeneous, therefore so are projection and identity. The sum or composition of homogeneous functions is homogeneous, therefore constant multiplication of a homogeneous function is also homogeneous.

We can now state our main results, first in terms of general families of homogeneous functions.

Definition 19 Given a finite set S , a subset $S' \subseteq S$, and an instance space X : we say that a model family $\{M_w\}_{w \in \mathbf{R}^S}$ is **homogeneous on the parameter set S'** if for all $w \in \mathbf{R}^S$, $x \in X$ and $\lambda \geq 0$,

$$M_v(x) = \lambda M_w(x), \quad (10)$$

where $v \in \mathbf{R}^S$, $v_s = \lambda w_s$ when $s \in S'$, and $v_s = w_s$ when $s \notin S'$.

Note that under this definition homogeneity is a property of a set of parameters not an individual parameter: if a model family is homogeneous on S' and on S'' , that does not imply it is homogeneous on $S' \cup S''$; however, if $S''' \subseteq S'$, it will be homogeneous on S''' .

Our main results show that homogeneity and totality allow a holographic feature set to be extracted, where an auxiliary linear model on the feature set can replicate the output of any model in the family over the training data.

Theorem 20 If a model family \mathcal{M} is homogeneous on the parameter set S' , $M \in \mathcal{M}$, \mathcal{L}' is the family of linear models associated with S' given M , and S' is total on the training data given M , then there exists $N \in \mathcal{L}'$ such that M and N are equivalent on the training data.

Theorem 21 Given a problem \mathcal{P} , if a model family $\{M_w\}_{w \in \mathbf{R}^S}$ is homogeneous on the parameter set S' , then for any $w \in \mathbf{R}^S$:

1. if \mathcal{L}' is the family of linear models associated with S' given M_w , S' is total on the training data given M_w , and $\frac{\partial L_{\mathcal{P}}(M_w)}{\partial w_s} = 0$ for all $s \in S'$, then M_w is equal on the training data to any optimal N' in \mathcal{L}' .
2. if \mathcal{L}' is the family of linear models associated with S given M_w , S' is total on the training data given M_w , and $\frac{\partial L_{\mathcal{P}}(M_w)}{\partial w_s} = 0$ for all $s \in S$, then M_w is equal on the training data to any optimal N' in \mathcal{L}' .

Theorem 21 can be proved simply by combining Theorem 20 with Lemmas 13 and 14. We thus reach the conclusion that homogeneity allows one to generate features and build an auxiliary linear model that can behave equivalently to the original model over the training data.

6.1 APPLICATION TO NEURAL NETWORKS

To apply these results to deep neural networks we need to consider the question of when a model family defined by a network specification can be homogeneous in the sense of Definition 19. Recall the definition from Section 5 where a feedforward neural network is defined in terms of a directed acyclic graph $G = (V, E)$. A **cut** of G is a partition of the vertices in V into two disjoint subsets, where a **cut set** is the set of edges in E between those two sets. We wish to partition G into two sets, $B \subseteq V$ (bottom) and $T \subseteq V$ (top), such that $B \cap T = \emptyset$, $I \subseteq B$ and $o^* \in T$. The cut set consists of edges (u, v) where $u \in B$ and $t \in T$. We will call such a cut an **ordered cut**, since $t \not\leq b$ for all $b \in B$, $t \in T$ under the partial order generated by G . Note that, since there are no incoming edges to I and $o^* \notin I$ is assumed, we have the following.

Lemma 22 For any feedforward neural network, there is an ordered cut, specifically $B = I$ and $T = V \setminus I$.

Features generated by the parameters on edges in the cut set turn out to be of critical importance. In particular, by considering a partition of the vertices in a feedforward network into an ordered cut—consisting of B , T and the cut set—we reach the deepest result in this paper: that the solutions to the deep network can be mimicked by the solutions to the corresponding linear model defined on the features generated by the ordered cut.

Theorem 23 Given a problem \mathcal{P} and a family of feedforward network models $\mathcal{D}(V, E, I, \{g_i\}_{i \in I}, \sigma^*, A) = \{M_w\}_{w \in \mathbf{R}^E}$, where all $a \in A$ are homogeneous: **if** B and T is an ordered cut of the feedforward network, such that E' is the cut set, **then**:

1. $\{M_w\}_{w \in \mathbf{R}^E}$ is homogeneous on E' ;
2. for some $w \in \mathbf{R}^E$, if E' is total on the training data given M_w , \mathcal{L}' is the family of linear models associated with E' given M_w , and $\frac{\partial L_{\mathcal{P}}(M_w)}{\partial w_e} = 0$ for all $e \in E'$, then M_w is equal on the training data to any optimal model in \mathcal{L}' ;
3. for some $w \in \mathbf{R}^E$, if \mathcal{L}' is the family of linear models associated with E given M_w , if E is total on the training data given M_w , and $\frac{\partial L_{\mathcal{P}}(M_w)}{\partial w_e} = 0$ for all $e \in E$, then M_w is equal on the training data to any optimal model in \mathcal{L}' .

This result is of particular significance to the problem of repeated training, where data is arriving over time and one needs to continually produce new models. In such a scenario, it is important that the models produced behave stably: roughly speaking, there needs to be some form of repeatability in the training process and consistency in the resulting model prediction errors. Such stability is important so that human beings, who must ultimately decide whether a machine learning system can be trusted in a continuous setting, can accept that a subsequent trained model will not make significantly worse mistakes than a trained model makes today. An auxiliary family of linear models improves this notion of stability in two key ways:

1. One can approximate an optimal model with techniques from convex optimization (when it exists).
2. All optimal models are semantically equivalent on the training data.

Thus, we obtain a new form of reproducibility in a continual learning process. Such stability does not imply that the model will remain unchanged in the future, and the data may change, but this provides a solid foundation.

Although the features associated with all parameters can be extracted, the fact that one can extract an equivalent linear model over a fraction of the features is a significant result. If one extracted all features, multiple distinct linear models would be optimal on the training data yet differ on new data, which would decrease stability.

It is pointed out in [21], Theorem 3.4, that there is a duality between entropy maximization and maximum likelihood. In fact, if we wish to maximize entropy subject to a set of constraints, then the member of an exponential family that satisfies those constraints maximizes entropy. Thus, when a linear model is calibrated on all of its input features, it is not only maximizing likelihood, it is maximizing entropy.⁵ If a deep network is at a local minimum where the derivative is equal to zero, it is equal to one of these linear models: thus, one can consider a locally optimal model in a deep network to be maximizing entropy subject to being calibrated on the extracted features.

7 OTHER ARCHITECTURES

Not all deep models are feedforward networks: recursive neural networks are not defined on a fixed graph, residual networks (ResNets) [8] fix some weights to be constant, and convolutional neural networks (CNNs) [12] have sets of edges with equal weights.

We can handle convolutional and residual neural networks by considering a function $w^* : \mathbf{R}^n \rightarrow \mathbf{R}^E$, where we choose $v \in \mathbf{R}^n$ and use $w = w^*(v)$ as the parameters for the network. Thus, for $\mathcal{D}(V, E, I, \{g_i\}_{i \in I}, \sigma^*, A)$ represented as $\{M_w\}_{w \in \mathbf{R}^E}$, we can write $\{Q_v\}_{v \in \mathbf{R}^n}$ such that $Q_v = M_{w^*(v)}$ for all $v \in \mathbf{R}^n$. For ResNets, we need to specify some $E^f \subseteq E$ and $w^f \in \mathbf{R}^{E^f}$, where w_e^f is always the weight of $(w^*(v))_e$. Define $E^d = E - E^f$. For CNNs, we need to specify a partition of the edges: it is easiest to do so with a function $\pi : E^d \rightarrow \{1 \dots n\}$. So, for any $v \in \mathbf{R}^n$ and all $e \in E$:

$$(w^*(v))_e = \begin{cases} w_e^f & \text{if } e \in E^f \\ v_{\pi(e)} & \text{otherwise} \end{cases} \quad (11)$$

Thus, we will consider a **family of RC neural networks**⁶ $\mathcal{RC}(V, E, I, \{g_i\}_{i \in I}, \sigma^*, A, w^*) = \{Q_v\}_{v \in \mathbf{R}^n}$. Notice that now, instead of having $|E|$ we have n parameters. Define $E_1 \dots E_n \subseteq E^d$ such that $E_s = \pi^{-1}(s)$, i.e. the conventional representation of the partition. Given a

⁵ There is a nuance: this particular kind of entropy maximization is with respect to a particular measure. For the Bernoulli or any multinomial family, this corresponds with the traditional concept of entropy (entropy with respect to the counting measure). However, conventional differential entropy, which is defined with respect to the Lebesgue measure, is maximized subject to an additional constraint on variance.

⁶RC stands for ResNet and CNN.

family of RC feedforward neural networks, a cut B, T is **well-behaved** if it is an ordered cut and there exists an $S \subseteq \{1 \dots n\}$ such that the cut set equals $\bigcup_{s \in S} E_s$.

Theorem 24 For a problem \mathcal{P} and family of RC feedforward network models $\mathcal{RC}(V, E, I, \{g_i\}_{i \in I}, o^*, A, w^*)$ denoted $\{Q_v\}_{v \in \mathbf{R}^n}$, where all $a \in A$ are homogeneous: if $E_1 \dots E_n$ is the partition of the dynamic parameters, and B and T are a well-behaved cut such that there is an $S \subseteq \{1 \dots n\}$ where $E' = \bigcup_{s \in S} E_s$ is the cut set; **then**

1. $\{Q_v\}_{v \in \mathbf{R}^n}$ is a homogeneous model family with respect to S ;
2. if for some $v \in \mathbf{R}^n$, $\frac{\partial L_{\mathcal{P}}(Q_v)}{\partial v_s} = 0$ for all $s \in S$, the set S is total on the training data given Q_v , and \mathcal{L}' is the family of linear models associated with S given Q_v , then Q_v is equal on the training data to any optimal N' in \mathcal{L}' ;
3. if for some $v \in \mathbf{R}^n$, $\frac{\partial L_{\mathcal{P}}(Q_v)}{\partial v_s} = 0$ for all $s \in \{1 \dots n\}$, the set $\{1 \dots n\}$ is total on the training data given Q_v , and \mathcal{L}' is the family of linear models associated with $\{1 \dots n\}$ given Q_v , then Q_v is equal on the training data to any optimal N' in \mathcal{L}' .

Thus, if the cut exists, we can extract all the features.

8 REGULARIZATION

To this point, we have assumed parameters can have any value, and there is no external cost to choosing one set of parameters over another. However, there is evidence that regularization yields better generalization [6, Chapter 7].

In our framework, regularization can be defined in a generic way. Given a model family \mathcal{M} , a regularization function $R : \mathcal{M} \rightarrow \mathbf{R}$ assigns a cost to each model, independent of how it fits the data.

Definition 25 Given a model family \mathcal{M} , we say that $M \in \mathcal{M}$ is **optimal in \mathcal{M} given the problem \mathcal{P} and a regularization function $R : \mathcal{M} \rightarrow \mathbf{R}$** , if for all $M' \in \mathcal{M}$, $L_{\mathcal{P}}(M) + R(M) \leq L_{\mathcal{P}}(M') + R(M')$.

This kind of optimality can sometimes be interpreted as maximum a posteriori rather than maximum likelihood optimization. For the remainder of this section, we assume we have a finite S and model family $\mathcal{M} = \{M_w\}_{w \in \mathbf{R}^S}$. For a regularization function $R : \mathcal{M} \rightarrow \mathbf{R}$, define $R^* : \mathbf{R}^S \rightarrow \mathbf{R}$ such that $R^*(w) = R(M_w)$. We assume R is (strictly) convex if R^* is (strictly) convex. Given a subset $S' \subseteq S$, we will say that a function $r : \mathbf{R}^S \rightarrow \mathbf{R}$ is **additively separable** with respect to S' if there exists $r^{S'} : \mathbf{R}^{S'} \rightarrow \mathbf{R}^S$, and $r^{S-S'} : \mathbf{R}^{S-S'} \rightarrow \mathbf{R}^S$ such that for all $w \in \mathbf{R}^S$,

$r(w) = r^{S'}(\pi^{S \rightarrow S'}(w)) + r^{S-S'}(\pi^{S \rightarrow S-S'}(w))$. R is additively separable with respect to S' if R^* is additively separable with respect to S' . We can consider when, for all $s \in S$, the partial derivative with respect to s is zero:

$$\frac{\partial L_{\mathcal{P}}(M_w) + R(M_w)}{\partial w_s} = 0. \quad (12)$$

This is a stationary point and might be a saddle point or a local minima. The most common regularizers are L1 (where $R(M_w) = \sum_{s \in S} |w_s|$) and L2 (where $R(M_w) = \sum_{s \in S} w_s^2$), but group lasso [24] and other regularizations are also possible. Note that L1 regularization and group lasso are convex, L2 regularization is strictly convex, L1 and L2 regularization are separable, and group L1 regularization is sometimes separable.

As before, we want to connect homogeneous families (such as some deep networks) with regularization to the associated linear family with regularization.

Theorem 26 Given a problem \mathcal{P} , a set S , a subset $S' \subseteq S$, a model family $\mathcal{M} = \{M_w\}_{w \in \mathbf{R}^S}$ that is homogeneous with respect to S' , a strictly convex regularization function $R : \mathcal{M} \rightarrow \mathbf{R}$ that is additively separable with respect to S' , and a model M_w where for all $s \in S'$:

$$\frac{\partial [L_{\mathcal{P}}(M_w) + R(M_w)]}{\partial w_s} = 0, \quad (13)$$

if S' is total on the training data given M_w , $\mathcal{L}' = \{N_v\}_{v \in \mathbf{R}^{S'}}$ is the family of linear models associated with S' given M_w , and a new regularizer $R' : \mathcal{L}' \rightarrow \mathbf{R}$ is defined such that $R'(N_v) = (R^*)^{S'}(v)$, **then** M_w must be equal on the training data to any optimal $N \in \mathcal{L}'$ given the supervised learning problem and regularization R' .

Thus, in all circumstances where we showed functions were homogeneous (for the generic feedforward networks and some CNNs and ResNets), we now know that for a separable regularizer we can map these to a linear model family. However, we cannot simply add back all the other generated features not present in the ordered cut, because there might be a parameterization across all the generated features that has the same predictions and lower regularization. This requires further study.

9 EXPERIMENTAL EVALUATION

We investigated the main assertions that:

1. In practice, the holographic features generated can faithfully recover the original classifier.
2. The holographic features generated are indeed calibrated with respect to the original classifier.
3. Training with the generated holographic features is more stable than training a neural network.

Data set	# examples	input dim	λ	layer width	hidden layers
Pima	750 of 768	9	1/720	4	2
Census	30000 of 32561	108	1/1490	10	3
MNIST-3v5	9000 of 11552	784	1/1490	100	3

Table 1: Data sets and neural network architectures used.

- The holographic features support generalization beyond the training set they were generated from.

We used three binary classification data sets, UCI Pima, UCI Census and MNIST (3 vs 5). The details of each data set and the respective neural network architectures used are given in Table 1. Since our experimental design required us to partition each data set into three disjoint subsets, we extracted the number of examples indicated from the initial portion of the original data then split the chosen data into three equal sized subsets. In particular, we consider training the neural network on the first third, training the extracted linear models on the middle third, and using the final third to assess generalization performance. We also normalized the input features in each case by subtracting the means and dividing by the range.

Faithfulness First we evaluate whether a linear classifier trained on holographic features can faithfully reproduce the predictions of the original neural network, whereas using activation features fails to do so. Specifically, we trained the neural network on the first data partition, extracted the holographic and activation feature sets for each layer-wise cut of the neural network, trained a linear model over these feature sets on the same data, then compared the predictions on the same training data against those made by the source neural network. Figures 1 to 3 show the agreement plots for the Pima, Census and MNIST data sets respectively. The plots show the results achieved by trained linear models over the activation and holographic feature sets respectively, as well as using exact weights for the holographic features. There are minor differences due to the neural network not being at a true local minimum, but the assertion appears to be verified.

Calibration Next, we investigate the assertion that the holographic features must be calibrated. To do so, we plot $\sum_{i=1}^m f(x_i)\mu_p(M(x_i))$ against $\sum_{i=1}^m f(x_i)y_i$ for each extracted set of holographic features over the training data. Figure 4 shows that the holographic features are indeed well calibrated, regardless of which cut was used to generate them, or which data set is considered.

Stability To determine whether a linearized representation improves learning stability compared to re-training a neural network, we conducted the following experiment. We trained a neural network on the first partition, extracted distinct holographic and activation feature sets

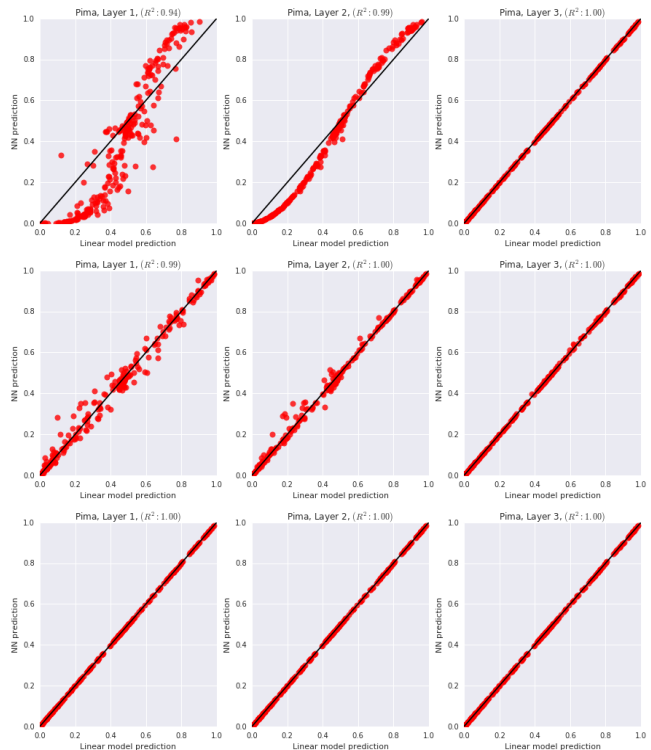


Figure 1: Classifier comparison for *Pima*. Training on activation features (row 1), holographic features (row 2), and closed form solution to holographic weights (row 3).

Model type	Pima	Census	MNIST-3v5
Neural net	5.9×10^{-5}	2.8×10^{-5}	1.1×10^{-4}
Activation	$< 10^{-15}$	$< 10^{-15}$	$< 10^{-15}$
Holographic	$< 10^{-15}$	$< 10^{-15}$	$< 10^{-15}$

Table 2: Average KL divergence between posteriors for different re-trainings of the same model on the same data.

for each layer-wise cut, then repeatedly re-trained the linear models with different random initializations on the second data partition, gathering the predictions made on the third data partition. We then re-trained the neural network on the original training data with different random initializations, and gathered the predictions made on the third data partition. To evaluate learning stability, we measured the average symmetrized KL divergence between the predictions of the different re-trained models for each type. The results are reported in Table 2. As expected, re-training produces nearly identical predictions for the linear models since the problem is convex. However, neural network re-training results in non-negligible variability between the different learned models.

Generalization Finally, we investigate whether holographic features support generalization to data from

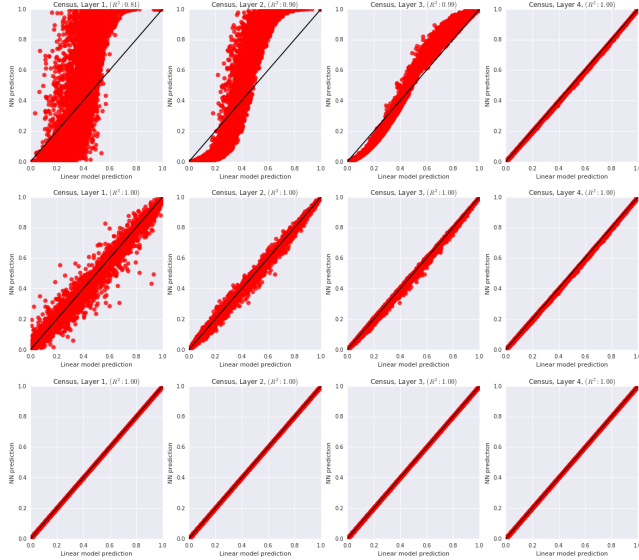


Figure 2: Classifier comparison for *Census*. Training on activation features (row 1), holographic features (row 2), and closed form solution to holographic weights (row 3).

	Pima		Census		MNIST-3v5	
Layer	Activ	Holo	Activ	Holo	Activ	Holo
1	.529	.462	.477	.331	.375	.133
2	.536	.455	.494	.339	.249	.085
3	.474	.471	.380	.345	.155	.082

Table 3: Average negative log-likelihood on third partition for linear classifier trained on second partition (after being extracted from first partition).

which they were not inferred. In this case, we generated both holographic and activation features from a neural network that was trained on the first data partition, as above. Then we trained the linear models defined over these extracted feature sets on the second data partition, and assessed their negative log-likelihood on the third data partition. The results are given in Table 3. For comparison, the generalization performance of the neural network is Pima: 0.421, Census: 0.360, and MNIST-3v5: 0.086 (smaller is better). It is clear that the holographic features support better generalization than activation features, while achieving competitive generalization to the original neural network on the larger data sets.

10 CONCLUSION

We have introduced the concept of generated (“holographic”) features from a model family, specified by the gradient of the prediction with respect to an appropriate subset of parameters. We have shown that, for any model family, if a model is at a local minimum, then

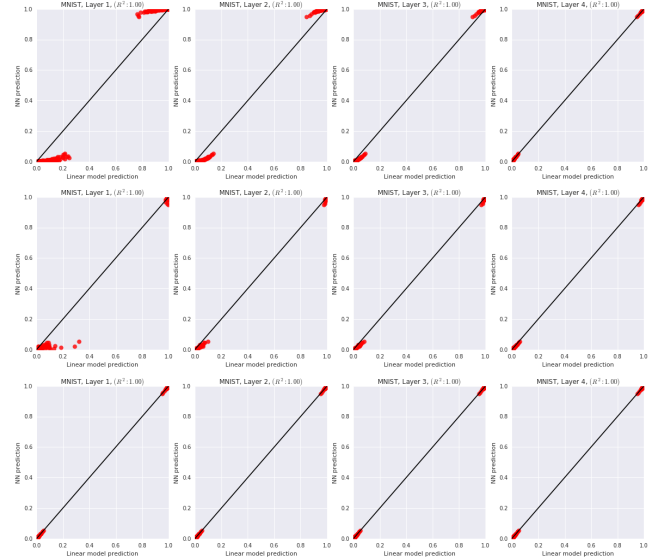


Figure 3: Classifier comparison for *MNIST*. Training on activation features (row 1), holographic features (row 2), and closed form solution to holographic weights (row 3).

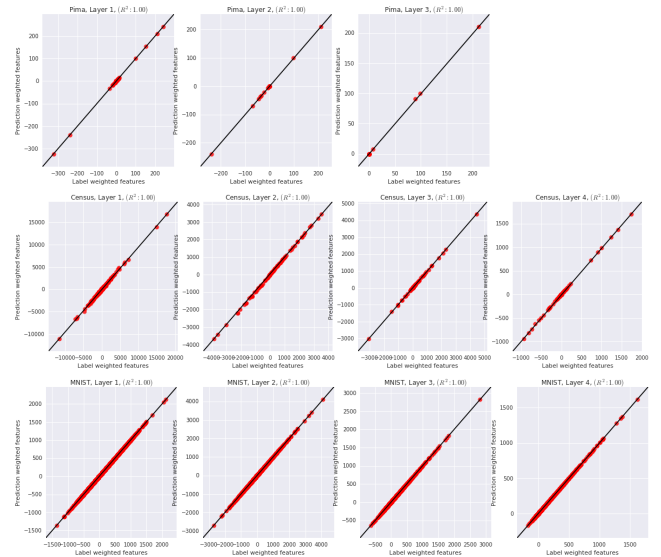


Figure 4: Calibration plots for extracted holographic features from each layer (columns) on each data set (rows).

the model is calibrated with respect to features generated from the parameters with respect to the model. Moreover, we have shown that for many standard feedforward networks, an optimal linear model on the generated features will make the same predictions as the original feedforward network. We show that in practice, building linear models with the generated “holographic” features better replicates the original network than building linear models based on intermediate activations.

References

- [1] G. Alain and Y. Bengio. Understanding intermediate layers using linear classifier probes. *arXiv e-prints*, abs/1610.01644, October 2016.
- [2] L. D. Brown. Fundamentals of statistical exponential families with applications in statistical decision theory. *Lecture Notes-Monograph Series*, 9:i-279, 1986.
- [3] Y. Dauphin, R. Pascanu, Ç. Gülçehre, K. Cho, S. Ganguli, and Y. Bengio. Identifying and attacking the saddle point problem in high-dimensional non-convex optimization. *CoRR*, abs/1406.2572, 2014.
- [4] D. Erhan, Y. Bengio, A. A. Courville, and P. Vincent. Visualizing higher-layer features of a deep network. Technical report, University of Montreal, 2009.
- [5] D. Erhan, P.-A. Manzagol, Y. Bengio, S. Bengio, and P. Vincent. The difficulty of training deep architectures and the effect of unsupervised pre-training. In *AISTATS*, pages 153–160, 2009.
- [6] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [7] A. Graves, A. Mohamed, and G. E. Hinton. Speech recognition with deep recurrent neural networks. In *ICASSP*, 2013.
- [8] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [9] J. Hirschberg and C. Manning. Advances in natural language processing. *Science*, 349(6245):261–266, Jul 2015.
- [10] A. Krizhevsky, I. Sutskever, and G. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012.
- [11] L. D. Kudryavtsev. Homogeneous function. In M. Hazewinkel, editor, *Encyclopedia of Mathematics*. Springer, 2001.
- [12] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [13] A. L. Maas, A. Y. Hannun, and A. Y. Ng. Rectifier nonlinearities improve neural network acoustic models. In *ICML*, 2013.
- [14] T. Mikolov, I. Sutskever, K. Chen, G. Corrado, and J. Dean. Distributed representations of words and phrases and their compositionality. In *NIPS*, 2013.
- [15] A. Mordvintsev, C. Olah, and M. Tyka. Deepdream - a code example for visualizing neural networks. <https://research.googleblog.com/2015/06/inceptionism-going-deeper-into-neural.html>, July 2015.
- [16] J. Nelder and R. Wedderburn. Generalized linear models. *Journal of the Royal Statistical Society Series A (General)*, 135(3):370–384, 1972.
- [17] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015.
- [18] D. Schuurmans and M. Zinkevich. Deep learning games. In *NIPS*, 2016.
- [19] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. LeCun. Overfeat: Integrated recognition, localization and detection using convolutional networks. In *ICLR*, 2014.
- [20] K. Simonyan, A. Vedaldi, and A. Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *CoRR*, abs/1312.6034, 2013.
- [21] M. J. Wainwright and M. I. Jordan. Foundations and trends® in machine learning. 1(1–2):1–305, 2008.
- [22] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey, J. Klingner, A. Shah, M. Johnson, X. Liu, . Kaiser, S. Gouws, Y. Kato, T. Kudo, H. Kazawa, K. Stevens, G. Kurian, N. Patil, W. Wang, C. Young, J. Smith, J. Riesa, A. Rudnick, O. Vinyals, G. Corrado, M. Hughes, and J. Dean. Google’s neural machine translation system: Bridging the gap between human and machine translation. *CoRR*, abs/1609.08144, 2016.
- [23] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson. How transferable are features in deep neural networks? In *NIPS*, pages 3320–3328, 2014.
- [24] M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society B*, 68:49–67, 2006.
- [25] M. Zeiler, G. Taylor, and R. Fergus. Adaptive deconvolutional networks for mid and high level feature learning. In *ICCV*, 2011.
- [26] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional neural networks. In *ECCV*, 2014.

A EXPONENTIAL FAMILIES

This section reviews the standard results on exponential families in the literature [16, 2, 21]. A 1-dimensional standard exponential family [2, page 2] at its core can be represented by a reference measure $(\mathbf{R}, \mathcal{B}, \rho_p)$, where the set of outcomes (whether possible or not: see the definition of Ω_p below) is given by \mathbf{R} , the real numbers; the set of measurable sets is given by \mathcal{B} , the Borel algebra of \mathbf{R} ; and $\rho_p : \mathcal{B} \rightarrow [0, \infty]$ is a positive measure (often not a probability measure). Define $A_p : \mathbf{R} \rightarrow [-\infty, +\infty]$ to be a normalization function where:

$$A_p(\theta) = \ln \left(\int \exp(\theta\omega) d\rho_p(\omega) \right). \quad (14)$$

If A_p is finite for all $\theta \in \mathbf{R}$, then p is defined to be full [2]. If p is full, then for each $\theta \in \mathbf{R}$, there is a distribution in the family of the form:

$$\forall B \in \mathcal{B}, \Pr[x \in B | \theta] = \int_B \exp(\theta\omega - A_p(\theta)) \rho_p(d\omega), \quad (15)$$

Full exponential families are also “regular” [2, page 2]. Define Ω_p to be the support of ρ_p , the minimal closed set $\mathcal{W} \subseteq \mathbf{R}$ such that $\rho_p(\mathbf{R} - \mathcal{W}) = 0$, i.e. the possible outcomes. If p is full and $|\Omega_p| > 1$, it is “minimal” [2, page 2].⁷ A full, minimal, 1-dimensional standard exponential family p has a strictly convex A_p [2, Theorem 1.13] that is infinitely differentiable everywhere [2, Theorem 2.2],

The likelihood of ω given θ is $\exp(\theta\omega - A_p(\theta))$. The negative log likelihood of ω given θ is $\ell_p(\omega, \theta) = A_p(\theta) - \theta\omega$. Notice that if A_p is strictly convex, then ℓ_p is strictly convex in its second argument. The mean $\mu_p : \mathbf{R} \rightarrow \mathbf{R}$ is:

$$\mu_p(\theta) = \int \exp(\theta\omega - A_p(\theta)) \rho_p(d\omega). \quad (16)$$

It is useful to define $\lambda_p : \mathbf{R} \rightarrow \mathbf{R}$ to be:

$$\lambda_p(\theta) = \int \exp(\theta\omega) \rho_p(d\omega). \quad (17)$$

Now given these standard definitions, we can prove Lemma 1:

⁷The definition of minimal in [2, page 2] states that p is minimal if the dimension of the convex hull of the support equals the dimension of the set of parameters where A_p is finite, but since we are dealing with full, 1-dimensional, standard exponential families, that complexity is unnecessary, as the dimension of the set of parameters where A_p is finite is always 1, and the dimension of the convex hull of the support is zero if the support has 1 point, and 1 if the support has two or more points.

Lemma 1 *The Bernoulli family, Gaussian family (fixed variance), and Poisson family are full, minimal, 1-dimensional standard exponential families. For a full, minimal, 1-dimensional standard exponential family p :*

1. for every $\theta \in \mathbf{R}$, $\omega \in \Omega_p$, $\ell_p(\omega, \theta)$ is differentiable with respect to θ and $\frac{\partial \ell_p(\omega, \theta)}{\partial \theta} = \mu_p(\theta) - \omega$; and
2. ℓ_p is **strictly convex in its second argument**: for every $\theta, \theta' \in \mathbf{R}$, if $\theta \neq \theta'$, then for all $\lambda \in (0, 1)$ we have $\ell_p(\omega, \lambda\theta + (1 - \lambda)\theta') < \lambda\ell_p(\omega, \theta) + (1 - \lambda)\ell_p(\omega, \theta')$.

Proof: As we stated before, A_p is strictly convex if p is a full, minimal, 1 dimensional standard exponential family, and this implies that ℓ_p is strictly convex in its second argument. If p is a standard exponential family that is full and minimal, then λ_p is infinitely differentiable everywhere [2, Theorem 2.2], and by [2, page 34]:

$$\lambda'_p(\theta) = \int \omega \exp(\theta\omega) \rho_p(d\omega) \quad (18)$$

We then normalize:

$$\frac{\lambda'_p(\theta)}{\lambda_p(\theta)} = \frac{\int \omega \exp(\theta\omega) \rho_p(d\omega)}{\lambda_p(\theta)} \quad (19)$$

$$= \frac{\int \omega \exp(\theta\omega) \rho_p(d\omega)}{\exp(A_p(\theta))} \quad (20)$$

$$= \int \omega \exp(\theta\omega - A_p(\theta)) \rho_p(d\omega) \quad (21)$$

$$= \mu_p(\theta). \quad (22)$$

So, for $\ell_p(\omega, \theta)$:

$$\frac{\partial \ell_p(\omega, \theta)}{\partial \theta} = A'_p(\theta) - \omega \quad (23)$$

$$= \frac{\lambda'_p(\theta)}{\lambda_p(\theta)} - \omega \quad (24)$$

$$= \mu_p(\theta) - \omega. \quad (25)$$

Equation 25 is [16, Equation 3].

We now just have to show that $|\Omega_p| > 1$ and A_p is finite everywhere for the families mentioned. It is natural to think of the definition of ρ_p in terms of the possible outcomes, but Ω_p is defined in terms of ρ_p . So, instead we define \mathcal{W} as a set (that will turn out to be the possible outcomes), define ρ_p in terms of \mathcal{W} , and then show $\Omega_p = \mathcal{W}$. Note that if \mathcal{W} is finite, to prove $\Omega_p = \mathcal{W}$, we need only prove that for all $\omega \in \mathcal{W}$, $\rho_p(\omega) > 0$, and $\rho_p(\mathbf{R} - \mathcal{W}) = 0$. If $\mathcal{W} = \mathbf{R}$, then we need to show any finite, nonempty, open interval has positive measure to prove $\Omega_p = \mathcal{W}$ (this is because $\mathbf{R} - \Omega_p$ is open, and if there is some $\omega \in \mathbf{R} - \Omega_p$, then there must be a neighborhood N of ω (a finite, nonempty, open interval) in $\mathbf{R} - \Omega_p$ where $\rho_p(N) = 0$).

1. For the Bernoulli family of distributions, $\mathcal{W} = \{0, 1\}$ and $\rho_p(B) = |B \cap \mathcal{W}|$. Thus, $\rho_p(\{0\}) = \rho_p(\{1\}) = 1$, and $\rho_p(\mathbf{R} - \mathcal{W}) = |(\mathbf{R} - \mathcal{W}) \cap \mathcal{W}| = 0$, so $\Omega_p = \mathcal{W}$. Also, $A_p(\theta) = -\ln(1 + \exp(\theta))$, $\Pr[\omega|\theta] = \frac{\exp(\theta\omega)}{1 + \exp(\theta)}$, and $\ell_p(\omega, \theta) = -\theta\omega + \ln(1 + \exp(\theta))$, and $\mu_p(\theta) = \frac{1}{1 + \exp(-\theta)}$. Since $A_p(\theta)$ is finite everywhere, the Bernoulli family is full, and since $|\Omega_p| > 1$, the Bernoulli family is minimal.
2. For the Gaussian family of distributions with fixed mean $\sigma^2 = 1$ we have $\mathcal{W} = \mathbf{R}$, but we also need to define a particular ρ_p . Specifically, define some $h(\omega) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{\omega^2}{2})$. For any Borel measurable set, define $\rho_p(B) = \int_B h(\omega) d\omega$, where the right hand side is the standard Lebesgue integral. Note ρ_p itself is a Gaussian with mean zero and variance one. Since $h(\omega) > 0$ and $h(\omega)$ is continuous, on any finite, closed interval $[a, b]$ it has a minimum $m > 0$, and therefore on any finite, nonempty, open interval (a, b) , $\rho_p((a, b)) \geq (b - a)m$, so $\Omega_p = \mathcal{W}$. Thus, $A_p(\theta) = \frac{\theta^2}{2}$, $\Pr[\omega \in B|\theta] = \frac{1}{\sqrt{2\pi}} \int_B \exp(-\frac{(\omega - \theta)^2}{2}) d\omega$, $\mu_p(\theta) = \theta$, the likelihood⁸ is $\exp(\theta\omega - \frac{\theta^2}{2})$, and $\ell_p(\theta, \omega) = \frac{\theta^2}{2} - \theta\omega = \frac{1}{2}(\theta - \omega)^2 + \frac{\omega^2}{2}$. Notice that this loss is off by $\frac{\omega^2}{2}$ from squared loss, and this term is independent of θ . Finally, $|\Omega_p| > 1$, A_p is finite everywhere, implying ρ_p is a full, minimal, standard, 1-dimensional family.
3. For the Poisson family of distributions $\mathcal{W} = \{0, 1, 2, \dots\}$. Define $\rho_p(B) = \sum_{\omega \in \mathcal{W} \cap B} \frac{1}{\omega!}$, so $A_p(\theta) = \exp(\theta)$. Moreover, for any non-negative integer ω , $\rho_p(\{\omega\}) = \frac{1}{\omega!}$, and $\rho_p(\mathbf{R} - \mathcal{W})$ is the sum over an empty set, and therefore zero. Finally, $\Pr[\omega|\theta] = \frac{1}{\omega!} \exp(\theta\omega - \exp(\theta))$, and the likelihood⁹ is $\exp(\theta\omega - \exp(\theta))$. $\mu_p(\theta) = \exp(\theta)$ and $\ell_p(\omega, \theta) = \exp(\theta) - \theta\omega$. Again, $|\Omega_p| > 1$, and A_p is finite everywhere.

■

B RELATION TO TRADITIONAL LAYERED NEURAL NETWORKS

A more conventional way to represent a network involves writing layers, by interleaving fixed activation functions with learned affine functions. We imagine a sequence of integers $n_0 \dots n_k$, representing the number of nodes in

⁸Note the distinction here between the conventional density defined with respect to the Lebesgue measure, and this likelihood defined with respect to ρ . However, since we are tuning θ and ω is given, this is simply a constant in ℓ_p .

⁹As before, there is a slight distinction between the conventional probability mass and the likelihood as defined here.

each layer, with 0 being the input layer (with $X = \mathbf{R}^{n_0}$), and k being the output layer (with $n_k = 1$). We choose an activation function $a_i : \mathbf{R} \rightarrow \mathbf{R}$ for each layer $\{0, \dots, k - 1\}$, and define $A^i : \mathbf{R}^{n_i} \rightarrow \mathbf{R}^{n_i}$ such that $(A^i(v))_j = a_i(v_j)$ for all $v \in \mathbf{R}^{n_i}$, for all $j \in \{1 \dots n_i\}$. Often a_0 is the identity, and $a_1 \dots a_{k-1}$ are *relu*, sigmoid, or some other standard function.

Then, in-between these layers, we have a matrix $W^i \in \mathbf{R}^{n_i \times n_{i-1}}$ and a vector $w^i \in \mathbf{R}^{n_i}$. We define $h^0 : \mathbf{R}^{n_0} \rightarrow \mathbf{R}^{n_0}$ to be the identity. We can then define $h^i : \mathbf{R}^{n_0} \rightarrow \mathbf{R}^{n_i}$ for $i \in \{1 \dots k\}$ recursively such that $h^i(x) = W^i A^{i-1}(h^{i-1}(x)) + w^i$. The model in this form is $M_{\mathcal{W}, w}(x) = h^k(x)$. As we will show, these can be easily modeled in our graph representation. However, we will see that learning affine (as opposed to linear) functions show that the conventional concept of “layer” is not as crisp and neat as one would like.

B.1 REPRESENTATION AS A GRAPH

First of all, this kind of model can be represented in the graph network that we describe in Section 5. Most oddities of the representation come from the bias features.

1. For each $i \in \{1 \dots k\}$, define $V_i = \{v_{i,1} \dots v_{i,n_i}\}$, and let v_c be a special vertex (which will be a special input node that will always equal 1).
2. Define $V = \{v_c\} \cup \bigcup_{i=0}^k V_i$.
3. Define $I = V_0 \cup \{v_c\}$.
4. Define $o^* = v_{k,1}$,
5. For $i \in \{0, \dots, k - 1\}$, for all $j \in \{1 \dots n_i\}$, let $a_{v_{i,j}} = a^i$.
6. Define a_{o^*} and a_{v_c} to be the identity.
7. For all $x \in X$, for all $j \in 1 \dots n_0$, let $g_{v_0,j}(x) = x_j$.
8. For all $x \in X$, $g_{v_c}(x) = 1$.
9. The edges associated with W^i are $E_i = V_{i-1} \times V_i$.
10. The edges associated with w^i are $E_i^c = \{v_c\} \times V_i$.
11. Define $E = \bigcup_{i=1}^k E_i \cup E_i^c$.

$D = (V, E, I, \{g_i\}_{i \in I}, o^*, \{a_v\}_{v \in V})$ is a neural network, equivalent to the layered form we described in the previous section.

In the next section, we will show how the edges in E_i^c (associated with the bias parameters) play an unusual role in our work.

B.2 AN ORDERED CUT IN A LAYERED NEURAL NETWORK

We introduced ordered cuts and cut sets in Section 6. The most obvious cut would be $B_i = \{v_c\} \cup \bigcup_{j=0}^{i-1} V_j$ and

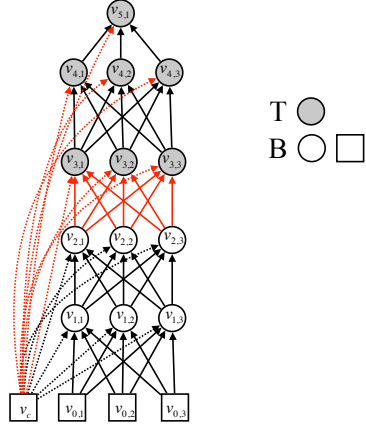


Figure 5: A traditional neural network represented as a graph, with 4 fully connected hidden layers, and a bias parameter for each node. The input nodes are squares, and the internal nodes are circles, with the output node on top. The ordered cut is indicated with the white nodes and the black nodes, and the edges in the cut set are red. Note that while this ordered cut naturally separates the second and third hidden layers, bias parameters from the third, fourth, and fifth layer are in the cut set.

$T_i = \bigcup_{j=i}^k V_j$, and these are the ones we use in the experiments.

Suppose $k = 5$, and we consider the cut B_2, T_2 , then the cut set E' is the set of edges with one endpoint in B_2 and one endpoint in T_2 (see Figure 5). Obviously, E_2 is a subset of the cut set, as is E_2^c . Less obviously, E_3^c, E_4^c , and E_5^c are also subsets of the cut set.¹⁰ However, upon reflection this makes sense: the proof in Appendix F relies on the network after the cut set to be a homogeneous function, and affine functions with nonzero offsets are certainly not homogeneous functions.

This is the reason that we use the graph representation to introduce our results. Cut sets have very counterintuitive properties in the conventional representation, but once the traditional representations are reduced to a network, they make perfect sense.

¹⁰Keep in mind, while E_5^c is in the cut set, it only contains one parameter, and the generated (holographic) feature is always 1.

B.3 EXPLORING THE NATURE OF GENERATED FEATURES

So, as we consider these conventional networks, it is natural to ask, what does a generated feature look like? How does it relate to a conventional activation feature? Let us break this down into features generated from weight matrices W (or $E_1 \dots E_k$), and features generated from bias vectors w (or $E_1^c \dots E_k^c$). In this section, for simplicity we assume that partial derivatives exist where necessary: see Appendix I for a deeper discussion of differentiability in deep networks.

We consider the bias features first for some layer $i \in \{1 \dots k-1\}$. For any $j \geq i$, we can recursively define $h^{j,i} : \mathbf{R}^{n_i} \rightarrow \mathbf{R}^{n_j}$ such that $h^{i,i}(v) = v$ and for all $j \in \{i+1 \dots k\}$, $h^{j,i}(v) = W^j A^{j-1}(h^{j-1,i}(v)) + w^j$. This is an accumulation of the transforms after the input of layer i , and for all $x \in X$:

$$M_{W,w}(x) = h^{k,i}(h^i(x)) \quad (26)$$

$$M_{W,w}(x) = h^{k,i}(W^i A^{i-1}(h^{i-1}(x)) + w^i). \quad (27)$$

For some $q \in \{1 \dots n_i\}$, if we define $f_q^i : X \rightarrow \mathbf{R}$ to be the feature generated from $(v_c, v_{i,q})$, then we can write:

$$\begin{aligned} f_q^i(x) &= \frac{\partial M_{W,w}(x)}{\partial w_q^i} \\ &= \left. \frac{\partial h^{k,i}(v)}{\partial v_q} \right|_{v=W^i A^{i-1}(h^{i-1}(x)) + w^i}. \end{aligned} \quad (28) \quad (29)$$

Notice that this is the partial derivative of the prediction with respect to the input of node $v_{i,q}$.

Next, we consider features generated from a weight matrix. Consider some $i \in \{1 \dots k\}$, some $p \in \{1 \dots n_{i-1}\}$ and some $q \in \{1 \dots n_i\}$. Then $(v_{i-1,p}, v_{i,q}) \in E_i$ is the edge related to parameter $W_{q,p}^i$. Define $F_{q,p}^i : X \rightarrow \mathbf{R}$ to be the feature generated from $(v_{i-1,p}, v_{i,q})$. Using $h^{k,i}$ again:

$$F_{q,p}^i(x) = \frac{\partial M_{W,w}(x)}{\partial W_{q,p}^i} \quad (30)$$

$$= \left. \frac{\partial h^{k,i}(v)}{\partial v_q} \right|_{v=W^i A^{i-1}(h^{i-1}(x)) + w^i} A_p^{i-1}(h^{i-1}(x)) \quad (31)$$

$$= f_q^i(x) a^{i-1}(h_p^{i-1}(x)). \quad (32)$$

So, the generated feature from $(v_{i-1,p}, v_{i,q})$ is the activation feature of $v_{i-1,p}$ times the generated feature of $(v_c, v_{i,q})$.

In summary, the generated features of conventional layers have a nice form, and one can take a conventional layered network and translate it into a graph. However,

from a mathematical perspective, it is much easier to reason about graphs, because of the clean concept of cuts and cut sets. Moreover, the bias features work in highly counterintuitive ways, and presenting a theory without them tells an incomplete story.

C CONVEXITY

In this section, we will state some known results about convexity.

Fact 27 *Given a supervised learning problem $\mathcal{P} = (p, X, \{x_1 \dots x_m\}, \{y_1 \dots y_m\})$, given two models $M : X \rightarrow \mathbf{R}$ and $M' : X \rightarrow \mathbf{R}$ that are equal on the training data, if M is calibrated on a feature $f : X \rightarrow \mathbf{R}$, then M' is calibrated on f .*

One can think about the model as mapping a matrix of inputs to a vector of predictions, which is then composed with a loss function that maps a vector of predictions to a single loss. Next, we show when this second mapping will be (strictly) convex.

Lemma 28 *Given convex sets $C_1 \dots C_m \subseteq \mathbf{R}$, for each $i \in \{1 \dots m\}$ a function $L_i : C_i \rightarrow \mathbf{R}$, then if $C = \times_{i=1}^m C_i$, and there is a function $L : C \rightarrow \mathbf{R}$ such that for all $x \in C$:*

$$L(x) = \sum_{i=1}^m L_i(x_i), \quad (33)$$

then:

1. if for all i , L_i is convex, then L is convex.
2. if for all i , L_i is strictly convex, then L is strictly convex.

Proof: Consider $x, y \in C$, and $\lambda \in [0, 1]$. Without loss of generality, assume $x \neq y$, and $\lambda \in (0, 1)$. By convexity, for all i , $\lambda L_i(x_i) + (1 - \lambda)L_i(y_i) \geq L_i(\lambda x_i + (1 - \lambda)y_i)$, so:

$$\begin{aligned} L(\lambda x + (1 - \lambda)y) &= \sum_{i=1}^m L_i(\lambda x_i + (1 - \lambda)y_i) \end{aligned} \quad (34)$$

$$\leq \sum_{i=1}^m (\lambda L_i(x_i) + (1 - \lambda)L_i(y_i)) \quad (35)$$

$$\leq \lambda \sum_{i=1}^m L_i(x_i) + (1 - \lambda) \sum_{i=1}^m L_i(y_i) \quad (36)$$

$$\leq \lambda L(x) + (1 - \lambda)L(y). \quad (37)$$

To prove the result for strong convexity, we need to be a little more careful. Since $x \neq y$, there exists a $j \in$

$\{1 \dots m\}$ where $x_j \neq y_j$. So $L_j(\lambda x_j + (1 - \lambda)y_j) < \lambda L_j(x_j) + (1 - \lambda)L_j(y_j)$. Thus:

$$\begin{aligned} L(\lambda x + (1 - \lambda)y) &= L_j(\lambda x_j + (1 - \lambda)y_j) \\ &+ \sum_{i \neq j} L_i(\lambda x_i + (1 - \lambda)y_i) \end{aligned} \quad (38)$$

$$\begin{aligned} &< \lambda L_j(x_j) + (1 - \lambda)L_j(y_j) \\ &+ \sum_{i \neq j} L_i(\lambda x_i + (1 - \lambda)y_i) \end{aligned} \quad (39)$$

$$\begin{aligned} &< \lambda L_j(x_j) + (1 - \lambda)L_j(y_j) \\ &+ \sum_{i \neq j} (\lambda L_i(x_i) + (1 - \lambda)L_i(y_i)) \end{aligned} \quad (40)$$

$$< \sum_{i=1}^m (\lambda L_i(x_i) + (1 - \lambda)L_i(y_i)) \quad (41)$$

$$< \lambda \sum_{i=1}^m L_i(x_i) + (1 - \lambda) \sum_{i=1}^m L_i(y_i) \quad (42)$$

$$< \lambda L(x) + (1 - \lambda)L(y). \quad (43)$$

■

D PROOFS OF CALIBRATION ON GENERATED FEATURES

Next we show that if the partial derivative of the loss with respect to a parameter is zero, then the feature generated by that parameter is calibrated.

Theorem 7 *For a problem \mathcal{P} and model family $\mathcal{M} = \{M_w\}_{w \in \mathbf{R}^S}$, given a $w \in \mathbf{R}^S$ and $s \in S$ such that f_s is the feature generated from parameter s of model M_w : if $\{s\}$ is total on the training data given M_w and $\frac{\partial L_{\mathcal{P}}(M_w)}{\partial w_s} = 0$, then M_w is calibrated with respect to f_s .*

Proof: Define $p, m, X, x_1 \dots x_m, y_1 \dots y_m$ such that $\mathcal{P} = (p, X, \{x_1 \dots x_m\}, \{y_1 \dots y_m\})$. We begin with the partial derivative of $L_{\mathcal{P}}$. Note that from Lemma 1, ℓ_p is partially differentiable with respect to its second argument, and since for any $i \in \{1 \dots m\}$ $\{s\}$ is total for x_i given M_w , then $M_w(x_i)$ is partially

differentiable with respect to w_s . Therefore:

$$0 = \frac{\partial L_{\mathcal{P}}(M_w)}{\partial w_s} \quad (44)$$

$$= \frac{\partial}{\partial w_s} \sum_{i=1}^m \ell_p(y_i, M_w(x_i)) \quad (45)$$

$$= \sum_{i=1}^m \frac{\partial \ell_p(y_i, M_w(x_i))}{\partial w_s} \quad (46)$$

$$= \sum_{i=1}^m \left. \frac{\partial \ell_p(y_i, \hat{y}_i)}{\partial \hat{y}_i} \right|_{\hat{y}_i=M_w(x_i)} \frac{\partial M_w(x_i)}{\partial w_s} \quad (47)$$

$$= \sum_{i=1}^m (\mu_p(M_w(x_i)) - y_i) \frac{\partial M_w(x_i)}{\partial w_s}. \quad (48)$$

The last step is because of Lemma 1. For any $i \in \{1 \dots m\}$, since $M_w(x_i)$ is partially differentiable with respect to w_s , we can use $f_s(x_i) = \frac{\partial^+ M_w(x_i)}{\partial w_s} = \frac{\partial M_w(x_i)}{\partial w_s}$ to get:

$$0 = \sum_{i=1}^m (\mu_p(M_w(x_i)) - y_i) f_s(x_i) \quad (49)$$

$$\sum_{i=1}^m y_i f_s(x_i) = \sum_{i=1}^m \mu_p(M_w(x_i)) f_s(x_i). \quad (50)$$

■

Lemma 9 For finite S , given a set of features $\{f_s\}_{s \in S}$, the family of linear models $\mathcal{L}(\{f_s\}_{s \in S})$, and $N_w \in \mathcal{L}$: the set S is total for all $x \in X$ given N_w , and the feature generated from parameter s by model N_w is f_s .

Proof: First note that, by definition, a linear model is a linear function of w , and therefore a differentiable function of w regardless of the input or w . Thus, for any $x \in X$, for any w , the set S is total. Notice that, for any $x \in X$:

$$\frac{\partial M_w(x)}{\partial w_s} = \frac{\partial}{\partial w_s} \sum_{t \in S} w_t f_t(x) \quad (51)$$

$$= \sum_{t \in S} \frac{\partial}{\partial w_s} (w_t f_t(x)) \quad (52)$$

$$= f_s(x). \quad (53)$$

■

Another known key result about linear models is that any two linear models that minimize loss will produce the same predictions on the training data.

Theorem 29 Given a supervised learning problem $\mathcal{P} = (p, X, \{x_1 \dots x_m\}, \{y_1 \dots y_m\})$, a set of features

$f_1 \dots f_n : X \rightarrow \mathbf{R}$, and a family of linear models $\mathcal{L}(\{f_1 \dots f_n\})$, if two models $M, M' \in \mathcal{L}$ are optimal in \mathcal{L} for \mathcal{P} , then they are equal on the training data.

Proof: Define $L^* : \mathbf{R}^m \rightarrow \mathbf{R}$ such that for all $\hat{y} \in \mathbf{R}^m$:

$$L^*(\hat{y}) = \sum_{i=1}^m \ell_p(y_i, \hat{y}_i). \quad (54)$$

By Lemma 1, each ℓ_p is strictly convex in its second argument. By Lemma 28,¹¹ L^* is strictly convex. We define $P : \mathbf{R}^n \rightarrow \mathbf{R}^m$, such that for all $v \in \mathbf{R}^n$, for all $i \in \{1 \dots m\}$, $P_i(v) = M_v(x_i)$. Therefore, $L^*(P(v)) = L_{\mathcal{P}}(M_v)$. We need to prove $P(w) = P(w')$.

Consider $w'' = \frac{w+w'}{2}$. Since the models are linear, we know that for all $x \in X$, $M_{w''}(x) = \frac{M_w(x) + M_{w'}(x)}{2}$. Thus, $P(w'') = \frac{P(w) + P(w')}{2}$. Assume for the sake of contradiction, $P(w) \neq P(w')$. Since L^* is strictly convex:

$$L^*(P(w'')) < \frac{L^*(P(w)) + L^*(P(w'))}{2}; \quad (55)$$

since $L^*(P(v)) = L_{\mathcal{P}}(M_v)$:

$$L_{\mathcal{P}}(M_{w''}) < \frac{L_{\mathcal{P}}(M_w) + L_{\mathcal{P}}(M_{w'})}{2}; \quad (56)$$

since $L_{\mathcal{P}}(M_w) = L_{\mathcal{P}}(M_{w'})$:

$$L_{\mathcal{P}}(M_{w''}) < L_{\mathcal{P}}(M_w). \quad (57)$$

This contradicts the original hypothesis that both are minima. Thus, $P(w) = P(w')$, which is the same as saying, for all $i \in \{1 \dots m\}$, $M_w(x_i) = M_{w'}(x_i)$. ■

Lemma 10 (c.f. [16, Equation 10]) Given a problem \mathcal{P} , a finite set S and a set of features $\{f_s\}_{s \in S}$: a model N in the family of linear models $\mathcal{L}(\{f_s\}_{s \in S})$ is optimal if and only if N is calibrated with respect to f_s for all $s \in S$.

Proof: Define $p, m, X, x_1 \dots x_m, y_1 \dots y_m$ such that $\mathcal{P} = (p, X, \{x_1 \dots x_m\}, \{y_1 \dots y_m\})$. Define $\{N_w\}_{w \in \mathbf{R}^S} = \mathcal{L}(\{f_s\}_{s \in S})$ where $N_w(x) = \sum_{s \in S} w_s f_s(x)$. Define $w^* \in \mathbf{R}^S$ such that $N_{w^*} = N$. If N (and therefore N_{w^*}) is calibrated, for all $s \in S$:

$$\sum_{i=1}^m \mu_p(N_{w^*}(x_i)) f_s(x_i) = \sum_{i=1}^m y_i f_s(x_i) \quad (58)$$

$$0 = \sum_{i=1}^m (y_i - \mu_p(N_{w^*}(x_i))) f_s(x_i). \quad (59)$$

¹¹Formally, for all $i \in \{1 \dots m\}$, we could define $L_i : \mathbf{R} \rightarrow \mathbf{R}$ such that $L_i(\hat{y}_i) = \ell_p(y_i, \hat{y}_i)$.

By Lemma 1:

$$0 = \sum_{i=1}^m \frac{\partial \ell_p(y_i, \hat{y}_i)}{\partial \hat{y}_i} \Big|_{\hat{y}_i = N_{w^*}(x_i)} f_s(x_i) \quad (60)$$

$$0 = \sum_{i=1}^m \frac{\partial \ell_p(y_i, N_{w^*}(x_i))}{\partial w_s^*} \quad (61)$$

$$0 = \frac{\partial L_{\mathcal{P}}(N_{w^*})}{\partial w_s^*}. \quad (62)$$

If we define $L^* : \mathbf{R}^n \rightarrow \mathbf{R}$ as $L^*(w) = L_{\mathcal{P}}(N_w)$, it is convex, then N_{w^*} (and therefore N) is optimal.

To prove the converse, suppose that there is some $s \in S$ that is not calibrated. Then $\frac{\partial L_{\mathcal{P}}(N_{w^*})}{\partial w_s^*} \neq 0$, implying that there is a model with lower loss. ■

Lemma 11 *Given a problem \mathcal{P} , a finite set S , a subset $S' \subseteq S$, and a set of features $\{f_s\}_{s \in S}$: if a model N is optimal in the family of linear models $\mathcal{L}(\{f_s\}_{s \in S'})$ and N is calibrated with respect to f_s for all $s \in S - S'$, then N is also an optimal model in $\mathcal{L}(\{f_s\}_{s \in S})$.*

Proof: Define $p, m, X, x_1 \dots x_m, y_1 \dots y_m$ such that $\mathcal{P} = (p, X, \{x_1 \dots x_m\}, \{y_1 \dots y_m\})$. Note that any model in $\mathcal{L}(\{f_s\}_{s \in S'})$ (and specifically N) is in $\mathcal{L}(\{f_s\}_{s \in S})$. Since N is optimal in $\mathcal{L}(\{f_s\}_{s \in S'})$, by Lemma 10, it is calibrated with respect to $\{f_s\}_{s \in S'}$. Moreover, by assumption it is calibrated with respect to $\{f_s\}_{s \in S - S'}$, and therefore it is calibrated with respect to $\{f_s\}_{s \in S}$, and again by Lemma 10, it is optimal with respect to $\mathcal{L}(\{f_s\}_{s \in S})$. ■

Lemma 13 *Given a problem \mathcal{P} , a finite set S , a subset $S' \subseteq S$, a model family $\{M_w\}_{w \in \mathbf{R}^S}$, and a $w \in \mathbf{R}^S$ such that $\frac{\partial L_{\mathcal{P}}(M_w)}{\partial w_s} = 0$ for all $s \in S'$: if \mathcal{L}' is the family of linear models associated with S' given M_w , S' is total on the training data given M_w , and M_w is equal on the training data to some $N' \in \mathcal{L}'$, then N' is optimal in \mathcal{L}' .*

Proof: Define $p, m, X, x_1 \dots x_m, y_1 \dots y_m$ such that $\mathcal{P} = (p, X, \{x_1 \dots x_m\}, \{y_1 \dots y_m\})$. For all $s \in S'$, define $f_s : X \rightarrow \mathbf{R}$ to be the feature generated by s given M_w . By Theorem 7, M_w is calibrated with respect to f_s . Moreover, the family of linear models associated with S' given M_w is $\mathcal{L}' = \mathcal{L}(\{f_s\}_{s \in S'})$. Since N' and M_w are equal on the training data, by Fact 27, N' is calibrated with respect to all $\{f_s\}_{s \in S'}$. Thus, by Lemma 10, N' is optimal in \mathcal{L}' . ■

Lemma 14 *Given a problem \mathcal{P} , a finite set S , a subset $S' \subseteq S$, a model family $\{M_w\}_{w \in \mathbf{R}^S}$, and a $w \in \mathbf{R}^S$ such that $\frac{\partial L_{\mathcal{P}}(M_w)}{\partial w_s} = 0$ for all $s \in S$: if \mathcal{L}' is the family of*

linear models associated with S' given M_w , S' is total, M_w is equal on the training data to some $N' \in \mathcal{L}'$, and \mathcal{L}'' is the family of linear models associated with S , then any optimal $N'' \in \mathcal{L}''$ equals M_w on the training data.

Proof: For all $s \in S$, define $f_s : X \rightarrow \mathbf{R}$ to be the feature generated by s given M_w . By Theorem 7, M_w is calibrated with respect to f_s , and by Fact 27, N' is calibrated with respect to f_s . By Lemma 13, N' is optimal in $\mathcal{L}' = \mathcal{L}(\{f_s\}_{s \in S'})$. By Lemma 11, N' is optimal in $\mathcal{L}'' = \mathcal{L}(\{f_s\}_{s \in S})$. Thus, by Theorem 29, any optimal $N'' \in \mathcal{L}''$ is equal on the training data to N' , and transitively to M_w . ■

E EULER'S HOMOGENEOUS FUNCTION THEOREM

This appendix can be skipped, as Lemma 17 is a well known result. However, this section provides a deeper discussion of total differentiability and partial differentiability that can be helpful in understanding the rest of the paper.

Although a variant of Lemma 17 is in [11], we prove the specific variant here.

Lemma 17 (Euler's Homogeneous Function Theorem) (Degree 1 Case) *If a homogeneous function $f : \mathbf{R}^p \rightarrow \mathbf{R}^q$ is differentiable at $x \in \mathbf{R}^p$, then for all $k \in \{1 \dots q\}$,*

$$f_k(x) = \sum_{j=1}^p \frac{\partial f_k(x)}{\partial x_j} x_j. \quad (9)$$

Proof: If $x = 0$, then for any $k \in \{1 \dots q\}$:

$$\sum_{j=1}^p \frac{\partial f_k(x)}{\partial x_j} x_j = \sum_{j=1}^p \frac{\partial f_k(x)}{\partial x_j} \times 0 \quad (63)$$

$$= 0 \quad (64)$$

$$= 0 \times f_k(x) \quad (65)$$

$$= f_k(0x) \quad (66)$$

$$= f_k(x). \quad (67)$$

For the remainder, assume $x \neq 0$. Consider a particular x where $f(x)$ is differentiable, and the derivative is a matrix $G \in \mathbf{R}^{q \times p}$ where:

$$\lim_{h \rightarrow 0} \frac{\|f(x+h) - f(x) - Gh\|}{\|h\|} = 0. \quad (68)$$

Specifically, we can write:

$$\lim_{\epsilon \rightarrow 0} \frac{\|f(x+\epsilon x) - f(x) - G\epsilon x\|}{\|\epsilon x\|} = 0. \quad (69)$$

In the limit, $\epsilon > -1$, so by the homogeneous property:

$$\lim_{\epsilon \rightarrow 0} \frac{\|(1 + \epsilon)f(x) - f(x) - G\epsilon x\|}{\|\epsilon x\|} = 0 \quad (70)$$

$$\lim_{\epsilon \rightarrow 0} \frac{\|\epsilon f(x) - G\epsilon x\|}{\epsilon \|x\|} = 0 \quad (71)$$

$$\lim_{\epsilon \rightarrow 0} \frac{\|f(x) - Gx\|}{\|x\|} = 0 \quad (72)$$

$$\frac{\|f(x) - Gx\|}{\|x\|} = 0 \quad (73)$$

$$\|f(x) - Gx\| = 0 \quad (74)$$

$$f(x) = Gx. \quad (75)$$

Since $G_{k,j} = \frac{\partial f_k(x)}{\partial x_j}$, the result follows. ■

One might wonder, is it possible to prove this for functions that are only partially (and not totally) differentiable? Yes and no: for one or two dimensions partial differentiability is sufficient for Euler's function to hold, but for three or more dimensions, there is no such guarantee, and we can demonstrate this with a counterexample.

Lemma 30 *If a homogeneous function $f : \mathbf{R}^p \rightarrow \mathbf{R}^q$ is partially differentiable at a point $x \in \mathbf{R}^p$, and $p \leq 2$, then for all $k \in \{1 \dots q\}$:*

$$f_k(x) = \sum_{j=1}^p \frac{\partial f_k(x)}{\partial x_j} x_j. \quad (76)$$

Proof: Since partial differentiability implies differentiability when $p = 1$, the proof for $p = 1$ follows directly from Lemma 17. So, we assume $p = 2$. Consider a specific point (x, y) where f is partially differentiable, i.e. $\frac{\partial f(x,y)}{\partial x}$ and $\frac{\partial f(x,y)}{\partial y}$ exist. If $x = 0$ or $y = 0$, then it is basically analogous to the one dimensional case, as the function along the axis can be considered a 1 dimensional homogeneous function that is differentiable at that point.

Next, assume $x > 0$ and $y > 0$, i.e. a point in the positive quadrant as in Figure 6: near the end of the proof, we will show how to reduce a problem in a different quadrant to one in the positive quadrant.

We can define $g : (0, x) \rightarrow (0, \infty)$ and $h : (0, x) \rightarrow (0, \infty)$ such that $g(\epsilon) = \frac{y\epsilon}{x-\epsilon}$ and $h(\epsilon) = \frac{y+g(\epsilon)}{y}$.

Then, for any $\epsilon \in (0, x)$, observe that $h(\epsilon)(x - \epsilon, y) = (x, y + g(\epsilon))$, so $h(\epsilon)f(x - \epsilon, y) = f(x, y + g(\epsilon))$. Since

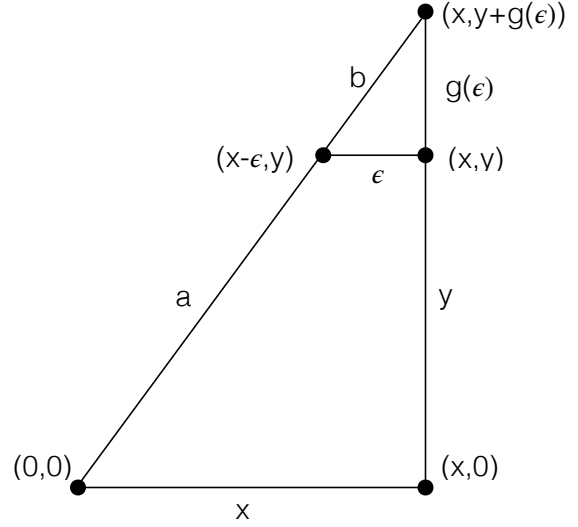


Figure 6: For homogeneous functions with a 2-dimensional domain, the partial derivatives are inextricably linked. Note that the definition of $g(\epsilon)$ can be derived from $\frac{g(\epsilon)+y}{x} = \frac{g(\epsilon)}{\epsilon}$. The function values at $(x - \epsilon, y)$ and $(x, y + g(\epsilon))$ are connected through the homogeneous property.

$\lim_{\epsilon \rightarrow 0^+} g(\epsilon) = 0$, we know that:

$$\frac{\partial f(x, y)}{\partial y} = \lim_{\epsilon \rightarrow 0^+} \frac{f(x, y + g(\epsilon)) - f(x, y)}{g(\epsilon)} \quad (77)$$

$$\frac{\partial f(x, y)}{\partial y} = \lim_{\epsilon \rightarrow 0^+} \frac{f(x - \epsilon, y)h(\epsilon) - f(x, y)}{g(\epsilon)} \quad (78)$$

$$\frac{\partial f(x, y)}{\partial y} = \lim_{\epsilon \rightarrow 0^+} \frac{f(x - \epsilon, y)h(\epsilon) - f(x, y)h(\epsilon)}{g(\epsilon)} \quad (79)$$

$$\begin{aligned} &+ \lim_{\epsilon \rightarrow 0^+} \frac{f(x, y)h(\epsilon) - f(x, y)}{g(\epsilon)} \\ \frac{\partial f(x, y)}{\partial y} &= \lim_{\epsilon \rightarrow 0^+} \frac{f(x - \epsilon, y) - f(x, y)}{\epsilon} \lim_{\epsilon \rightarrow 0^+} \frac{h(\epsilon)\epsilon}{g(\epsilon)} \\ &+ f(x, y) \lim_{\epsilon \rightarrow 0^+} \frac{h(\epsilon) - 1}{g(\epsilon)}. \end{aligned} \quad (80)$$

First, observe that $\lim_{\epsilon \rightarrow 0^+} \frac{f(x - \epsilon, y) - f(x, y)}{\epsilon} = -\frac{\partial f(x, y)}{\partial x}$. Notice that for $\epsilon \in (0, x)$:

$$\frac{h(\epsilon) - 1}{g(\epsilon)} = \frac{\frac{y+g(\epsilon)}{y} - 1}{g(\epsilon)} = \frac{1}{y}. \quad (81)$$

Also, since $\lim_{\epsilon \rightarrow 0^+} h(\epsilon) = 1$, for $\epsilon \in (0, x)$:

$$\lim_{\epsilon \rightarrow 0^+} \frac{h(\epsilon)\epsilon}{g(\epsilon)} = \lim_{\epsilon \rightarrow 0^+} \frac{\epsilon}{g(\epsilon)} \quad (82)$$

$$= \lim_{\epsilon \rightarrow 0^+} \frac{\epsilon(x - \epsilon)}{y\epsilon} \quad (83)$$

$$= \frac{x}{y}. \quad (84)$$

So, from Equation 80:

$$\frac{\partial f(x, y)}{\partial y} = -\frac{\partial f(x, y)}{\partial x} \frac{x}{y} + f(x, y) \frac{1}{y} \quad (85)$$

$$\frac{\partial f(x, y)}{\partial x} x + \frac{\partial f(x, y)}{\partial y} y = f(x, y), \quad (86)$$

which is Euler's function.

Now if $x < 0$ or $y < 0$, we could flip the function on one axis or on both without affecting homogeneity, partial differentiability, or Euler's function. Specifically, note that if we defined $f^* : \mathbf{R}^2 \rightarrow \mathbf{R}$ such that for all $x', y' \in \mathbf{R}$, $f^*(x', y') = f(-x', y')$, then $\frac{\partial f^*(-x, y)}{\partial x} = -\frac{\partial f(x, y)}{\partial x}$, $\frac{\partial f^*(-x, y)}{\partial y} = \frac{\partial f^*(-x, y)}{\partial y}$, and if $f^*(-x, y) = \frac{\partial f^*(-x, y)}{\partial x}(-x) + \frac{\partial f^*(-x, y)}{\partial y}y$ then:

$$f(x, y) = f^*(-x, y) \quad (87)$$

$$= \frac{\partial f^*(-x, y)}{\partial x}(-x) + \frac{\partial f^*(-x, y)}{\partial y}y \quad (88)$$

$$= \left(-\frac{\partial f(x, y)}{\partial x}\right)(-x) + \frac{\partial f(x, y)}{\partial y}y \quad (89)$$

$$= \frac{\partial f(x, y)}{\partial x}x + \frac{\partial f(x, y)}{\partial y}y. \quad (90)$$

Similarly for flipping on the x axis. ■

However, if we add a third dimension, things get incredibly complex. Consider the following function:

$$f(x, y, z) = \begin{cases} 7z - x & \text{if } x + y - 2z \geq 0 \text{ and } x - y < 0 \\ 7z - y & \text{if } x - y \geq 0 \text{ and } x + y - 2z > 0 \\ x + 5z & \text{if } x + y - 2z \leq 0 \text{ and } x - z > 0 \\ 7z - x & \text{if } x - z \leq 0 \text{ and } x - y > 0 \\ 7z - y & \text{if } x - y \leq 0 \text{ and } y - z < 0 \\ y + 5z & \text{if } y - z \geq 0 \text{ and } x + y - 2z < 0 \\ y + 5z & \text{if } x = y = z \end{cases} \quad (91)$$

So, we wish to establish five things:

1. f is well-defined.
2. f is continuous.
3. f is homogeneous.
4. f has partial derivatives at $(1, 1, 1)$.

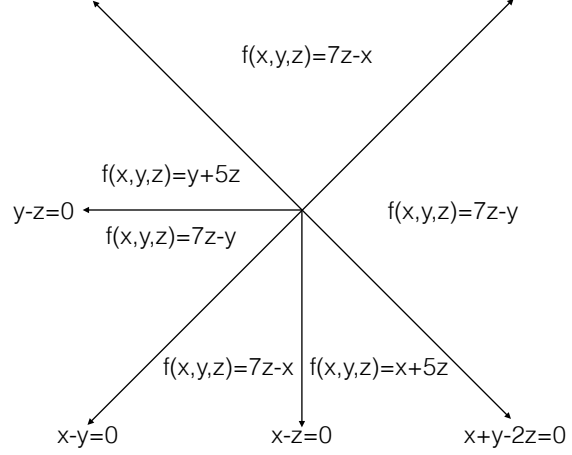


Figure 7: A visual representation of the eight regions of the piecewise linear homogeneous function $f : \mathbf{R}^3 \rightarrow \mathbf{R}$. We are looking along the axis $x = y = z$.

5. Euler's formula does not hold at $(1, 1, 1)$.

To confirm f is well-defined, we must confirm that it is defined everywhere. If one considers Figure 7, one will notice that we have defined each region, starting from the top region and going clockwise, and then defined the center. Each region contains its counterclockwise edge, but not its clockwise edge, and none contain the center. Thus, each point is in exactly one of the cases.

Next, we must establish continuity. First, we establish it at the center. If $x = y = z$, then clearly $7z - x = 7z - y$, and $y + 5z = x + 5z$. If $x = y = z$, then $x + y - 2z = 0$, so $7z - x = 7z - x - (x + y - 2z) = x + 5z$.

Thus, all four linear functions equal each other along the line $x = y = z$. Simple algebra shows us, along the plane $x - y = 0$, $7z - x = 7z - y$, because $7z - x = 7z - x + (x - y) = 7z - y$. Similarly, we can look at each boundary, and prove equality.

Notice that homogeneity is pretty straightforward. Since each constraint is in itself a linear inequality with no constants, given any (x, y, z) , $(\lambda x, \lambda y, \lambda z)$ has the exact same constraints hold. Since within each region, the function is linear, then it is also homogeneous.

At the point $(1, 1, 1)$, notice that:

1. if you move in either direction along the x axis, $f(x, y, z) = 7z - y$, so $\left. \frac{\partial f(x, 1, 1)}{\partial x} \right|_{x=1} = 0$,
2. if you move in either direction along the y axis, $f(x, y, z) = 7z - x$, so $\left. \frac{\partial f(1, y, 1)}{\partial y} \right|_{y=1} = 0$,

3. and if you move in either direction along the z axis, you will stay where $x - y = 0$, so $f(x, y, z) = 7z - x$, so $\left. \frac{\partial f(1,1,z)}{\partial z} \right|_{z=1} = 7$.

Note that $f(1, 1, 1) = 6$. However:

$$\begin{aligned} \left. \frac{\partial f(x, 1, 1)}{\partial x} \right|_{x=1} \times 1 + \left. \frac{\partial f(1, y, 1)}{\partial y} \right|_{y=1} \times 1 + \left. \frac{\partial f(1, 1, z)}{\partial z} \right|_{z=1} \times 1 \\ = 0 \times 1 + 0 \times 1 + 7 \times 1 = 7. \end{aligned} \quad (92)$$

Thus, since $6 \neq 7$, this function, while homogeneous and continuous, has a point where it is partially differentiable, but Euler's formula does not hold.

With further effort, we can also show that the function can be constructed using *relu* gates. Thus, in a very crucial way, Euler's function requires total differentiability, justifying the large role this concept of total features has in our theoretical analysis. However, as our empirical analysis shows, we can pretty much ignore this concern in practice. In Appendix I, we give a sufficient condition (based on a complex proof) to determine if the model is total on an input.

F PROOFS FOR HOMOGENEOUS FUNCTIONS

Before continuing, it is important to note that when defining the homogeneous property of a model family, we directly appealed to a property of differentiability due to the technical issues described in higher dimensional homogeneous functions in Appendix E.

We now prove Theorem 20:

Theorem 20 *If a model family \mathcal{M} is homogeneous on the parameter set S' , $M \in \mathcal{M}$, \mathcal{L}' is the family of linear models associated with S' given M , and S' is total on the training data given M , then there exists $N \in \mathcal{L}'$ such that M and N are equivalent on the training data.*

Proof: Define $\mathcal{M} = \{M'_w\}_{w \in \mathbf{R}^S}$, and $w \in \mathbf{R}^S$ such that $M'_w = M$. Given some x in the training data, since \mathcal{M} is homogeneous with respect to S' , and S' is total on x given M'_w , then by Lemma 17:

$$M'_w(x) = \sum_{s \in S'} \frac{\partial M'_w(x)}{\partial w_s} w_s. \quad (93)$$

Since $f_s(x) = \frac{\partial^+ M'_w(x)}{\partial w_s}$ is the feature generated by s given M , if x is in the training data, then $\{s\} \subseteq S'$ is total on x given M'_w , and $f_s(x) = \frac{\partial M'_w(x)}{\partial w_s}$, so:

$$M'_w(x) = \sum_{s \in S'} w_s f_s(x). \quad (94)$$

Note $\mathcal{L}' = \mathcal{L}(\{f_s\}_{s \in S'})$, and if $\{N'_w\}_{w \in \mathbf{R}^{S'}} = \mathcal{L}'$, then for the $w' \in \mathbf{R}^{S'}$ where $w'_s = w_s$ for all $s \in S'$:

$$N'_{w'}(x) = \sum_{s \in S'} w'_s f_s(x) \quad (95)$$

$$= \sum_{s \in S'} w_s f_s(x) \quad (96)$$

$$= M'_w(x) = M(x), \quad (97)$$

so $M(x) = N'_{w'}(x)$ on the training data, and $N'_{w'}$ is in the family of linear models associated with S' given M . ■

We now prove Theorem 23:

Theorem 23 *Given a problem \mathcal{P} and a family of feedforward network models $\mathcal{D}(V, E, I, \{g_i\}_{i \in I}, \sigma^*, A) = \{M_w\}_{w \in \mathbf{R}^E}$, where all $a \in A$ are homogeneous: if B and T is an ordered cut of the feedforward network, such that E' is the cut set, then:*

1. $\{M_w\}_{w \in \mathbf{R}^E}$ is homogeneous on E' ;
2. for some $w \in \mathbf{R}^E$, if E' is total on the training data given M_w , \mathcal{L}' is the family of linear models associated with E' given M_w , and $\frac{\partial L_{\mathcal{P}}(M_w)}{\partial w_e} = 0$ for all $e \in E'$, then M_w is equal on the training data to any optimal model in \mathcal{L}' ;
3. for some $w \in \mathbf{R}^E$, if \mathcal{L}' is the family of linear models associated with E given M_w , if E is total on the training data given M_w , and $\frac{\partial L_{\mathcal{P}}(M_w)}{\partial w_e} = 0$ for all $e \in E$, then M_w is equal on the training data to any optimal model in \mathcal{L}' .

Proof: We need only prove that the model family is homogeneous: the other two results follow from Theorem 20, Lemma 13, and Lemma 14. We need to construct a function from the inputs passing through the cut set to the second part of the function. First, define $E'' = E - E'$, and given w , define $w'' \in \mathbf{R}^{E''}$ such that for all $e \in E''$, $w''_e = w_e$. Thus, for all $t \in T$, define $d_{t,w''} : \mathbf{R}^{E'} \rightarrow \mathbf{R}$ recursively (using the partial ordering of the vertices in the network) such that for all $z \in \mathbf{R}^{E'}$:

$$\begin{aligned} d_{t,w''}(z) = a_t \left(\sum_{u:(u,t) \in E'} z_{(u,t)} \right. \\ \left. + \sum_{u:(u,t) \in E''} d_{u,w''}(z) w''_{(u,t)} \right). \end{aligned} \quad (98)$$

We subscript this by w'' to indicate that $d_{t,w''}$ is independent of the weights in E' . An important point to note is that for any $x \in X$, for any $b \in B$, $c_{b,w}(x)$ does not depend upon the weight of any edge in E' . For any

$(b, t) \in E'$, we can write $g_{w''}(x)_{(b,t)} = c_{b,w}(x)$ to formalize this idea. Define $w' \in \mathbf{R}^{E'}$ such that for all $e \in \mathbf{R}^{E'}$, $w'_e = w_e$. Finally, for any $z, z' \in \mathbf{R}^{E'}$ we denote the **Hadamard (entrywise) product** as $(z \circ z')_e = z_e z'_e$. We can prove recursively, for all $t \in T$:

$$c_{t,w}(x) = d_{t,w''}(w' \circ g_{w''}(x)). \quad (99)$$

Importantly, this means:

$$M_w(x) = c_{o^*,w}(x) \quad (100)$$

$$= d_{t,w''}(w' \circ g_{w''}(x)). \quad (101)$$

We can use this formulation to prove $M_w(x)$ is homogeneous in w' . First, we can prove recursively that $d_{t,w''}(z)$ is a homogeneous function of z : since all $a \in A$ are homogeneous, it is a homogeneous function of the sum of a set of projections (and all projections are homogeneous) and a set of homogeneous functions multiplied by a constant. To prove that the overall function is homogeneous, we introduce $\lambda \geq 0$, $v \in \mathbf{R}^E$, $v' \in \mathbf{R}^{E'}$, and $v'' \in \mathbf{R}^{E''}$, such that $v'_s = v_s$ for all $s \in E'$, $v''_s = v_s$ for all $s \in E''$, $v'' = w''$, and $v' = \lambda w'$. Note that:

$$M_v(x) = c_{o^*,v}(x) \quad (102)$$

$$= d_{t,v''}(v' \circ g_{v''}(x)) \quad (103)$$

$$= d_{t,w''}((\lambda w') \circ g_{w''}(x)) \quad (104)$$

$$= d_{t,w''}(\lambda(w' \circ g_{w''}(x))). \quad (105)$$

Since $d_{t,w''}$ is homogeneous:

$$M_v(x) = \lambda d_{t,w''}(w' \circ g_{w''}(x)) \quad (106)$$

$$= \lambda M_w(x). \quad (107)$$

We have established that $M_w(x)$ is homogeneous with respect to w' . Thus, from Theorem 20, Lemma 13, and Lemma 14, the other results hold. ■

G PROOFS FOR RESNETS AND CNNS

Recall the main theorem about ResNets and CNNs:

Theorem 24 *For a problem \mathcal{P} and family of RC feedforward network models $\mathcal{RC}(V, E, I, \{g_i\}_{i \in I}, o^*, A, w^*)$ denoted $\{Q_v\}_{v \in \mathbf{R}^n}$, where all $a \in A$ are homogeneous: if $E_1 \dots E_n$ is the partition of the dynamic parameters, and B and T are a well-behaved cut such that there is an $S \subseteq \{1 \dots n\}$ where $E' = \bigcup_{s \in S} E_s$ is the cut set; then*

1. $\{Q_v\}_{v \in \mathbf{R}^n}$ is a homogeneous model family with respect to S ;
2. if for some $v \in \mathbf{R}^n$, $\frac{\partial L_{\mathcal{P}}(Q_v)}{\partial v_s} = 0$ for all $s \in S$, the set S is total on the training data given Q_v , and \mathcal{L}' is the family of linear models associated with S given Q_v , then Q_v is equal on the training data to any optimal N' in \mathcal{L}' ;

3. if for some $v \in \mathbf{R}^n$, $\frac{\partial L_{\mathcal{P}}(Q_v)}{\partial v_s} = 0$ for all $s \in \{1 \dots n\}$, the set $\{1 \dots n\}$ is total on the training data given Q_v , and \mathcal{L}' is the family of linear models associated with $\{1 \dots n\}$ given Q_v , then Q_v is equal on the training data to any optimal N' in \mathcal{L}' .

Proof: We need only prove that the model family is homogeneous: as with Theorem 23, the other two results follow from Theorem 20, Lemma 13, and Lemma 14. We can define $S'' = \{1 \dots n\} - S$, and $v'' \in \mathbf{R}^{S''}$ such that $v''_s = v_s$ for all $s \in S''$. We can define $v' \in \mathbf{R}^S$ such that $v'_s = v_s$ when $s \in S$. We can define $E'' = E - E'$, and $w'' : \mathbf{R}^{S''} \rightarrow \mathbf{R}^{E''}$ such that for all $e \in E''$, $w''_e(v'') = w''_e$ if $e \in E''$, and $w''_e = v''_{\pi(e)}$ otherwise. As before, we will define $d_{t,v''} : \mathbf{R}^{E''} \rightarrow \mathbf{R}$, and we will create a matrix $G^{v''} : X \rightarrow \mathbf{R}^{E' \times S}$. As before, for all $t \in T$, for all $z \in \mathbf{R}^{E'}$, we define $d_{t,v''}$ recursively:

$$d_{t,v''}(z) = a_t \left(\sum_{u:(u,t) \in E'} z(u,t) + \sum_{u:(u,t) \in E''} d_{u,v''}(z) w''_{(u,t)}(v'') \right) \quad (108)$$

Since, for all $b \in B$, $c_{b,w^*(v)}(x)$ does not depend upon parameters in the cut set, we can write $(G^{v''}(x))_{e,s} = 0$ if $\pi(e) \neq s$, and otherwise $(G^{v''}(x))_{(b,t),s} = c_{b,w^*(v)}(x)$. Crucially, neither d nor G is a function of the parameters in S . We can now write the activation energy of a node in the top as a combination of d , G , and v' :

$$c_{t,w^*(v)}(x) = d_{t,v''}(G^{v''}(x)v') \quad (109)$$

Notice that $G^{v''}(x)v'$ is the product of a matrix and a vector, and $(G^{v''}(x)v')_{(b,t)} = v'_{\pi(b,t)} c_{b,w^*(v)}(x)$. Moreover, this relies on the fact that the cut set does not include E^f , as this would make the activation energies on the cut set an affine function of v' instead of a linear one. Considering the output of the model is $c_{t,w^*(v)}(x)$ yields:

$$Q_v(x) = d_{o^*,v''}((G^{v''}(x))v'). \quad (110)$$

Consider any $\lambda > 0$. Define $y \in \mathbf{R}^n$, where for all $s \in S''$, $y_s = v_s$, and for all $s \in S$, $y_s = \lambda v_s$. If we define $y' \in \mathbf{R}^S$ such that $y'_s = y_s$ for all $s \in S$, then we can write:

$$Q_y(x) = d_{o^*,v''}((G^{v''}(x))y') \quad (111)$$

$$Q_y(x) = d_{o^*,v''}((G^{v''}(x))(\lambda v')) \quad (112)$$

$$Q_y(x) = d_{o^*,v''}(\lambda(G^{v''}(x))v') \quad (113)$$

As we argued in the proof of Theorem 23, $d_{o^*,v''}(z)$ is a homogeneous function of z . So:

$$Q_y(x) = \lambda d_{o^*,v''}((G^{v''}(x))v') \quad (114)$$

$$Q_y(x) = \lambda Q_v(x) \quad (115)$$

Thus, the RC class is homogeneous with respect to S . ■

H REGULARIZATION AND RESTRICTIONS

Now we prove Theorem 26.

Theorem 26 *Given a problem \mathcal{P} , a set S , a subset $S' \subseteq S$, a model family $\mathcal{M} = \{M_w\}_{w \in \mathbf{R}^S}$ that is homogeneous with respect to S' , a strictly convex regularization function $R : \mathcal{M} \rightarrow \mathbf{R}$ that is additively separable with respect to S' , and a model M_w where for all $s \in S'$:*

$$\frac{\partial[L_{\mathcal{P}}(M_w) + R(M_w)]}{\partial w_s} = 0, \quad (13)$$

if S' is total on the training data given M_w , $\mathcal{L}' = \{N_v\}_{v \in \mathbf{R}^{S'}}$ is the family of linear models associated with S' given M_w , and a new regularizer $R' : \mathcal{L}' \rightarrow \mathbf{R}$ is defined such that $R'(N_v) = (R^)^{S'}(v)$, then M_w must be equal on the training data to any optimal $N \in \mathcal{L}'$ given the supervised learning problem and regularization R' .*

Proof: Define $v' \in \mathbf{R}^{S'}$ such that $v'_s = w_s$ for all $s \in S'$. For all $s \in S'$, define f_s to be the feature generated by s given M_w . We will prove the result by showing that:

1. For all $s \in S'$, for all x in the training data, $f_s(x) = \frac{\partial M_w(x)}{\partial w_s}$.
2. For all x in the training data, $N_{v'}(x) = M_w(x)$.
3. For all $s \in S'$, $\frac{\partial L_{\mathcal{P}}(N_{v'})}{\partial v'_s} = \frac{\partial L_{\mathcal{P}}(M_w)}{\partial w_s}$.
4. For all $s \in S'$, $\frac{\partial R'(N_{v'})}{\partial v'_s} = \frac{\partial R(M_w)}{\partial w_s}$.
5. Now, we know that $\frac{\partial[L_{\mathcal{P}}(N_{v'}) + R'(N_{v'})]}{\partial v'_s} = \frac{\partial[L_{\mathcal{P}}(M_w) + R(M_w)]}{\partial w_s} = 0$. Then, we use this to prove that $N_{v'}$ is the unique optimal solution in \mathcal{L}' .

We begin by proving item 1. By definition, $f_s(x) = \frac{\partial^+ M_w(x)}{\partial w_s}$. Since S' is total on the training data given M_w , M_w is partially differentiable with respect to w_s , so $f_s(x) = \frac{\partial M_w(x)}{\partial w_s}$.

Next we prove item 2 (c.f. Theorem 20). Note that, for any $x \in X$:

$$N_{v'}(x) = \sum_{s \in S'} v'_s f_s(x). \quad (116)$$

By the definition of v' :

$$N_{v'}(x) = \sum_{s \in S'} w_s f_s(x). \quad (117)$$

By item 1, if x is in the training data:

$$N_{v'}(x) = \sum_{s \in S'} w_s \frac{\partial M_w(x)}{\partial w_s}. \quad (118)$$

Because the model family \mathcal{M} is homogeneous, and S' is total on the training data given M_w :

$$N_{v'}(x) = M_w(x). \quad (119)$$

Next, we prove item 3. Choose an arbitrary $s \in S'$:

$$\frac{\partial L_{\mathcal{P}}(N_{v'})}{\partial v'_s} = \sum_{i=1}^m \frac{\partial L_{\mathcal{P}}(y, \hat{y}'_i)}{\partial \hat{y}'_i} \Big|_{\hat{y}'_i = N_{v'}(x_i)} f_s(x_i). \quad (120)$$

By item 1:

$$\frac{\partial L_{\mathcal{P}}(N_{v'})}{\partial v'_s} = \sum_{i=1}^m \frac{\partial L_{\mathcal{P}}(y, \hat{y}'_i)}{\partial \hat{y}'_i} \Big|_{\hat{y}'_i = N_{v'}(x_i)} \frac{\partial M_w(x_i)}{\partial w_s}. \quad (121)$$

By item 2, $\hat{y}'_i = N_{v'}(x_i) = M_w(x_i)$:

$$\frac{\partial L_{\mathcal{P}}(N_{v'})}{\partial v'_s} = \sum_{i=1}^m \frac{\partial L_{\mathcal{P}}(y, \hat{y}'_i)}{\partial \hat{y}'_i} \Big|_{\hat{y}'_i = M_w(x_i)} \frac{\partial M_w(x_i)}{\partial w_s} \quad (122)$$

$$= \frac{\partial L_{\mathcal{P}}(M_w)}{\partial w_s}. \quad (123)$$

Now, we prove item 4. Choose an arbitrary $s \in S'$. Note that $\frac{\partial R(M_w)}{\partial w_s} = \frac{\partial R^*(w)}{\partial w_s} = \frac{\partial (R^*)^{S'}(\pi^{S \rightarrow S'}(w))}{\partial w_s}$. Since $v' = \pi^{S \rightarrow S'}(w)$, $\frac{\partial (R^*)^{S'}(\pi^{S \rightarrow S'}(w))}{\partial w_s} = \frac{\partial (R^*)^{S'}(v')}{\partial v'_s}$. Also, by the definition of R' , $\frac{\partial (R^*)^{S'}(v')}{\partial v'_s} = \frac{\partial R'(N_{v'})}{\partial v'_s}$. So $\frac{\partial R(M_w)}{\partial w_s} = \frac{\partial R'(N_{v'})}{\partial v'_s}$.

Finally, we must prove item 5. By item 3 and item 4, for any $s \in S'$:

$$\frac{\partial[L_{\mathcal{P}}(N_{v'}) + R'(N_{v'})]}{\partial v'_s} = \frac{\partial[L_{\mathcal{P}}(M_w) + R(M_w)]}{\partial w_s}. \quad (124)$$

Then, by assumption:

$$\frac{\partial[L_{\mathcal{P}}(N_{v'}) + R'(N_{v'})]}{\partial v'_s} = \frac{\partial[L_{\mathcal{P}}(M_w) + R(M_w)]}{\partial w_s} = 0. \quad (125)$$

Now, define a function $g : \mathbf{R}^{S'} \rightarrow \mathbf{R}$ such that for all $v \in \mathbf{R}^{S'}$, $g(v) = L_{\mathcal{P}}(N_v) + R'(N_v)$. Notice that the first part is convex, and the second part is strictly convex, implying g is strictly convex. Notice that $\nabla g(v') = 0$. Thus, v' is a minimum of g . Moreover, since g is strictly convex, it can have no more than one minimum. So $N_{v'}$ is the unique optimal model, and it is equivalent to M_w . ■

I SUFFICIENT CONDITIONS FOR TOTAL FEATURE SETS

In this section, we want to focus on when the concept of totality applies in deep networks. Specifically, are there easy rules of thumb to determine whether a feature set is total?

I.1 TOTALITY AND DIFFERENTIABILITY

While we invented the term "totality", it is very similar to the concept of differentiability. Specifically, consider $g : \mathbf{R}^{S'} \rightarrow \mathbf{R}$, where for all $h \in \mathbf{R}^{S'}$, $g(h) = M_{w+h}(x)$. Then, $M_w(x)$ is total if $g(h)$ is differentiable at zero.

I.2 TOTALITY AND FEEDFORWARD NETWORKS

Let's focus on feedforward networks with *relu* gates, as those are relatively simple and have nice homogeneous properties. Specifically, assume that each input and the output have identity transformations, and the internal nodes are *relu* gates.

Notice that the features will always exist for feedforward networks. If one considers the output of a feedforward network as a function of weights of the edges given a fixed input, the output is a piecewise polynomial function.

The *relu* function has a single point of non-differentiability at zero. If we "avoid" this point, we will have a total model on an example. More formally, let's take apart deep networks in a different way. Specifically, for all $v \in V - I$, $w \in \mathbf{R}^E$, define $k_{v,w} : X \rightarrow \mathbf{R}$ such that:

$$k_{v,w}(x) = \sum_{u:\{u,v\} \in E} c_{u,w}(x)w_{(u,v)}. \quad (126)$$

Notice that, for all $v \in V - I$:

$$c_{v,w}(x) = a_v(k_{v,w}(x)). \quad (127)$$

In this section we assume a_v is the identity for $v \in I \cup \{o^*\}$ and a_v is a *relu* function elsewhere. Observe¹² that, for some $w \in \mathbf{R}^E$, for some $x \in X$, if $k_{v,w}(x) \neq 0$ for all $v \in V - I$, then E is total for x given M_w . However, notice that if $k_{v,w}(x) = 0$ for some $v \in V - I$, that does not mean that the features are not total. Specifically, if $k_{v,w}(x) = 0$, but for all $u \in V$ where $(v, u) \in E$, $w_{(v,u)} = 0$, then effectively the node v has no effect. Thus, we will define a node to be **soft** given M_w and x

¹²The following may not be obvious: however, it is a corollary of Lemma 31.

if $k_{v,w}(x) \neq 0$ or $\frac{\partial^+ M_w(x)}{\partial c_{v,w}(x)} = 0$. Otherwise, we will say v is **hard**. Notice that any node where all weights on outgoing edges are zero is soft.

Lemma 31 *Given a set $E' \subseteq E$, given a model M_w and an example $x \in X$, if for all $(u, v) \in E'$, for all $v' \geq v$ (in the sense that there exists a path from v to v'), v' is soft, then E' is total on M_w .*

Proof: At a high level, we construct differentiable functions, starting with one that maps the input of o^* to its output, and then incorporating more and more nodes and edges until we have incorporated all of E' . These functions will not be differentiable everywhere, simply where we need them to be.⁴

Consider V' to be the set of all vertices $v' \in V$ where there is a $(u, v) \in E'$ where $v \leq v'$. Now, without loss of generality, assume $E' = \{(u, v) \in E : v \in V'\}$. Then, we know that all $v' \in V'$ are soft. Without loss of generality, assume $o^* \in V'$.

Now, we can arrange a total ordering on the vertices in V' (in the opposite direction of the partial ordering induced by the graph), beginning with o^* , and we will denote the ordering $v_1 \dots v_j$, such that there is no path from v_i to v_j for all $i > j$, and $v_1 = o^*$. We will denote $V_i = \{v_1, \dots, v_i\}$ (such that $V_j = V'$). and $E_i = \{(u, v) \in E : u \in V_i\}$. Define $w^i \in \mathbf{R}^{E_i}$ such that $w_{(u,v)}^i = w_{(u,v)}$ if $(u, v) \in E_i$. Define $\alpha_{v_1}^1 = k_{v_1,w}(x)$. We define $d^1 : \mathbf{R}^{V_1} \times \mathbf{R}^{E_1} \rightarrow \mathbf{R}$ such that¹³ $d^1(\alpha, \emptyset) = a_{v_1}(\alpha) = \alpha$. Thus, $d^1(\alpha_{v_1}^1) = c_{v_1,w}(x)$, and d^1 is differentiable with respect to its arguments.

Now, we recursively define $d^2 \dots d^j$. Given d^i , we define $m^i : \mathbf{R}^{V_{i+1}} \times \mathbf{R}^{E_{i+1}} \rightarrow \mathbf{R}_v^i$ such that $(m^i(\alpha, w'))_v = \alpha_v + a(v_{i+1})w'_{v_{i+1},v}$ if $(v_{i+1}, v) \in E$ and $m^i(\alpha, w')_v = \alpha_v$ if $(v_{i+1}, v) \notin E$. Moreover, for any sets S , and $T \subseteq S$, for any $v \in \mathbf{R}^S$, we define $\Pi^T(v)$ such that $\Pi^T(v)_i = v_i$ for all $i \in T$. Then, we can define $d^{i+1} : \mathbf{R}^{V_{i+1}} \times \mathbf{R}^{E_{i+1}} \rightarrow \mathbf{R}$ such that $d^{i+1}(\alpha, w') = d^i(m^i(\alpha, w'), \Pi^{E_i}(w'))$.

Recursively, we define $\alpha^{i+1} \in \mathbf{R}^{V_{i+1}}$ such that $\alpha_{v_{i+1}} = k_{v_{i+1},w}(x)$ and for all $v \in V_i$, $\alpha_v^{i+1} = \alpha_v^i - a(k_{v_{i+1},w}(x))w_{(v_{i+1},v)}$ if $(v_{i+1}, v) \in E$, and $\alpha_v^{i+1} = \alpha_v^i$ otherwise. Thus, observe that $d^{i+1}(\alpha^{i+1}, w^{i+1}) = d^i(\alpha^i, w^i)$, so recursively $d^{i+1}(\alpha^{i+1}, w^{i+1}) = c_{v_1,w}(x)$. Moreover, recursively one can establish that for any $w' \in \mathbf{R}^E$ where $\Pi^{E-E_{i+1}}(w') = \Pi^{E-E_{i+1}}(w)$, $d^{i+1}(\alpha^{i+1}, \Pi^{E_i}(w')) = c_{v_1,w'}(x)$.

Now, we must show differentiability. The inductive hypothesis is d^i is differentiable, or more formally, for

¹³Note that E_1 is the emptyset, as there are no outgoing edges from o_1^* .

$h \in \mathbf{R}^{V_i} \times \mathbf{R}^{E_i}$:

$$\lim_{h \rightarrow 0} \frac{d^i(\alpha^i + \Pi^{V_i}(h), w^i + \Pi^{E_i}(h))}{\|h\|} = 0 \quad (128)$$

First, we know that v_i is soft. If $k_{v_{i+1}, w}(x) \neq 0$, then $\alpha_{v_{i+1}}^{i+1} \neq 0$. If $\alpha_{v_{i+1}}^{i+1} < 0$, then in a region around α^{i+1}, w^{i+1} , $m^i((\alpha^{i+1}, w^{i+1}) + h) = \alpha^i + \Pi^{V_i}(h)$, i.e. is linear. If $\alpha_{v_{i+1}}^{i+1} > 0$, then in a region around α^{i+1}, w^{i+1} , $m^i((\alpha^{i+1}, w^{i+1}) + h)_v = \alpha_v^i + h_v + h_{v_{i+1}}(w_{(v_{i+1}, v)}^{i+1} + h_{(v_{i+1}, v)})$ if $(v_{i+1}, v) \in E$, and $m^i((\alpha^{i+1}, w^{i+1}) + h)_v = \alpha_v^i + h_v$ otherwise (again, linear).

Thus, if $k_{v_{i+1}, w}(x) \neq 0$, then $\alpha_{v_{i+1}}^{i+1} \neq 0$, in a region around α^{i+1} and w^i , m^i is linear. Then d^{i+1} is the composition of a differentiable function, a linear function, and a projection, and therefore differentiable at α^{i+1} and w^i .

On the other hand, if $k_{v_{i+1}, w}(x) = 0$, then $\alpha_{v_{i+1}}^{i+1} = 0$. Since v_i is soft, then $\frac{\partial^+ M_w(x)}{\partial c_{v_{i+1}, w}(x)} = 0$. This implies $\frac{\partial^+ d^{i+1}(\alpha^{i+1}, w^{i+1})}{\partial \alpha_{v_{i+1}}^{i+1}} = 0$. Using the chain rule for directional derivatives:

$$0 = \frac{\partial^+ d^{i+1}(\alpha^{i+1}, w^{i+1})}{\partial \alpha_{v_{i+1}}^{i+1}} \quad (129)$$

$$0 = \partial_g d^i(\alpha^i, w^i) \quad (130)$$

Where $g = \frac{\partial^+ m^i(\alpha^{i+1}, w^{i+1})}{\partial \alpha_{v_{i+1}}^{i+1}}$, so $g_v = w_{v_{i+1}, v}^{i+1}$ if $(v_{i+1}, v) \in E$, and $g_v = 0$ otherwise. If J^i is the derivative of d^i at α^i, w^i , then, because v_{i+1} is soft and $\alpha_{v_{i+1}}^{i+1} = 0$:

$$\sum_{v \in V_i} J_v^i g_v = 0 \quad (131)$$

$$\sum_{(v_{i+1}, v) \in E} J_v^i w_{v_{i+1}, v} = 0 \quad (132)$$

In order to prove differentiability of d^{i+1} , we introduce a function $\epsilon^i : \mathbf{R}^{V_i} \times \mathbf{R}^{E_i} \rightarrow \mathbf{R}$ quantifying the error of the derivative of d^i , such that for any $h^\alpha \in \mathbf{R}^{V_i}$, $h^w \in \mathbf{R}^{E_i}$, $\epsilon(h^\alpha, h^w) = d^i(\alpha^i + h^\alpha, w^i + h^w) - (\sum_{v \in V_i} J_v^i h_v^\alpha + \sum_{e \in E_i} J_e^i h_e^w)$. Thus, $\lim_{h^\alpha, h^w \rightarrow 0} \frac{\epsilon(h^\alpha, h^w)}{\|h^\alpha, h^w\|} = 0$, where $\|h^\alpha, h^w\| = \sqrt{\|h^\alpha\|^2 + \|h^w\|^2}$.

First, although m^i is not differentiable when $\alpha_{v_{i+1}}^{i+1} = 0$, the derivative is ‘‘almost’’ the linear projection operator Π^{V_i} , because the partial derivative of $\alpha_{v_{i+1}}^{i+1} = 0$. We can prove this using a proof similar to the proof of the chain rule. We will denote the following η , and will endeavor

to prove it exists and is zero.

$$\eta = \lim_{h^\alpha, h^w \rightarrow 0} \frac{d^{i+1}(\alpha^{i+1} + h^\alpha, w^{i+1} + h^w)}{\|h^\alpha, h^w\|} \quad (133)$$

$$= \frac{d^{i+1}(\alpha^{i+1}, w^{i+1}) + \sum_{v \in V_i} J_v^i h_v^\alpha + \sum_{e \in E_i} J_e^i h_e^w}{\|h^\alpha, h^w\|} \quad (134)$$

The last sum is a linear operator: notice that in this case, we are effectively showing $J^{i+1} = J^i$. First, we note that $d^{i+1}(\alpha^{i+1}, w^{i+1}) = d^i(\alpha^i, w^i)$:

$$\eta = \lim_{h^\alpha, h^w \rightarrow 0} \frac{d^{i+1}(\alpha^{i+1} + h^\alpha, w^{i+1} + h^w)}{\|h^\alpha, h^w\|} \quad (135)$$

$$= \frac{d^i(\alpha^i, w^i) + \sum_{v \in V_i} J_v^i h_v^\alpha + \sum_{e \in E_i} J_e^i h_e^w}{\|h^\alpha, h^w\|} \quad (136)$$

Furthermore, we study the first term separately:

$$\begin{aligned} & d^{i+1}(\alpha^{i+1} + h^\alpha, w^{i+1} + h^w) \\ &= d^i(m^i(\alpha^{i+1} + h^\alpha, w^{i+1} + h^w), \Pi^{E_i}(w^{i+1} + h^w)) \\ &= d^i(m^i(\alpha^{i+1} + h^\alpha, w^{i+1} + h^w), w^i + \Pi^{E_i}(h^w)) \end{aligned} \quad (137)$$

Write $a = a_{v_{i+1}}$, which is a relu function, because $v_{i+1} \neq o^*$, and $v_{i+1} \notin I$. Using ϵ we get:

$$\begin{aligned} & d^{i+1}(\alpha^{i+1} + h^\alpha, w^{i+1} + h^w) \\ &= d^i(\alpha^i, w^i) + \epsilon^i(m^i(\alpha^{i+1} + h^\alpha, w^{i+1} + h^w), \Pi^{E_i}(w^{i+1} + h^w)) \\ & \quad + \sum_{v \in V_i} J_v^i h_v^\alpha \\ & \quad + \sum_{(v_{i+1}, v) \in E} J_v^i a(h_{v_{i+1}}^\alpha)(w_{v_{i+1}, v} + h_{v_{i+1}, v}^w) \\ & \quad + \sum_{e \in E_i} J_e^i h_e^w \end{aligned} \quad (138)$$

Focusing on the most complex term:

$$\begin{aligned} & \sum_{(v_{i+1}, v) \in E} J_v^i a(h_{v_{i+1}}^\alpha)(w_{v_{i+1}, v} + h_{v_{i+1}, v}^w) \\ &= a(h_{v_{i+1}}^\alpha) \sum_{(v_{i+1}, v) \in E} J_v^i w_{v_{i+1}, v} \\ & \quad + a(h_{v_{i+1}}^\alpha) \sum_{(v_{i+1}, v) \in E} J_v^i h_{v_{i+1}, v}^w \end{aligned} \quad (139)$$

Note that $\sum_{v_{i+1}, v} J_v^i w_{v_{i+1}, v} = 0$, so:

$$\begin{aligned} & \sum_{(v_{i+1}, v) \in E} J_v^i a(h_{v_{i+1}}^\alpha)(w_{v_{i+1}, v} + h_{v_{i+1}, v}^w) \\ &= a(h_{v_{i+1}}^\alpha) \sum_{(v_{i+1}, v) \in E} J_v^i h_{v_{i+1}, v}^w \end{aligned} \quad (140)$$

So, now we consider the limit:

$$\begin{aligned} & \lim_{h^\alpha, h^w \rightarrow 0} \frac{\sum_{(v_{i+1}, v) \in E} J_v^i a(h_{v_{i+1}}^\alpha)(w_{v_{i+1}, v} + h_{v_{i+1}, v}^w)}{\|h^\alpha, h^w\|} \\ &= \lim_{h^\alpha, h^w \rightarrow 0} \frac{a(h_{v_{i+1}}^\alpha) \sum_{(v_{i+1}, v) \in E} J_v^i h_{v_{i+1}, v}^w}{\|h^\alpha, h^w\|} \quad (141) \end{aligned}$$

Taking the absolute value:

$$\begin{aligned} & \lim_{h^\alpha, h^w \rightarrow 0} \frac{\left| \sum_{(v_{i+1}, v) \in E} J_v^i a(h_{v_{i+1}}^\alpha)(w_{v_{i+1}, v} + h_{v_{i+1}, v}^w) \right|}{\|h^\alpha, h^w\|} \\ &= \lim_{h^\alpha, h^w \rightarrow 0} \frac{|a(h_{v_{i+1}}^\alpha)| \sum_{(v_{i+1}, v) \in E} |J_v^i| |h_{v_{i+1}, v}^w|}{\|h^\alpha, h^w\|} \quad (142) \end{aligned}$$

Since $|a(h_{v_{i+1}}^\alpha)| \leq |h_{v_{i+1}}^\alpha|$:

$$\begin{aligned} & \lim_{h^\alpha, h^w \rightarrow 0} \frac{\left| \sum_{(v_{i+1}, v) \in E} J_v^i a(h_{v_{i+1}}^\alpha)(w_{v_{i+1}, v} + h_{v_{i+1}, v}^w) \right|}{\|h^\alpha, h^w\|} \\ &= \lim_{h^\alpha, h^w \rightarrow 0} \frac{|h_{v_{i+1}}^\alpha| \sum_{(v_{i+1}, v) \in E} |J_v^i| |h_{v_{i+1}, v}^w|}{\|h^\alpha, h^w\|} \quad (143) \end{aligned}$$

The quadratic terms in the numerator mean that the limit is zero, because for any $i, j \in \{1 \dots n\}$ $\lim_{z \rightarrow 0} \frac{z_i z_j}{\|z\|} = 0$.

$$\lim_{h^\alpha, h^w \rightarrow 0} \frac{\sum_{(v_{i+1}, v) \in E} J_v^i a(h_{v_{i+1}}^\alpha)(w_{v_{i+1}, v} + h_{v_{i+1}, v}^w)}{\|h^\alpha, h^w\|} = 0$$

We next consider the term $\epsilon^i(m^i(\alpha^{i+1} + h^\alpha, w^{i+1} + h^w), \Pi^{E_i}(w^{i+1} + h^w))$. In order to bound this, we need to understand how fast m^i is approaching α^i .

$$\begin{aligned} & \left\| m^i(\alpha^{i+1} + h^\alpha, w^{i+1} + h^w) - \alpha^i \right\|^2 \\ &= \sum_{v \in V_i: (v_{i+1}, v) \in E} (h_v^\alpha + a(h_{v_{i+1}}^\alpha)(w_{v_{i+1}, v}^{i+1} + h_{v_{i+1}, v}^w))^2 \\ &\quad + \sum_{v \in V_i: (v_{i+1}, v) \notin E} (h_v^\alpha)^2 \\ &\leq \sum_{v \in V_i: (v_{i+1}, v) \in E} 2(h_v^\alpha)^2 \\ &\quad + \sum_{v \in V_i: (v_{i+1}, v) \in E} 4(a(h_{v_{i+1}}^\alpha))^2 (w_{v_{i+1}, v}^{i+1})^2 \\ &\quad + \sum_{v \in V_i: (v_{i+1}, v) \in E} 4(a(h_{v_{i+1}}^\alpha))^2 (h_{v_{i+1}, v}^w)^2 \\ &\quad + \sum_{v \in V_i: (v_{i+1}, v) \notin E} (h_v^\alpha)^2 \quad (145) \end{aligned}$$

Define $W^{i+1} = \sum_{v \in V_i: (v_{i+1}, v) \in E} (w_{v_{i+1}, v})^2$.

$$\begin{aligned} & \left\| m^i(\alpha^{i+1} + h^\alpha, w^{i+1} + h^w) - \alpha^i \right\|^2 \\ &\leq 2 \|h^\alpha\|^2 + 4(a(h_{v_{i+1}}^\alpha))^2 W^{i+1} + 4(a(h_{v_{i+1}}^\alpha))^2 \|h^w\|^2 \quad (146) \end{aligned}$$

In the limit $|h_{v_{i+1}}^\alpha| < 1$, so $(a(h_{v_{i+1}}^\alpha))^2 < 1$, and we can write:

$$\begin{aligned} & \left\| m^i(\alpha^{i+1} + h^\alpha, w^{i+1} + h^w) - \alpha^i \right\|^2 \\ &\leq 2 \|h^\alpha\|^2 + 4(a(h_{v_{i+1}}^\alpha))^2 W^{i+1} + 4 \|h^w\|^2 \quad (147) \end{aligned}$$

Moreover, $a(h_{v_{i+1}}^\alpha)^2 \leq (h_{v_{i+1}}^\alpha)^2 \leq \|h^\alpha, h^w\|^2$:

$$\begin{aligned} & \left\| m^i(\alpha^{i+1} + h^\alpha, w^{i+1} + h^w) - \alpha^i \right\|^2 \\ &\leq (4 + 4W^{i+1}) \|h^\alpha, h^w\|^2 \quad (148) \end{aligned}$$

Since $\Pi^{E_i}(w^{i+1} - h^w) - w^i = \Pi^{E_i}(h^w)$, then $\|\Pi^{E_i}(w^{i+1} - h^w) - w^i\|^2 \leq \|h^w\|^2 \leq \|h^\alpha, h^w\|^2$, so:

$$\begin{aligned} & \left\| m^i(\alpha^{i+1} + h^\alpha, w^{i+1} + h^w) - \alpha^i, \Pi^{E_i}(w^{i+1} - h^w) - w^i \right\| \\ &\leq \sqrt{5 + 4W^{i+1}} \|h^\alpha, h^w\| \quad (149) \end{aligned}$$

So, considering the limit of the absolute value:

$$\begin{aligned} & \lim_{h^\alpha, h^w \rightarrow 0} \frac{|\epsilon^i(m^i(\alpha^{i+1} + h^\alpha, w^{i+1} + h^w), \Pi^{E_i}(w^{i+1} + h^w))|}{\|h^\alpha, h^w\|} \\ &= \sqrt{5 + 4W^{i+1}} \lim_{h^\alpha, h^w \rightarrow 0} \frac{|\epsilon^i(m^i(\alpha^{i+1} + h^\alpha, w^{i+1} + h^w), \Pi^{E_i}(w^{i+1} + h^w))|}{\sqrt{5 + 4W^{i+1}} \|h^\alpha, h^w\|} \quad (144) \end{aligned}$$

Note that since the denominator is larger than the norm of the vector inside ϵ^i , then the whole thing approaches zero.

$$\lim_{h^\alpha, h^w \rightarrow 0} \frac{|\epsilon^i(m^i(\alpha^{i+1} + h^\alpha, w^{i+1} + h^w), \Pi^{E_i}(w^{i+1} + h^w))|}{\|h^\alpha, h^w\|} = 0 \quad (151)$$

Now we have established that:

$$\begin{aligned} & \lim_{h^\alpha, h^w \rightarrow 0} \frac{d^{i+1}(\alpha^{i+1} + h^\alpha, w^{i+1} + h^w)}{\|h^\alpha, h^w\|} \\ &= \lim_{h^\alpha, h^w \rightarrow 0} \frac{d^i(\alpha^i, w^i)}{\|h^\alpha, h^w\|} \\ &\quad + \frac{\sum_{v \in V_i} J_v^i h_v^\alpha}{\|h^\alpha, h^w\|} \\ &\quad + \frac{\sum_{e \in E_i} J_e^i h_e^w}{\|h^\alpha, h^w\|} \quad (152) \end{aligned}$$

Plugging this into η yields $\eta = 0$, implying that d^{i+1} is differentiable. Recursively, we have established that d^j is differentiable. However, E_j is a subset of E' : by construction, $V' = V_j$ are the nodes that are at the end of the edges in E' , not at the beginning, and E_j are the edges that can be reached from V_j . The final step is to construct a function $f : \mathbf{R}^{E'} \rightarrow \mathbf{R}$ as a function defined on all the relevant edges. Define $E^* = E' - E_j$. For all $(u, v) \in E^*$, define $\beta_{(u,v)} = c_{u,w}(x)$. Define $w^* \in \mathbf{R}^{E'}$ such that $w_e^* = w_e$ for all $e \in E'$. Then, we define $m^* : \mathbf{R}^{E'} \rightarrow \mathbf{R}^{V_j}$ such that for all $v \in V_j$, for all $w' \in \mathbf{R}^{E'}$:

$$m^*(w') = \sum_{(u,v) \in E^*} w'_{(u,v)} \beta_{(u,v)} \quad (153)$$

Now we can define $f(w') = d^j(m^*(w'), \Pi^{E_j}(w'))$. Note that $f(w^*) = M_w(x)$: moreover, for any $w' \in \mathbf{R}^{E'}$, if $\Pi^{E-E'}(w') = \Pi^{E-E'}(w)$, then $f(\Pi^{E'}(w')) = M_{w'}(x)$. Finally, since d^j is differentiable and m^* is linear, then f is differentiable at w^* . This implies that E' is total on x given $M_w(x)$. ■